

000 001 002 003 004 005 BOOSTING PERTURBED GRADIENT ASCENT 006 FOR LAST-ITERATE CONVERGENCE IN GAMES 007 008 009

010 **Anonymous authors**
 011 Paper under double-blind review
 012
 013
 014
 015
 016
 017
 018
 019
 020
 021
 022
 023

ABSTRACT

024 This paper presents a payoff perturbation technique, introducing a strong convexity to players' payoff functions in games. This technique is specifically designed
 025 for first-order methods to achieve last-iterate convergence in games where the gradient of the payoff functions is monotone in the strategy profile space, potentially
 026 containing additive noise. Although perturbation is known to facilitate the convergence of learning algorithms, the magnitude of perturbation requires careful
 027 adjustment to ensure last-iterate convergence. Previous studies have proposed a scheme in which the magnitude is determined by the distance from a periodically re-initialized anchoring or reference strategy. Building upon this, we propose Gradient Ascent with Boosting Payoff Perturbation, which incorporates a novel perturbation into the underlying payoff function, maintaining the periodically re-initializing anchoring strategy scheme. This innovation empowers us to provide faster last-iterate convergence rates against the existing payoff perturbed
 028 algorithms, even in the presence of additive noise.
 029
 030
 031
 032

1 INTRODUCTION

033 This study considers online learning in monotone games, where the gradient of the payoff function is
 034 monotone in the strategy profile space. Monotone games encompassed diverse well-studied games
 035 as special instances, such as concave-convex games, zero-sum polymatrix games (Cai & Daskalakis,
 036 2011; Cai et al., 2016), λ -cocoercive games (Lin et al., 2020), and Cournot competition (Monderer
 037 & Shapley, 1996). Due to their wide-ranging applications, there has been growing interest in developing
 038 learning algorithms to compute Nash equilibria in monotone games.
 039
 040

041 Typical learning algorithms such as Gradient Ascent (Zinkevich, 2003) and Multiplicative Weights
 042 Update (Bailey & Piliouras, 2018) have been extensively studied and shown to converge to equilibria
 043 in an average-iterate sense, which is termed *average-iterate convergence*. However, averaging the
 044 strategies can be undesirable because it can lead to additional memory or computational costs in
 045 the context of training Generative Adversarial Networks (Goodfellow et al., 2014) and preference-
 046 based fine-tuning of large language models (Munos et al., 2024; Swamy et al., 2024). In contrast,
 047 *last-iterate convergence*, in which the updated strategy profile itself converges to a Nash equilibrium,
 048 has emerged as a stronger notion than average-iterate convergence.

049 Payoff-perturbed algorithms have recently been regaining attention in this context (Sokota et al.,
 050 2023; Liu et al., 2023). Payoff perturbation, a classical technique referenced in Facchinei & Pang
 051 (2003), introduces a strongly convex penalty to the players' payoff functions to stabilize learning.
 052 This leads to convergence toward approximate equilibria, not only in the *full feedback* setting where
 053 the perfect gradient vector of the payoff function can be used to update strategies, but also in the
 054 *noisy feedback* setting where the gradient vector is contaminated by noise.

055 However, to ensure convergence toward a Nash equilibrium of the underlying game, the magnitude
 056 of perturbation requires careful adjustment. As a remedy, it is adjusted by the distance from an
 057 anchoring or reference strategy. Koshal et al. (2010) and Tatarenko & Kamgarpour (2019) simply
 058 decay the magnitude in each iteration, and their methods asymptotically converge, since the per-
 059 turbated function gradually loses strong convexity. In response to this, recent studies (Perolat et al.,
 060 2021; Abe et al., 2023; 2024) re-initialize the anchoring strategies periodically, or in a predefined
 061 interval, so that they keep the perturbed function strongly convex and achieve non-asymptotic con-
 062 vergence.
 063

We should also mention the *optimistic* family of learning algorithms, which incorporates recency bias and exhibits last-iterate convergence (Daskalakis et al., 2018; Daskalakis & Panageas, 2019; Mertikopoulos et al., 2019; Wei et al., 2021). Unfortunately, the property has mainly been proven in the full feedback setting. Although it might empirically work with noisy feedback, the convergence is slower, as demonstrated in Section 6. The fast convergence in the noisy feedback setting is another reason why payoff-perturbed algorithms have been gaining renewed interest.

The most recent payoff-perturbed algorithm, *Adaptively Perturbed Mirror Descent* (APMD) (Abe et al., 2024), achieves $\tilde{\mathcal{O}}(1/\sqrt{T})^1$ and $\tilde{\mathcal{O}}(1/T^{1/10})$ last-iterate convergence rates in the full/noisy feedback setting, respectively. The motivation of this study lies in improving these convergence rates. We propose an elegant one-line modification of APMD, which effectively accelerates convergence. In fact, we just add the difference between the current anchoring strategy and the initial anchoring strategy to the payoff perturbation function in APMD.

Our contributions are manifold. Firstly, we propose a novel payoff-perturbed learning algorithm named *Gradient Ascent with Boosting Payoff Perturbation* (GABP). This method incorporates a unique perturbation payoff function, enabling it to achieve fast convergence. Subsequently, we prove that GABP exhibits accelerated $\tilde{\mathcal{O}}(1/T)$ and $\tilde{\mathcal{O}}(1/T^{1/7})$ last-iterate convergence rates to a Nash equilibrium with full and noisy feedback, respectively². To derive these rates, we utilize the concept of the potential function used in Cai & Zheng (2023). Specifically, the potential function we employ is customized for handling noisy feedback. We further show that each player’s individual regret is at most $\mathcal{O}((\ln T)^2)$ in the full feedback setting, provided all players play according to GABP. Finally, through our experiments, we demonstrate the competitive or superior performance of GABP over the Optimistic Gradient algorithm (Daskalakis et al., 2018; Wei et al., 2021), the Accelerated Optimistic Gradient algorithm (Cai & Zheng, 2023), and APMD in concave-convex games, irrespective of the presence of noise.

2 PRELIMINARIES

Monotone games. In this study, we focus on a continuous multi-player game, which is denoted as $([N], (\mathcal{X}_i)_{i \in [N]}, (v_i)_{i \in [N]})$. $[N] = \{1, 2, \dots, N\}$ denotes the set of N players. Each player $i \in [N]$ chooses a *strategy* π_i from a d_i -dimensional compact convex strategy space \mathcal{X}_i , and we write $\mathcal{X} = \prod_{i \in [N]} \mathcal{X}_i$. Each player i aims to maximize her payoff function $v_i : \mathcal{X} \rightarrow \mathbb{R}$, which is differentiable on \mathcal{X} . We denote $\pi_{-i} \in \prod_{j \neq i} \mathcal{X}_j$ as the strategies of all players except player i , and $\pi = (\pi_i)_{i \in [N]} \in \mathcal{X}$ as the *strategy profile*. This paper particularly studies learning in *smooth monotone games*, where the gradient operator $V(\cdot) = (\nabla_{\pi_i} v_i(\cdot))_{i \in [N]}$ of the payoff functions is monotone: $\forall \pi, \pi' \in \mathcal{X}$,

$$\langle V(\pi) - V(\pi'), \pi - \pi' \rangle \leq 0, \quad (1)$$

and L -Lipschitz for $L > 0$,

$$\|V(\pi) - V(\pi')\| \leq L \|\pi - \pi'\|, \quad (2)$$

where $\|\cdot\|$ denotes the ℓ_2 -norm.

Many common and well-studied games, such as concave-convex games, zero-sum polymatrix games (Cai et al., 2016), λ -cocoercive games (Lin et al., 2020), and Cournot competition (Monderer & Shapley, 1996), are included in the class of monotone games.

Example 2.1 (Concave-convex games). Consider a game defined by $(\{1, 2\}, (\mathcal{X}_1, \mathcal{X}_2), (v, -v))$, where $v : \mathcal{X}_1 \times \mathcal{X}_2 \rightarrow \mathbb{R}$. In this game, player 1 wishes to maximize v , while player 2 aims to minimize v . If v is concave in $\pi_1 \in \mathcal{X}_1$ and convex in $\pi_2 \in \mathcal{X}_2$, the game is called a concave-convex game or minimax optimization problem, and it is not hard to see that this game is a special case of monotone games.

Example 2.2 (Cournot Competition). Consider a Cournot competition model with a linear price function. There are N firms in competition, and each independently and simultaneously chooses a quantity $\pi_i \in \mathcal{X}_i := [0, C_i]$ to produce certain goods, where C_i is a constant greater than 0. The

¹We use $\tilde{\mathcal{O}}$ to denote a Landau notation that disregards a polylogarithmic factor.

²For a more detailed comparison of our rates with other works, please refer to Table 2 in Appendix E.2.

108 price of the goods is determined by a linear function $P(\pi) = a - b \sum_{i \in [N]} \pi_i$, where a and b are
 109 constants greater than 0. The payoff for each firm i is calculated as the total revenue from producing
 110 π_i units of the goods, minus the associated production cost, i.e., $v_i(\pi) = \pi_i P(\pi) - c_i \pi_i$. This game
 111 has been shown to satisfy the property of monotone games as defined in Eq. (1) (Bravo et al., 2018).
 112

113 **Nash equilibrium and gap function.** A *Nash equilibrium* (Nash, 1951) is a widely used solution
 114 concept for a game, which is a strategy profile where no player can gain by changing her own
 115 strategy. Formally, a strategy profile $\pi^* \in \mathcal{X}$ is called a Nash equilibrium, if and only if π^* satisfies
 116 the following condition:

$$117 \quad \forall i \in [N], \forall \pi_i \in \mathcal{X}_i, v_i(\pi_i^*, \pi_{-i}^*) \geq v_i(\pi_i, \pi_{-i}^*).$$

119 We define the set of all Nash equilibria to be Π^* . It has been shown that there exists at least one
 120 Nash equilibrium (Debreu, 1952) for any smooth monotone games.

121 To quantify the proximity to Nash equilibrium for a given strategy profile $\pi \in \mathcal{X}$, we use the *gap*
 122 function, which is defined as:

$$124 \quad \text{GAP}(\pi) := \max_{\tilde{\pi} \in \mathcal{X}} \langle V(\pi), \tilde{\pi} - \pi \rangle.$$

126 Additionally, we use another measure of proximity to Nash equilibrium, referred to as the *tangent*
 127 *residual*. This measure is defined as:

$$128 \quad r^{\tan}(\pi) := \min_{a \in N_{\mathcal{X}}(\pi)} \| -V(\pi) + a \|,$$

130 where $N_{\mathcal{X}}(\pi) = \{(a_i)_{i \in [N]} \in \prod_{i=1}^N \mathbb{R}^{d_i} \mid \sum_{i=1}^N \langle a_i, \pi'_i - \pi_i \rangle \leq 0, \forall \pi' \in \mathcal{X}\}$ is the normal cone
 131 at $\pi \in \mathcal{X}$. It is easy to see that $\text{GAP}(\pi) \geq 0$ (resp. $r^{\tan}(\pi) \geq 0$) for any $\pi \in \mathcal{X}$, and the equality
 132 holds if and only if π is a Nash equilibrium. Defining $D := \sup_{\pi, \pi' \in \mathcal{X}} \|\pi - \pi'\|$ as the diameter of
 133 \mathcal{X} , the gap function for any given strategy profile $\pi \in \mathcal{X}$ is upper bounded by its tangent residual:

134 **Lemma 2.3** (Lemma 2 of Cai et al. (2022a)). *For any $\pi \in \mathcal{X}$, we have:*

$$136 \quad \text{GAP}(\pi) \leq D \cdot r^{\tan}(\pi).$$

138 The gap function and the tangent residual are standard measures of proximity to Nash equilibrium;
 139 e.g., it has been used in Cai & Zheng (2023); Abe et al. (2024).

141 **Problem setting.** This study focuses on the online learning setting in which the following process
 142 repeats from iterations $t = 1$ to T : (i) Each player $i \in [N]$ chooses her strategy $\pi_i^t \in \mathcal{X}_i$, based on
 143 previously observed feedback; (ii) Each player i receives the (noisy) gradient vector $\hat{\nabla}_{\pi_i} v_i(\pi^t)$ as
 144 feedback. This study examines two feedback models: *full feedback* and *noisy feedback*. In the full
 145 feedback setting, each player observes the perfect gradient vector $\hat{\nabla}_{\pi_i} v_i(\pi^t) = \nabla_{\pi_i} v_i(\pi^t)$. In the
 146 noisy feedback setting, each player's gradient feedback $\nabla_{\pi_i} v_i(\pi^t)$ is contaminated by an additive
 147 noise vector ξ_i^t , i.e., $\hat{\nabla}_{\pi_i} v_i(\pi^t) = \nabla_{\pi_i} v_i(\pi^t) + \xi_i^t$, where $\xi_i^t \in \mathbb{R}^{d_i}$. Throughout the paper, we
 148 assume that ξ_i^t is the zero-mean and bounded-variance noise vector at each iteration t .

149 **Payoff-perturbed learning algorithms.** To facilitate the convergence in the online learning setting,
 150 recent studies have utilized a *payoff perturbation* technique, where payoff functions are per-
 151 turbed by strongly convex functions (Sokota et al., 2023; Liu et al., 2023; Abe et al., 2022). However,
 152 while the addition of these strongly convex functions leads learning algorithms to converge to a sta-
 153 tionary point, this stationary point may be significantly distant from a Nash equilibrium. Therefore,
 154 the magnitude of perturbation requires careful adjustment. Perolat et al. (2021); Abe et al. (2023;
 155 2024) have introduced a scheme in which the magnitude is determined by the distance (or divergence
 156 function) from an anchoring strategy σ_i , which is periodically re-initialized. Specifically, Adaptively
 157 Perturbed Mirror Descent (APMD) (Abe et al., 2024) perturbs each player's payoff function by a
 158 strongly convex divergence function $G(\pi_i, \sigma_i) : \mathcal{X}_i \times \mathcal{X}_i \rightarrow [0, \infty)$, where the anchoring strategy
 159 σ_i is periodically replaced by the current strategy π_i^t every predefined iterations T_σ .

160 Let us denote the number of updates of σ_i up to iteration t as $k(t)$, and the anchoring strategy
 161 after $k(t)$ updates as $\sigma_i^{k(t)}$. Since σ_i is overwritten every T_σ iterations, we can write $k(t) = \lfloor (t -$

162

163

Algorithm 1 GABP for player i .

164

Require: Learning rates $\{\eta_t\}_{t \geq 0}$, perturbation strength μ , update interval T_σ , initial strategy π_i^1

165

1: $k \leftarrow 1, \tau \leftarrow 0$

166

2: $\sigma_i^1 \leftarrow \pi_i^1$

167

3: **for** $t = 1, 2, \dots, T$ **do**

168

4: Receive the gradient feedback $\widehat{\nabla}_{\pi_i} v_i(\pi^t)$

169

5: Update the strategy by

170

$$\pi_i^{t+1} = \arg \max_{x \in \mathcal{X}_i} \left\{ \eta_t \left\langle \widehat{\nabla}_{\pi_i} v_i(\pi^t) - \mu \frac{\sigma_i^k - \sigma_i^1}{k+1} - \mu (\pi_i^t - \sigma_i^k), x \right\rangle - \frac{1}{2} \|x - \pi_i^t\|^2 \right\}$$

171

6: $\tau \leftarrow \tau + 1$

172

7: **if** $\tau = T_\sigma$ **then**

173

8: $k \leftarrow k + 1, \tau \leftarrow 0$

174

9: $\sigma_i^k \leftarrow \pi_i^{t+1}$

175

10: **end if**

176

11: **end for**

177

180

181

182 $1)/T_\sigma] + 1$ and $\sigma_i^{k(t)} = \pi_i^{T_\sigma(k(t)-1)+1}$. Except for the payoff perturbation and the update of the
183 anchor strategy, APMD updates each player i 's strategy in the same way as standard Mirror Descent
184 algorithms:

185

186

$$\pi_i^{t+1} = \arg \max_{x \in \mathcal{X}_i} \left\{ \eta_t \left\langle \widehat{\nabla}_{\pi_i} v_i(\pi^t) - \mu \nabla_{\pi_i} G(\pi_i^t, \sigma_i^{k(t)}), x \right\rangle - D_\psi(x, \pi_i^t) \right\},$$

187

188 where η_t is the learning rate at iteration t , $\mu \in (0, \infty)$ is the *perturbation strength*, and $D_\psi(\pi_i, \pi'_i) =$
189 $\psi(\pi_i) - \psi(\pi'_i) - \langle \nabla \psi(\pi'_i), \pi_i - \pi'_i \rangle$ is the Bregman divergence associated with a strictly convex
190 function $\psi : \mathcal{X}_i \rightarrow \mathbb{R}$. When both G and D_ψ is set to the squared ℓ^2 -distance, this algorithm can be
191 equivalently written as:

191

192

$$\pi_i^{t+1} = \arg \max_{x \in \mathcal{X}_i} \left\{ \eta_t \left\langle \widehat{\nabla}_{\pi_i} v_i(\pi^t) - \mu (\pi_i^t - \sigma_i^{k(t)}), x \right\rangle - \frac{1}{2} \|x - \pi_i^t\|^2 \right\}.$$

193

We refer to this version of APMD as Adaptively Perturbed Gradient Ascent (APGA). Abe et al. (2024) have shown that APGA exhibits the convergence rates of $\tilde{\mathcal{O}}(1/\sqrt{T})$ and $\tilde{\mathcal{O}}(1/T^{1/10})$ with full and noisy feedback, respectively.

194

195

3 GRADIENT ASCENT WITH BOOSTING PAYOFF PERTURBATION

196

197

This section proposes a novel payoff-perturbed learning algorithm, Gradient Ascent with Boosting Payoff Perturbation (GABP). The pseudo-code of GABP is outlined in Algorithm 1. At each iteration $t \in [T]$, GABP receives the gradient feedback $\widehat{\nabla}_{\pi_i} v_i(\pi^t)$, and updates each player's strategy by the following update rule:

198

199

$$\pi_i^{t+1} = \arg \max_{x \in \mathcal{X}_i} \left\{ \eta_t \left\langle \widehat{\nabla}_{\pi_i} v_i(\pi^t) - \underbrace{\mu \frac{\sigma_i^{k(t)} - \sigma_i^1}{k(t)+1}}_{(*)} - \mu (\pi_i^t - \sigma_i^{k(t)}), x \right\rangle - \frac{1}{2} \|x - \pi_i^t\|^2 \right\}. \quad (3)$$

200

201

σ_i is overwritten every T_σ iterations, and thus $\sigma_i^{k(t)}$ is define as $\sigma_i^{k(t)} = \pi_i^{T_\sigma(k(t)-1)+1}$. The term
202 $(*)$ in Eq. (3) is our proposed additional perturbation term. It shrinks as $k(t)$, the number of updates
203 of $\sigma_i^{k(t)}$, increases.

204

205

For a more intuitive explanation of the proposed perturbation term, we present the following update rule, which is equivalent to Eq. (3):

206

207

$$\pi_i^{t+1} = \arg \max_{x \in \mathcal{X}_i} \left\{ \eta_t \left\langle \widehat{\nabla}_{\pi_i} v_i(\pi^t) - \mu \left(\pi_i^t - \frac{k(t)\sigma_i^{k(t)} + \sigma_i^1}{k(t)+1} \right), x \right\rangle - \frac{1}{2} \|x - \pi_i^t\|^2 \right\}.$$

According this formula, it is evident that GABP perturbs the gradient vector $\widehat{\nabla}_{\pi} v_i(\pi^t)$ so that π^t approaches $\frac{k(t)\sigma_i^{k(t)} + \sigma_i^1}{k(t)+1}$, instead of $\sigma_i^{k(t)}$. This gradual evolution of the anchoring strategy in GABP, compared to $\sigma_i^{k(t)}$ itself, is anticipated to contribute to the stabilization of the learning dynamics. There is a tradeoff between the shrinking speed of the term (*) and the stabilizing impact on the last-iterate convergence rate of GABP. The shrinking speed of $1/(k(t)+1)$ achieves a faster convergence rate, and we believe that this represents the optimal balance for this trade-off. We remark that the term (*) bears a resemblance to the update rule of the Accelerated Optimistic Gradient (AOG) algorithm (Cai & Zheng, 2023). However, AOG differs in the sense that it actually modifies the proximal point in gradient ascent, instead of perturbing the gradient vector. A detailed comparison is discussed in Appendix E.1.

4 LAST-ITERATE CONVERGENCE RATES

This section provides the last-iterate convergence rates of GABP. Specifically, we derive $\tilde{\mathcal{O}}(1/T)$ and $\tilde{\mathcal{O}}(1/T^{\frac{1}{7}})$ rates for the full/noisy feedback setting, respectively.

4.1 FULL FEEDBACK SETTING

First, we demonstrate the last-iterate convergence rate of GABP with *full feedback* where each player receives the perfect gradient vector as feedback at each iteration t , i.e., $\widehat{\nabla}_{\pi_i} v_i(\pi^t) = \nabla_{\pi_i} v_i(\pi^t)$. Theorem 4.1 shows that the last-iterate strategy profile π^T updated by GABP converges to a Nash equilibrium with an $\tilde{\mathcal{O}}(1/T)$ rate in the full feedback setting.

Theorem 4.1. *If we use the constant learning rate $\eta_t = \eta \in (0, \frac{\mu}{(L+\mu)^2})$ and the constant perturbation strength $\mu > 0$, and set $T_\sigma = c \cdot \max(1, \frac{6 \ln 3(T+1)}{\ln(1+\eta\mu)})$ for some constant $c \geq 1$, then the strategy π^t updated by GABP satisfies for any $t \in [T]$:*

$$\text{GAP}(\pi^{t+1}) \leq D \cdot r^{\tan}(\pi^{t+1}) \leq \frac{17cD^2 \left(\frac{6 \ln 3(T+1)}{\ln(1+\eta\mu)} + 1 \right)}{t} \left(\mu + \frac{1 + \eta L}{\eta} \right).$$

This rate is competitive compared to the previous state-of-the-art rate of $\mathcal{O}(1/T)$ (Yoon & Ryu, 2021; Cai & Zheng, 2023). Note that the rate in Theorem 4.1 holds for any fixed $\mu > 0$.

4.2 NOISY FEEDBACK SETTING

Next, we establish the last-iterate convergence rate in the *noisy feedback* setting, where each player i observes a noisy gradient vector contaminated by an additive noise vector $\xi_i^t \in \mathbb{R}^{d_i}$: $\widehat{\nabla}_{\pi_i} v_i(\pi^t) = \nabla_{\pi_i} v_i(\pi^t) + \xi_i^t$. We assume that the noisy vector ξ_i^t is zero-mean and its variance is bounded. Formally, defining the sigma-algebra generated by the history of the observations as $\mathcal{F}_t := \sigma((\widehat{\nabla}_{\pi_i} v_i(\pi^1))_{i \in [N]}, \dots, (\widehat{\nabla}_{\pi_i} v_i(\pi^{t-1}))_{i \in [N]})$, $\forall t \geq 1$, the noisy vector ξ_i^t is assumed to satisfy the following conditions:

Assumption 4.2. *ξ_i^t satisfies the following properties for all $t \geq 1$ and $i \in [N]$: (a) Zero-mean: $\mathbb{E}[\xi_i^t | \mathcal{F}_t] = (0, \dots, 0)^\top$; (b) Bounded variance: $\mathbb{E}[\|\xi_i^t\|^2 | \mathcal{F}_t] \leq C^2$ with some constant $C > 0$.*

Assumption 4.2 is standard in online learning in games with noisy feedback (Mertikopoulos & Zhou, 2019; Hsieh et al., 2019; Abe et al., 2024) and stochastic optimization (Nemirovski et al., 2009; Nedić & Lee, 2014). Under Assumption 4.2 and a decreasing learning rate sequence η_t , we can obtain a faster last convergence rate $\tilde{\mathcal{O}}(1/T^{\frac{1}{7}})$ than $\tilde{\mathcal{O}}(1/T^{\frac{1}{10}})$ of APGA (Abe et al., 2024).

Theorem 4.3. *Let $\kappa = \frac{\mu}{2}, \theta = \frac{3\mu^2 + 8L^2}{2\mu}$. Suppose that Assumption 4.2 holds and $V(\pi) \leq \zeta$ for any $\pi \in \mathcal{X}$. We also assume that T_σ is set to satisfy $T_\sigma = c \cdot \max(T^{\frac{6}{7}}, 1)$ for some constant $c \geq 1$. If we use the constant perturbation strength $\mu > 0$ and the decreasing learning rate sequence $\eta_t = \frac{1}{\kappa(t - T_\sigma(k(t)-1) + 2\theta)}$, then the strategy π^{T+1} satisfies:*

$$\mathbb{E} [\text{GAP}(\pi^{T+1})] = \mathcal{O} \left(\frac{\ln T}{T^{\frac{1}{7}}} \right).$$

270 4.3 PROOF SKETCH OF THEOREMS 4.1 AND 4.3
 271

272 This section outlines the sketch of the proofs for Theorems 4.1 and 4.3. The complete proofs are
 273 placed in Appendix A and B.

274 We define the stationary point $\pi^{\mu, \sigma^{k(t)}}$, which satisfies the following condition: $\forall i \in [N]$,
 275

$$276 \quad \pi_i^{\mu, \sigma^{k(t)}} = \arg \max_{x \in \mathcal{X}_i} \left\{ v_i(x, \pi_{-i}^{\mu, \sigma^{k(t)}}) - \frac{\mu}{2} \|x - \hat{\sigma}_i^{k(t)}\|^2 \right\},$$

279 where $\hat{\sigma}_i^{k(t)} = \frac{k(t)\sigma_i^{k(t)} + \sigma_i^1}{k(t)+1}$. The primary technical challenge in deriving the last-iterate convergence
 280 rates lies in the construction of the following potential function $P^{k(t)}$, which can be utilized in the
 281 proofs for both full and noisy feedback settings:
 282

$$283 \quad P^{k(t)} := k(t)(k(t) + 1) \left(\frac{\|\pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1}\|^2}{2} + \langle \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \rangle \right).$$

287 Specifically, we demonstrate that this potential function is approximately non-increasing regardless of the presence of noise. Although the potential function $P^{k(t)}$ is inspired by one in Cai
 288 & Zheng (2023), their potential function contains the term $\eta^2 \sum_{i \in [N]} \|\hat{\nabla} v_i(\pi^t) - \hat{\nabla} v_i(\pi^{t-\frac{1}{2}})\|^2$,
 289 which could have a high value in the noisy feedback setting even if $\pi^t = \pi^{t-\frac{1}{2}}$ holds³. This complicates providing a last-iterate convergence result for the noisy feedback setting via their potential
 290 function. In contrast, our potential function $P^{k(t)}$ does not include the term dependent on $\hat{\nabla} v_i(\pi^t)$.
 291 This enables us to provide the last-iterate convergence rates even for the noisy feedback setting.
 292

295 **(1) Potential function for bounding the distance between $\pi^{\mu, \sigma^{k(t)}}$ and $\sigma^{k(t)}$.** As mentioned
 296 above, our main technical contribution is proving that $P^{k(t)}$ is approximately non-increasing (as
 297 shown in Lemma A.3). That is, we have for any $t \geq 1$ such that $k(t) \geq 2$:

$$299 \quad P^{k(t)+1} \leq P^{k(t)} + (k(t) + 1)^2 \cdot \mathcal{O} \left(\|\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)+1}\| + \|\pi^{\mu, \sigma^{k(t)-1}} - \sigma^{k(t)}\| \right). \quad (4)$$

301 By telescoping of Eq. (4) and the first-order optimality condition for $\pi^{\mu, \sigma^{k(t)}}$, we can derive the
 302 following upper bound on the distance between $\pi^{\mu, \sigma^{k(t)}}$ and $\hat{\sigma}^{k(t)}$:

$$304 \quad \frac{(k(t) + 1)(k(t) + 2)}{2} \|\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}\|^2 \leq \mathcal{O}(1) + (k(t) + 1)^2 \sum_{l=1}^{k(t)} \mathcal{O} \left(\|\pi^{\mu, \sigma^l} - \sigma^{l+1}\| \right).$$

307 Applying the definition of $\hat{\sigma}^{k(t)}$ and Cauchy-Schwarz inequality to this inequality, we obtain:

$$309 \quad \|\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}\|^2 \leq \frac{\mathcal{O}(\|\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}\|)}{k(t) + 1} + \mathcal{O} \left(\frac{1}{(k(t) + 1)^2} \right) + \sum_{l=1}^{k(t)} \mathcal{O} \left(\|\pi^{\mu, \sigma^l} - \sigma^{l+1}\| \right). \quad (5)$$

312 Note that the non-increasing property of our potential function, as described in Eq. (4), holds even
 313 in the noisy feedback setting. This implies that a similar proof technique for deriving Eq. (5) can be
 314 utilized to provide last-iterate convergence results both in full and noisy feedback settings.
 315

316 **(2) Convergence rate of $\sigma^{k(t)+1}$ to the stationary point $\pi^{\mu, \sigma^{k(t)}}$.** Leveraging the strong con-
 317 vexitity of the perturbation payoff function, $\frac{\mu}{2} \|x - \hat{\sigma}_i^{k(t)}\|^2$, we show that π^t converges to $\pi^{\mu, \sigma^{k(t)}}$
 318 exponentially fast in the full feedback setting (as shown in Lemma A.1). Specifically, we have for
 319 any $t \geq 1$:

$$321 \quad \|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 \leq \left(\frac{1}{1 + \eta\mu} \right)^{t-(k(t)-1)T_\sigma} \|\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}\|^2. \quad (6)$$

322 ³The comparison of the potential function of Cai & Zheng (2023) with ours can be found in Appendix E.1.

324 By using Eq. (6) and the assumption that $T_\sigma \geq \frac{6 \ln 3(T+1)}{\ln(1+\eta\mu)}$ in the full feedback setting, we can
 325 easily show that $\|\pi^{\mu,\sigma^l} - \sigma^{l+1}\| = \|\pi^{\mu,\sigma^l} - \pi^{T_\sigma l+1}\| \leq \mathcal{O}(k(t)^{-3})$ for any $l \leq k(t)$. Hence, from
 326 Eq. (5), we can derive the following convergence rate of the distance between $\pi^{\mu,\sigma^{k(t)}}$ and $\sigma^{k(t)}$
 327 with respect to $k(t)$ (as shown in Lemma A.2):
 328

$$\|\pi^{\mu,\sigma^{k(t)}} - \sigma^{k(t)}\| \leq \mathcal{O}(1/k(t)). \quad (7)$$

329
 330
 331 **(3) Decomposition of the gap function of the last-iterate strategy profile π^{T+1} .** Let us define
 332 $K := \lfloor T/T_\sigma \rfloor$. From Cauchy-Schwarz inequality and Lemma 2.3, we can decompose the gap
 333 function $\text{GAP}(\pi^{T+1})$ as follows:
 334

$$\begin{aligned} \text{GAP}(\pi^{T+1}) &\leq \text{GAP}(\pi^{\mu,\sigma^K}) + \mathcal{O}\left(\|\pi^{\mu,\sigma^K} - \pi^{T+1}\|\right) \\ &\leq D \cdot \min_{a \in N_{\mathcal{X}}(\pi^{\mu,\sigma^K})} \| -V(\pi^{\mu,\sigma^K}) + a \| + \mathcal{O}\left(\|\pi^{\mu,\sigma^K} - \pi^{T+1}\|\right). \end{aligned}$$

335 From the first-order optimality condition for π^{μ,σ^K} , we can see that $V(\pi^{\mu,\sigma^K}) - \mu(\pi^{\mu,\sigma^K} - \hat{\sigma}^K) \in$
 336 $N_{\mathcal{X}}(\pi^{\mu,\sigma^K})$. Thus, from the triangle inequality and L -smoothness of the gradient operator in Eq. (2),
 337 the gap function $\text{GAP}(\pi^{T+1})$ can be bounded as:
 338

$$\text{GAP}(\pi^{T+1}) \leq \mathcal{O}(1/K) + \mathcal{O}\left(\|\pi^{\mu,\sigma^K} - \sigma^K\|\right) + \mathcal{O}\left(\|\pi^{\mu,\sigma^K} - \pi^{T+1}\|\right). \quad (8)$$

339
 340 **(4) Putting it all together: last-iterate convergence rate of π^{T+1} .** By combining Eq. (6), Eq. (7),
 341 and Eq. (8), it holds that $\text{GAP}(\pi^{T+1}) \leq \mathcal{O}(1/K)$ in the full feedback setting. Hence, given $K =$
 342 $\lfloor T/T_\sigma \rfloor$, we can deduce that $\text{GAP}(\pi^{T+1}) \leq \mathcal{O}(T_\sigma/T)$. Finally, taking $T_\sigma = \Theta(\ln T)$, we obtain
 343 the upper bound on the gap function for the full feedback setting: $\text{GAP}(\pi^{T+1}) \leq \mathcal{O}(\ln T/T)$. Note
 344 that using a similar proof technique, we can also derive an upper bound on the tangent residual for
 345 the full feedback setting.

346 In the context of the noisy feedback setting, we achieve the following convergence rate to $\pi^{\mu,\sigma^{k(t)}}$
 347 instead of Eq. (6) (as shown in Lemma B.1):
 348

$$\mathbb{E} \left[\left\| \pi^{\mu,\sigma^{k(t)}} - \pi^{t+1} \right\|^2 \right] \leq \mathcal{O} \left(\frac{\ln t}{t - (k(t) - 1)T_\sigma} \right). \quad (9)$$

349 By using Eq. (9) and the assumption that $T_\sigma = \Theta(T^{\frac{6}{7}})$, we can still derive Eq. (7) and Eq. (8) for
 350 the noisy feedback setting. Therefore, we conclude that: $\mathbb{E} [\text{GAP}(\pi^{T+1})] \leq \mathcal{O}(\ln T/T^{\frac{1}{7}})$.
 351

352 5 INDIVIDUAL REGRET BOUND

353 In this section, we present an upper bound on an individual regret for each player. Specifically, our
 354 study examines two performance measures: the *external regret* and the *dynamic regret* (Zinkevich,
 355 2003). The external regret is a conventional measure in online learning. In online learning in games,
 356 the external regret for player i is defined as the gap between the player's realized cumulative payoff
 357 and the cumulative payoff of the best fixed strategy in hindsight:

$$\text{Reg}_i(T) := \max_{x \in \mathcal{X}_i} \sum_{t=1}^T (v_i(x, \pi_{-i}^t) - v_i(\pi^t)).$$

358 The dynamics regret is a much stronger performance metric, which is given by:
 359

$$\text{DynamicReg}_i(T) := \sum_{t=1}^T \left(\max_{x \in \mathcal{X}_i} v_i(x, \pi_{-i}^t) - v_i(\pi^t) \right).$$

360 We show in Theorem 5.1 that the individual regret is at most $\mathcal{O}((\ln T)^2)$ if each player $i \in [N]$
 361 plays according to GABP in the full feedback setting. The proof is given in Appendix C.

362 **Theorem 5.1.** *In the same setup of Theorem 4.1, we have for any player $i \in [N]$ and $T \geq 2$:*

$$\text{Reg}_i(T) \leq \text{DynamicReg}_i(T) \leq \mathcal{O}((\ln T)^2).$$

363 This regret bound is significantly superior to the $\mathcal{O}(\sqrt{T})$ regret bound of the Optimistic Gradient
 364 (OG) algorithm, and it is slightly inferior to the $\mathcal{O}(\ln T)$ regret bound of AOG (Cai & Zheng, 2023).
 365

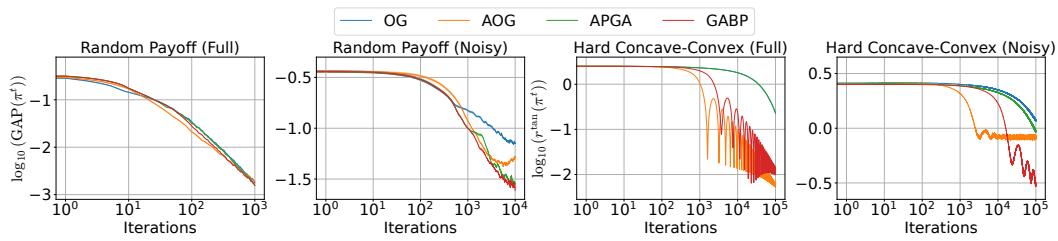


Figure 1: Performance of π^t for GABP, APGA, OG, and AOG with full and noisy feedback in the random payoff and hard concave-convex games, respectively. The shaded area represents the standard errors. Note that we report the gap function for the random payoff game, while the tangent residual is reported for the hard concave-convex game.

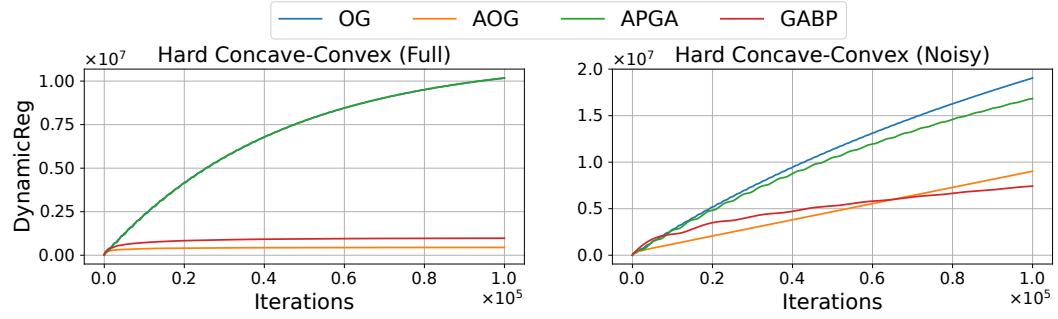


Figure 2: Dynamic regret for GABP, APGA, OG, and AOG with full and noisy feedback.

6 EXPERIMENTS

In this section, we present the empirical results of our GABP, comparing its performance with APGA (Abe et al., 2024), OG (Daskalakis et al., 2018; Wei et al., 2021), and AOG (Cai & Zheng, 2023). We conduct experiments on two classes of concave-convex games. One is random payoff games, which are two-player zero-sum normal-form games with payoff matrices of size d . Each player’s strategy space is represented by the d -dimensional probability simplex, i.e., $\mathcal{X}_1 = \mathcal{X}_2 = \Delta^d$. All entries of the payoff matrix are drawn independently from a uniform distribution over the interval $[-1, 1]$. We set $d = 50$ and the initial strategies are set to $\pi_1^1 = \pi_2^1 = \frac{1}{d}\mathbf{1}$. The other is a *hard concave-convex game* (Ouyang & Xu, 2021), formulated as the following max-min optimization problem: $\max_{x \in \mathcal{X}_1} \min_{y \in \mathcal{X}_2} f(x, y)$, where $f(x, y) = -\frac{1}{2}x^\top Hx + h^\top x + \langle Ax - b, y \rangle$. Following the setup in Cai & Zheng (2023), we choose $\mathcal{X}_1 = \mathcal{X}_2 = [-200, 200]^d$ with $d = 100$. The precise terms of $H \in \mathbb{R}^{d \times d}$, $A \in \mathbb{R}^{d \times d}$, $b \in \mathbb{R}^d$, and $h \in \mathbb{R}^d$ are provided in Appendix D.2. All algorithms are executed with initial strategies $\pi_1^1 = \pi_2^1 = \frac{1}{d}\mathbf{1}$. The detailed hyperparameters of the algorithms, tuned for best performance, are shown in Table 1 in Appendix D.3.

Figure 1 illustrates the logarithmic GAP or r^{tan} values per iteration for two games with each feedback. For the random payoff games with full or noisy feedback, 50 payoff matrices are generated using different random seeds. Likewise, for the hard concave-convex games, we use 10 different random seeds. We assume that the noise vector ξ_i^t is generated from the multivariate Gaussian distribution $\mathcal{N}(0, 0.1^2 \mathbf{I})$ in an i.i.d. manner for both games. In the former game with full feedback, GABP performs almost as well as the others. With noisy feedback, GABP outperforms the others, although the margin from APGA is slight. In the latter game, under the full feedback setting, GABP is competitive against AOG, whereas, under the noisy feedback setting, it demonstrates a substantial advantage over the others.

Figure 2 illustrates the dynamic regret in the hard concave-convex game. GABP exhibits lower regret than APGA and OG with both feedback, demonstrating its efficiency and robustness. Note that APGA and OG exhibit almost identical trajectories with full feedback, with their plots overlapping completely. In addition, GABP achieves competitive regret in comparison to AOG.

432 **7 RELATED LITERATURE**

433

434

435 No-regret learning algorithms have been extensively studied with the intent of achieving key objec-
 436 tives such as average-iterate convergence or last-iterate convergence. Recently, learning algorithms
 437 introducing optimism (Rakhlin & Sridharan, 2013a;b), such as optimistic Follow the Regularized
 438 Leader (Shalev-Shwartz & Singer, 2006) and optimistic Mirror Descent (Zhou et al., 2017; Hsieh
 439 et al., 2021), have been introduced to admit last-iterate convergence in a broad spectrum of game
 440 settings. These optimistic algorithms with full feedback have been shown to achieve last-iterate con-
 441 vergence in various classes of games, including bilinear games (Daskalakis et al., 2018; Daskalakis
 442 & Panageas, 2019; Liang & Stokes, 2019; de Montbrun & Renault, 2022), cocoercive games (Lin
 443 et al., 2020), and saddle point problems (Daskalakis & Panageas, 2018; Mertikopoulos et al., 2019;
 444 Golowich et al., 2020b; Wei et al., 2021; Lei et al., 2021; Yoon & Ryu, 2021; Lee & Kim, 2021;
 445 Cevher et al., 2023). Recent studies have provided finite convergence rates for monotone games
 446 (Golowich et al., 2020a; Cai et al., 2022a;b; Gorbunov et al., 2022; Cai & Zheng, 2023).

447

448 Compared to the full feedback setting, there are significant challenges in learning with noisy feed-
 449 back. For example, a learning algorithm must estimate the gradient from feedback that is contam-
 450 inated by noise. Despite the challenge, a vast literature has successfully achieved last-iterate con-
 451 vergence with noisy feedback in specific classes of games, including potential games (Cohen et al.,
 452 2017), strongly monotone games (Giannou et al., 2021b;a), and two-player zero-sum games (Abe
 453 et al., 2023). These results have often leveraged unique structures of their payoff functions, such as
 454 strict (or strong) monotonicity (Bravo et al., 2018; Kannan & Shanbhag, 2019; Hsieh et al., 2019;
 455 Anagnostides & Panageas, 2022) and strict variational stability (Mertikopoulos & Zhou, 2019; Mer-
 456 tikopoulos et al., 2019; 2022; Azizian et al., 2021). Without these restrictions, convergence is mainly
 457 demonstrated in an asymptotic manner, with no quantification of the rate (Hsieh et al., 2020; 2022;
 458 Abe et al., 2023). Consequently, an exceedingly large number of iterations might be necessary to
 459 reach an equilibrium.

460 There have been several studies focusing on payoff-regularized learning, where each player’s pay-
 461 off or utility function is perturbed or regularized via strongly convex functions (Cen et al., 2021;
 462 2023; Pattathil et al., 2023). Previous studies have successfully achieved convergence to stationary
 463 points, which are approximate equilibria. For instance, Sokota et al. (2023) have demonstrated that
 464 their perturbed mirror descent algorithm converges to a quantal response equilibrium (McKelvey &
 465 Palfrey, 1995; 1998). Similar results have been obtained with the Boltzmann Q-learning dynamics
 466 (Tuyls et al., 2006) and penalty-regularized dynamics (Coucheney et al., 2015) in continuous-
 467 time settings (Leslie & Collins, 2005; Abe et al., 2022; Hussain et al., 2023). To ensure convergence
 468 toward a Nash equilibrium of the underlying game, the magnitude of perturbation requires careful
 469 adjustment. Several learning algorithms have been proposed to gradually reduce the perturbation
 470 strength μ in response to this (Bernasconi et al., 2022; Liu et al., 2023; Cai et al., 2023). These
 471 include well-studied methods such as iterative Tikhonov regularization (Facchinei & Pang, 2003;
 472 Koshal et al., 2010; 2013; Yousefian et al., 2017; Tatarenko & Kamgarpour, 2019). Alternatively,
 473 Perolat et al. (2021) and Abe et al. (2023) have employed a payoff perturbation scheme, where the
 474 magnitude of perturbation is determined by the distance from an anchoring strategy, which is period-
 475 ically re-initialized by the current strategy. Recently, Abe et al. (2024) have established $\tilde{\mathcal{O}}(1/\sqrt{T})$
 476 and $\tilde{\mathcal{O}}(1/T^{\frac{1}{10}})$ last-iterate convergence rates for the payoff perturbation scheme in the full/noisy
 477 feedback setting, respectively. Our algorithm achieves faster $\tilde{\mathcal{O}}(1/T)$ and $\tilde{\mathcal{O}}(1/T^{\frac{1}{7}})$ last-iterate
 478 convergence rates by modifying the periodically re-initializing anchoring strategy scheme so that
 479 the anchoring strategy evolves more gradually.

480 **8 CONCLUSION**

481

482

483 This study proposes a novel payoff-perturbed algorithm, Gradient Ascent with Boosting Payoff Per-
 484 turbation, which achieves $\tilde{\mathcal{O}}(1/T)$ and $\tilde{\mathcal{O}}(1/T^{\frac{1}{7}})$ last-iterate convergence rates in monotone games
 485 with full/noisy feedback, respectively. Extending our results in settings where each player only
 486 observes bandit feedback is an intriguing and challenging future direction.

486 REFERENCES
487

- 488 Kenshi Abe, Mitsuki Sakamoto, and Atsushi Iwasaki. Mutation-driven follow the regularized leader
489 for last-iterate convergence in zero-sum games. In *UAI*, pp. 1–10, 2022.
- 490 Kenshi Abe, Kaito Ariu, Mitsuki Sakamoto, Kentaro Toyoshima, and Atsushi Iwasaki. Last-iterate
491 convergence with full and noisy feedback in two-player zero-sum games. In *AISTATS*, pp. 7999–
492 8028, 2023.
- 493 Kenshi Abe, Kaito Ariu, Mitsuki Sakamoto, and Atsushi Iwasaki. Adaptively perturbed mirror
494 descent for learning in games. In *ICML*, pp. 31–80, 2024.
- 495 Ioannis Anagnostides and Ioannis Panageas. Frequency-domain representation of first-order meth-
496 ods: A simple and robust framework of analysis. In *SOSA*, pp. 131–160, 2022.
- 497 Waïss Azizian, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. The last-iterate con-
498 vergence rate of optimistic mirror descent in stochastic variational inequalities. In *COLT*, pp.
499 326–358, 2021.
- 500 James P Bailey and Georgios Piliouras. Multiplicative weights update in zero-sum games. In *Eco-
501 nomics and Computation*, pp. 321–338, 2018.
- 502 Martino Bernasconi, Alberto Marchesi, and Francesco Trovò. Last-iterate convergence to trembling-
503 hand perfect equilibria. *arXiv preprint arXiv:2208.08238*, 2022.
- 504 Mario Bravo, David Leslie, and Panayotis Mertikopoulos. Bandit learning in concave N-person
505 games. In *NeurIPS*, pp. 5666–5676, 2018.
- 506 Yang Cai and Constantinos Daskalakis. On minmax theorems for multiplayer games. In *SODA*, pp.
507 217–234, 2011.
- 508 Yang Cai and Weiqiang Zheng. Doubly optimal no-regret learning in monotone games. In *ICML*,
509 pp. 3507–3524, 2023.
- 510 Yang Cai, Ozan Candogan, Constantinos Daskalakis, and Christos Papadimitriou. Zero-sum poly-
511 matrix games: A generalization of minmax. *Mathematics of Operations Research*, 41(2):648–
512 655, 2016.
- 513 Yang Cai, Argyris Oikonomou, and Weiqiang Zheng. Finite-time last-iterate convergence for learn-
514 ing in multi-player games. In *NeurIPS*, pp. 33904–33919, 2022a.
- 515 Yang Cai, Argyris Oikonomou, and Weiqiang Zheng. Tight last-iterate convergence of the extragra-
516 dient method for constrained monotone variational inequalities. *arXiv preprint arXiv:2204.09228*,
517 2022b.
- 518 Yang Cai, Haipeng Luo, Chen-Yu Wei, and Weiqiang Zheng. Uncoupled and convergent learning in
519 two-player zero-sum markov games with bandit feedback. In *NeurIPS*, pp. 36364–36406, 2023.
- 520 Shicong Cen, Yuting Wei, and Yuejie Chi. Fast policy extragradient methods for competitive games
521 with entropy regularization. In *NeurIPS*, pp. 27952–27964, 2021.
- 522 Shicong Cen, Yuejie Chi, Simon S Du, and Lin Xiao. Faster last-iterate convergence of policy
523 optimization in zero-sum Markov games. In *ICLR*, 2023.
- 524 Volkan Cevher, Georgios Piliouras, Ryann Sim, and Stratis Skoulakis. Min-max optimization made
525 simple: Approximating the proximal point method via contraction maps. In *Symposium on Sim-
526 plicity in Algorithms (SOSA)*, pp. 192–206, 2023.
- 527 Johanne Cohen, Amélie Héliou, and Panayotis Mertikopoulos. Learning with bandit feedback in
528 potential games. In *NeurIPS*, pp. 6372–6381, 2017.
- 529 Pierre Coucheney, Bruno Gaujal, and Panayotis Mertikopoulos. Penalty-regulated dynamics and
530 robust learning procedures in games. *Mathematics of Operations Research*, 40(3):611–633, 2015.

- 540 Constantinos Daskalakis and Ioannis Panageas. The limit points of (optimistic) gradient descent in
 541 min-max optimization. In *NeurIPS*, pp. 9256–9266, 2018.
- 542
- 543 Constantinos Daskalakis and Ioannis Panageas. Last-iterate convergence: Zero-sum games and
 544 constrained min-max optimization. In *ITCS*, pp. 27:1–27:18, 2019.
- 545 Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training GANs with
 546 optimism. In *ICLR*, 2018.
- 547
- 548 Étienne de Montbrun and Jérôme Renault. Convergence of optimistic gradient descent ascent in
 549 bilinear games. *arXiv preprint arXiv:2208.03085*, 2022.
- 550 Gerard Debreu. A social equilibrium existence theorem. *Proceedings of the National Academy of
 551 Sciences*, 38(10):886–893, 1952.
- 552
- 553 Francisco Facchinei and Jong-Shi Pang. *Finite-dimensional variational inequalities and comple-
 554 mentarity problems*. Springer, 2003.
- 555 Angeliki Giannou, Emmanuil Vasileios Vlatakis-Gkaragkounis, and Panayotis Mertikopoulos. Sur-
 556 vival of the strictest: Stable and unstable equilibria under regularized learning with partial infor-
 557 mation. In *COLT*, pp. 2147–2148, 2021a.
- 558 Angeliki Giannou, Emmanuil-Vasileios Vlatakis-Gkaragkounis, and Panayotis Mertikopoulos. On
 559 the rate of convergence of regularized learning in games: From bandits and uncertainty to opti-
 560 mism and beyond. In *NeurIPS*, pp. 22655–22666, 2021b.
- 561
- 562 Noah Golowich, Sarath Pattathil, and Constantinos Daskalakis. Tight last-iterate convergence rates
 563 for no-regret learning in multi-player games. In *NeurIPS*, pp. 20766–20778, 2020a.
- 564 Noah Golowich, Sarath Pattathil, Constantinos Daskalakis, and Asuman Ozdaglar. Last iterate is
 565 slower than averaged iterate in smooth convex-concave saddle point problems. In *COLT*, pp.
 566 1758–1784, 2020b.
- 567
- 568 Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
 569 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, pp. 2672–2680,
 570 2014.
- 571 Eduard Gorbunov, Adrien Taylor, and Gauthier Gidel. Last-iterate convergence of optimistic gradi-
 572 ent method for monotone variational inequalities. In *NeurIPS*, pp. 21858–21870, 2022.
- 573 Yu-Guan Hsieh, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. On the convergence
 574 of single-call stochastic extra-gradient methods. In *NeurIPS*, pp. 6938–6948, 2019.
- 575
- 576 Yu-Guan Hsieh, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. Explore aggres-
 577 sively, update conservatively: Stochastic extragradient methods with variable stepsize scaling. In
 578 *NeurIPS*, pp. 16223–16234, 2020.
- 579
- 580 Yu-Guan Hsieh, Kimon Antonakopoulos, and Panayotis Mertikopoulos. Adaptive learning in con-
 581 tinuous games: Optimal regret bounds and convergence to Nash equilibrium. In *COLT*, pp. 2388–
 582 2422, 2021.
- 583
- 584 Yu-Guan Hsieh, Kimon Antonakopoulos, Volkan Cevher, and Panayotis Mertikopoulos. No-regret
 585 learning in games with noisy feedback: Faster rates and adaptivity via learning rate separation. In
 586 *NeurIPS*, pp. 6544–6556, 2022.
- 587
- 588 Aamal Abbas Hussain, Francesco Belardinelli, and Georgios Piliouras. Asymptotic convergence
 589 and performance of multi-agent q-learning dynamics. In *AAMAS*, pp. 1578–1586, 2023.
- 590
- 591 Aswin Kannan and Uday V. Shanbhag. Optimal stochastic extragradient schemes for pseudomonotone
 592 stochastic variational inequality problems and their variants. *Computational Optimization and
 593 Applications*, 74(3):779–820, 2019.
- 594
- 595 Jayash Koshal, Angelia Nedić, and Uday V Shanbhag. Single timescale regularized stochastic ap-
 596 proximation schemes for monotone nash games under uncertainty. In *CDC*, pp. 231–236. IEEE,
 597 2010.

- 594 Jayash Koshal, Angelia Nedic, and Uday V. Shanbhag. Regularized iterative stochastic approxi-
 595 mation methods for stochastic variational inequality problems. *IEEE Transactions on Automatic
 596 Control*, 58(3):594–609, 2013.
- 597 Suseol Lee and Donghwan Kim. Fast extra gradient methods for smooth structured nonconvex-
 598 nonconcave minimax problems. In *NeurIPS*, pp. 22588–22600, 2021.
- 599 600 Qi Lei, Sai Ganesh Nagarajan, Ioannis Panageas, et al. Last iterate convergence in no-regret learning:
 601 constrained min-max optimization for convex-concave landscapes. In *AISTATS*, pp. 1441–1449,
 602 2021.
- 603 604 David S Leslie and Edmund J Collins. Individual q-learning in normal form games. *SIAM Journal
 605 on Control and Optimization*, 44(2):495–514, 2005.
- 606 607 Tengyuan Liang and James Stokes. Interaction matters: A note on non-asymptotic local convergence
 608 of generative adversarial networks. In *AISTATS*, pp. 907–915, 2019.
- 609 610 Tianyi Lin, Zhengyuan Zhou, Panayotis Mertikopoulos, and Michael Jordan. Finite-time last-iterate
 611 convergence for multi-agent learning in games. In *ICML*, pp. 6161–6171, 2020.
- 612 Mingyang Liu, Asuman Ozdaglar, Tiancheng Yu, and Kaiqing Zhang. The power of regularization
 613 in solving extensive-form games. In *ICLR*, 2023.
- 614 Richard D McKelvey and Thomas R Palfrey. Quantal response equilibria for normal form games.
 615 *Games and economic behavior*, 10(1):6–38, 1995.
- 616 Richard D McKelvey and Thomas R Palfrey. Quantal response equilibria for extensive form games.
 617 *Experimental economics*, 1:9–41, 1998.
- 618 Panayotis Mertikopoulos and Zhengyuan Zhou. Learning in games with continuous action sets and
 619 unknown payoff functions. *Mathematical Programming*, 173(1):465–507, 2019.
- 620 Panayotis Mertikopoulos, Bruno Lecouat, Houssam Zenati, Chuan-Sheng Foo, Vijay Chandrasekhar,
 621 and Georgios Piliouras. Optimistic mirror descent in saddle-point problems: Going
 622 the extra (gradient) mile. In *ICLR*, 2019.
- 623 Panayotis Mertikopoulos, Ya-Ping Hsieh, and Volkan Cevher. Learning in games from a stochastic
 624 approximation viewpoint. *arXiv preprint arXiv:2206.03922*, 2022.
- 625 Dov Monderer and Lloyd S Shapley. Potential games. *Games and economic behavior*, 14(1):124–
 626 143, 1996.
- 627 Remi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland,
 628 Zhaohan Daniel Guo, Yunhao Tang, Matthieu Geist, Thomas Mesnard, Côme Fiebel, Andrea
 629 Michi, Marco Selvi, Sertan Girgin, Nikola Momchev, Olivier Bachem, Daniel J Mankowitz,
 630 Doina Precup, and Bilal Piot. Nash learning from human feedback. In *ICML*, pp. 36743–36768,
 631 2024.
- 632 John Nash. Non-cooperative games. *Annals of mathematics*, pp. 286–295, 1951.
- 633 Angelia Nedić and Soomin Lee. On stochastic subgradient mirror-descent algorithm with weighted
 634 averaging. *SIAM Journal on Optimization*, 24(1):84–107, 2014.
- 635 A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to
 636 stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- 637 Yuyuan Ouyang and Yangyang Xu. Lower complexity bounds of first-order methods for convex-
 638 concave bilinear saddle-point problems. *Mathematical Programming*, 185(1):1–35, 2021.
- 639 Sarath Pattathil, Kaiqing Zhang, and Asuman Ozdaglar. Symmetric (optimistic) natural policy gra-
 640 dient for multi-agent learning with parameter convergence. In *AISTATS*, pp. 5641–5685, 2023.
- 641 Julien Perolat, Remi Munos, Jean-Baptiste Lespiau, Shayegan Omidshafiei, Mark Rowland, Pedro
 642 Ortega, Neil Burch, Thomas Anthony, David Balduzzi, Bart De Vylder, et al. From Poincaré re-
 643 currance to convergence in imperfect information games: Finding equilibrium via regularization.
 644 In *ICML*, pp. 8525–8535, 2021.

- 648 Alexander Rakhlin and Karthik Sridharan. Online learning with predictable sequences. In *COLT*,
 649 pp. 993–1019, 2013a.
 650
- 651 Sasha Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable se-
 652 quences. In *NeurIPS*, pp. 3066–3074, 2013b.
- 653 Shai Shalev-Shwartz and Yoram Singer. Convex repeated games and fenchel duality. In *NeurIPS*,
 654 pp. 1265–1272, 2006.
 655
- 656 Samuel Sokota, Ryan D’Orazio, J Zico Kolter, Nicolas Loizou, Marc Lanctot, Ioannis Mitliagkas,
 657 Noam Brown, and Christian Kroer. A unified approach to reinforcement learning, quantal re-
 658 sponse equilibria, and two-player zero-sum games. In *ICLR*, 2023.
- 659 Gokul Swamy, Christoph Dann, Rahul Kidambi, Zhiwei Steven Wu, and Alekh Agarwal. A
 660 minimaximalist approach to reinforcement learning from human feedback. *arXiv preprint*
 661 *arXiv:2401.04056*, 2024.
- 662 Tatiana Tatarenko and Maryam Kamgarpour. Learning Nash equilibria in monotone games. In *CDC*,
 663 pp. 3104–3109. IEEE, 2019.
 664
- 665 Karl Tuyls, Pieter Jan Hoen, and Bram Vanschoenwinkel. An evolutionary dynamical analysis
 666 of multi-agent learning in iterated games. *Autonomous Agents and Multi-Agent Systems*, 12(1):
 667 115–153, 2006.
- 668 Chen-Yu Wei, Chung-Wei Lee, Mengxiao Zhang, and Haipeng Luo. Linear last-iterate convergence
 669 in constrained saddle-point optimization. In *ICLR*, 2021.
 670
- 671 TaeHo Yoon and Ernest K Ryu. Accelerated algorithms for smooth convex-concave minimax prob-
 672 lems with $\mathcal{O}(1/k^2)$ rate on squared gradient norm. In *ICML*, pp. 12098–12109, 2021.
- 673 Farzad Yousefian, Angelia Nedić, and Uday V Shanbhag. On smoothing, regularization, and aver-
 674 aging in stochastic approximation methods for stochastic variational inequality problems. *Math-
 675 ematical Programming*, 165:391–431, 2017.
 676
- 677 Zhengyuan Zhou, Panayotis Mertikopoulos, Aris L Moustakas, Nicholas Bambos, and Peter Glynn.
 678 Mirror descent learning in continuous games. In *CDC*, pp. 5776–5783. IEEE, 2017.
- 679 Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In
 680 *ICML*, pp. 928–936, 2003.
 681
- 682
- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701

702 **A PROOFS FOR THEOREM 4.1**
 703

704 **A.1 PROOF OF THEOREM 4.1**
 705

706 *Proof of Theorem 4.1.* From the first-order optimality condition for π^t , we have for any $x \in \mathcal{X}$:

$$707 \quad \left\langle V(\pi^{t-1}) - \mu \left(\pi^{t-1} - \frac{k(t-1)\sigma^{k(t-1)} + \sigma^1}{k(t-1) + 1} \right) - \frac{1}{\eta} (\pi^t - \pi^{t-1}), \pi^t - x \right\rangle \geq 0,$$

710 and then $V(\pi^{t-1}) - \mu \left(\pi^{t-1} - \frac{k(t-1)\sigma^{k(t-1)} + \sigma^1}{k(t-1) + 1} \right) - \frac{1}{\eta} (\pi^t - \pi^{t-1}) \in N_{\mathcal{X}}(\pi^t)$. Thus, the tangent
 711 residual for π^t can be bounded as:

$$712 \quad r^{\tan}(\pi^t) = \min_{a \in N_{\mathcal{X}}(\pi^t)} \| -V(\pi^t) + a \| \\ 713 \quad \leq \left\| -V(\pi^t) + V(\pi^{t-1}) - \mu \left(\pi^{t-1} - \frac{k(t-1)\sigma^{k(t-1)} + \sigma^1}{k(t-1) + 1} \right) - \frac{1}{\eta} (\pi^t - \pi^{t-1}) \right\|.$$

717 Letting us define

$$718 \quad \pi_i^{\mu, \sigma^k} = \arg \max_{\pi_i \in \mathcal{X}_i} \left\{ v_i(\pi_i, \pi_{-i}^{\mu, \sigma^k}) - \frac{\mu}{2} \left\| \pi_i - \frac{k\sigma_i^k + \sigma_i^1}{k+1} \right\|^2 \right\},$$

721 then we get by triangle inequality:

$$722 \quad r^{\tan}(\pi^t) \leq \left\| -V(\pi^t) + V(\pi^{t-1}) - \frac{\mu}{k(t-1)+1} (\sigma^{k(t-1)} - \sigma^1) \right. \\ 723 \quad \quad \quad \left. - \mu(\pi^{\mu, \sigma^{k(t-1)}} - \pi^{\mu, \sigma^{k(t-1)}} + \pi^{t-1} - \sigma^{k(t-1)}) - \frac{1}{\eta} (\pi^t - \pi^{t-1}) \right\| \\ 724 \quad \leq \left\| -V(\pi^t) + V(\pi^{t-1}) \right\| + \frac{\mu}{k(t-1)+1} \left\| \sigma^{k(t-1)} - \sigma^1 \right\| \\ 725 \quad \quad \quad + \mu \left\| \pi^{\mu, \sigma^{k(t-1)}} - \sigma^{k(t-1)} \right\| + \mu \left\| \pi^{\mu, \sigma^{k(t-1)}} - \pi^{t-1} \right\| + \frac{1}{\eta} \left\| \pi^t - \pi^{t-1} \right\| \\ 726 \quad \leq \frac{1+\eta L}{\eta} \left\| \pi^t - \pi^{t-1} \right\| + \frac{\mu D}{k(t-1)+1} \\ 727 \quad \quad \quad + \mu \left\| \pi^{\mu, \sigma^{k(t-1)}} - \sigma^{k(t-1)} \right\| + \mu \left\| \pi^{\mu, \sigma^{k(t-1)}} - \pi^{t-1} \right\| \\ 728 \quad \leq \frac{1+\eta L}{\eta} \left\| \pi^{\mu, \sigma^{k(t-1)}} - \pi^t \right\| + \frac{\mu D}{k(t-1)+1} + \mu \left\| \pi^{\mu, \sigma^{k(t-1)}} - \sigma^{k(t-1)} \right\| \\ 729 \quad \quad \quad + \left(\mu + \frac{1+\eta L}{\eta} \right) \left\| \pi^{\mu, \sigma^{k(t-1)}} - \pi^{t-1} \right\|. \tag{10}$$

730 In terms of upper bound on $\left\| \pi^{\mu, \sigma^{k(t-1)}} - \pi^t \right\|$ and $\left\| \pi^{\mu, \sigma^{k(t-1)}} - \pi^{t-1} \right\|$, we introduce the following
 731 lemma:

732 **Lemma A.1.** *If we use the constant learning rate $\eta_t = \eta \in (0, \frac{\mu}{(L+\mu)^2})$, we have for any $t \geq 1$:*

$$733 \quad \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^t \right\|^2 \leq \left(\frac{1}{1+\eta\mu} \right)^{t-(k(t)-1)T_\sigma-1} \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2, \\ 734 \quad \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{t+1} \right\|^2 \leq \left(\frac{1}{1+\eta\mu} \right)^{t-(k(t)-1)T_\sigma} \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2.$$

750 Combining Eq. (10) and Lemma A.1, we have:

$$751 \quad r^{\tan}(\pi^t) \leq 2 \left(\mu + \frac{1+\eta L}{\eta} \right) \left\| \pi^{\mu, \sigma^{k(t-1)}} - \sigma^{k(t-1)} \right\| + \frac{\mu D}{k(t-1)+1}. \tag{11}$$

755 Next, we derive the following upper bound on $\left\| \pi^{\mu, \sigma^{k(t-1)}} - \sigma^{k(t-1)} \right\|$:

756 **Lemma A.2.** If we set $\eta_t = \eta \in (0, \frac{\mu}{(L+\mu)^2})$ and $T_\sigma \geq \max(1, \frac{6\ln 3(T+1)}{\ln(1+\eta\mu)})$, we have for any $t \geq 1$:

$$759 \quad \|\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}\| \leq \frac{8D}{k(t) + 1}. \\ 760$$

762 By combining Eq. (11) and Lemma A.2, we get:

$$764 \quad r^{\tan}(\pi^t) \leq \frac{16D}{k(t-1) + 1} \left(\mu + \frac{1 + \eta L}{\eta} \right) + \frac{\mu D}{k(t-1) + 1} \\ 765 \\ 766 \\ 767 \\ 768 \quad \leq \frac{17D}{k(t-1) + 1} \left(\mu + \frac{1 + \eta L}{\eta} \right). \\ 769$$

770 Therefore, since $k(t) = \lfloor \frac{t-1}{T_\sigma} \rfloor + 1$, it holds that:

$$773 \quad r^{\tan}(\pi^t) \leq \frac{17DT_\sigma}{t + T_\sigma - 2} \left(\mu + \frac{1 + \eta L}{\eta} \right). \\ 774 \\ 775$$

777 Finally, taking $T_\sigma = c \cdot \max(1, \frac{6\ln 3(T+1)}{\ln(1+\eta\mu)})$, we have:

$$779 \quad r^{\tan}(\pi^t) \leq \frac{17cD \left(\frac{6\ln 3(T+1)}{\ln(1+\eta\mu)} + 1 \right)}{t - 1} \left(\mu + \frac{1 + \eta L}{\eta} \right). \\ 780 \\ 781 \\ 782$$

□

786 A.2 PROOF OF LEMMA A.1

788 *Proof of Lemma A.1.* First, we have for any three vectors a, b, c :

$$790 \quad \frac{1}{2} \|a - b\|^2 - \frac{1}{2} \|a - c\|^2 + \frac{1}{2} \|b - c\|^2 = \langle c - b, a - b \rangle. \\ 791 \\ 792$$

793 Thus, we have for any $t \geq 1$:

$$795 \quad \frac{1}{2} \|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 - \frac{1}{2} \|\pi^{\mu, \sigma^{k(t)}} - \pi^t\|^2 + \frac{1}{2} \|\pi^{t+1} - \pi^t\|^2 = \langle \pi^t - \pi^{t+1}, \pi^{\mu, \sigma^{k(t)}} - \pi^{t+1} \rangle. \\ 796 \\ 797 \quad (12)$$

800 Here, let us define $\hat{\sigma}^{k(t)} = \frac{k(t)\sigma^{k(t)} + \sigma^1}{k(t) + 1}$. Then, from the first-order optimality condition for π^{t+1} ,
801 we have for any $t \geq 1$:

$$803 \quad \langle \eta \left(V(\pi^t) - \mu \left(\pi^t - \hat{\sigma}^{k(t)} \right) \right) - \pi^{t+1} + \pi^t, \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle \geq 0. \quad (13) \\ 804 \\ 805$$

806 Similarly, from the first-order optimality condition for $\pi^{\mu, \sigma^{k(t)}}$, we get:

$$807 \quad \langle V(\pi^{\mu, \sigma^{k(t)}}) - \mu \left(\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right), \pi^{\mu, \sigma^{k(t)}} - \pi^{t+1} \rangle \geq 0. \quad (14) \\ 808 \\ 809$$

810 Combining Eq. (12), Eq. (13), and Eq. (14) yields:
 811

$$\begin{aligned}
 & \frac{1}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{t+1} \right\|^2 - \frac{1}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^t \right\|^2 + \frac{1}{2} \left\| \pi^{t+1} - \pi^t \right\|^2 \\
 & \leq \eta \left\langle V(\pi^t) - \mu \left(\pi^t - \hat{\sigma}^{k(t)} \right), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
 & = \eta \left\langle V(\pi^{t+1}) - \mu \left(\pi^{t+1} - \hat{\sigma}^{k(t)} \right), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
 & \quad + \eta \left\langle V(\pi^t) - V(\pi^{t+1}) - \mu \left(\pi^t - \pi^{t+1} \right), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
 & \leq \eta \left\langle V(\pi^{\mu, \sigma^{k(t)}}) - \mu \left(\pi^{t+1} - \hat{\sigma}^{k(t)} \right), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
 & \quad + \eta \left\langle V(\pi^t) - V(\pi^{t+1}) - \mu \left(\pi^t - \pi^{t+1} \right), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
 & = \eta \left\langle V(\pi^{\mu, \sigma^{k(t)}}) - \mu \left(\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\rangle - \eta \mu \left\| \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 \\
 & \quad + \eta \left\langle V(\pi^t) - V(\pi^{t+1}) - \mu \left(\pi^t - \pi^{t+1} \right), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
 & \leq -\eta \mu \left\| \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 + \eta \left\langle V(\pi^t) - V(\pi^{t+1}) - \mu \left(\pi^t - \pi^{t+1} \right), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\rangle, \quad (15)
 \end{aligned}$$

829 where the second inequality follows from Eq. (1). From Cauchy-Schwarz inequality and Young's
 830 inequality, the second term in the right-hand side of this inequality can be bounded by:
 831

$$\begin{aligned}
 & \eta \left\langle V(\pi^t) - V(\pi^{t+1}) - \mu \left(\pi^t - \pi^{t+1} \right), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
 & = \eta \left\langle V(\pi^t) - V(\pi^{t+1}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\rangle - \eta \mu \left\langle \pi^t - \pi^{t+1}, \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
 & \leq \eta (\|V(\pi^t) - V(\pi^{t+1})\| + \mu \|\pi^t - \pi^{t+1}\|) \cdot \left\| \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\| \\
 & \leq \eta (L + \mu) \|\pi^t - \pi^{t+1}\| \cdot \left\| \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\| \\
 & \leq \frac{1}{2} \|\pi^t - \pi^{t+1}\|^2 + \frac{\eta^2 (L + \mu)^2}{2} \left\| \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 \\
 & \leq \frac{1}{2} \|\pi^t - \pi^{t+1}\|^2 + \frac{\eta \mu}{2} \left\| \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2,
 \end{aligned} \quad (16)$$

844 where the second inequality follow from Eq. (2), and the last inequality follows from the assumption
 845 that $\eta \leq \frac{\mu}{(L + \mu)^2}$. By combining Eq. (15) and Eq. (16), we get:
 846

$$\begin{aligned}
 & \frac{1}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{t+1} \right\|^2 - \frac{1}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^t \right\|^2 + \frac{1}{2} \left\| \pi^{t+1} - \pi^t \right\|^2 \\
 & \leq -\frac{\eta \mu}{2} \left\| \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 + \frac{1}{2} \left\| \pi^t - \pi^{t+1} \right\|^2.
 \end{aligned}$$

851 Thus,
 852

$$\frac{1 + \eta \mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{t+1} \right\|^2 \leq \frac{1}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^t \right\|^2.$$

855 Therefore, we have for any $t \geq 1$:

$$\left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{t+1} \right\|^2 \leq \frac{1}{1 + \eta \mu} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^t \right\|^2.$$

860 Furthermore, since $k(s) = k(t)$ for $s \in [(k(t) - 1)T_\sigma + 1, t]$, we have for such s that:
 861

$$\left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{s+1} \right\|^2 \leq \frac{1}{1 + \eta \mu} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^s \right\|^2.$$

Therefore, by applying this inequality from $t, t-1, \dots, (k(t)-1)T_\sigma + 1$, we get for any $t \geq 1$:

$$\begin{aligned} \|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 &\leq \left(\frac{1}{1+\eta\mu}\right)^{t-(k(t)-1)T_\sigma} \|\pi^{\mu, \sigma^{k(t)}} - \pi^{(k(t)-1)T_\sigma + 1}\|^2 \\ &= \left(\frac{1}{1+\eta\mu}\right)^{t-(k(t)-1)T_\sigma} \|\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}\|^2. \end{aligned} \quad (17)$$

Here, since $k(t) = k(t+1)$ when t satisfies that $t \neq T_\sigma \left\lfloor \frac{t}{T_\sigma} \right\rfloor$, we have for such t that:

$$\|\pi^{\mu, \sigma^{k(t+1)}} - \pi^{t+1}\|^2 \leq \left(\frac{1}{1+\eta\mu}\right)^{t-(k(t+1)-1)T_\sigma} \|\pi^{\mu, \sigma^{k(t+1)}} - \sigma^{k(t+1)}\|^2. \quad (18)$$

On the other hand, when t satisfies that $t = T_\sigma \left\lfloor \frac{t}{T_\sigma} \right\rfloor$:

$$\begin{aligned} k(t+1) &= \left\lfloor \frac{T_\sigma \left\lfloor \frac{t}{T_\sigma} \right\rfloor + 1 - 1}{T_\sigma} \right\rfloor + 1 = \left\lfloor \frac{t}{T_\sigma} \right\rfloor + 1 \\ \Rightarrow (k(t+1) - 1)T_\sigma &= T_\sigma \left\lfloor \frac{t}{T_\sigma} \right\rfloor = t \\ \Rightarrow \pi^{t+1} &= \pi^{(k(t+1)-1)T_\sigma + 1} = \sigma^{k(t+1)}. \end{aligned}$$

Therefore, we have for any $t \geq 1$ such that $t = T_\sigma \left\lfloor \frac{t}{T_\sigma} \right\rfloor$:

$$\begin{aligned} \|\pi^{\mu, \sigma^{k(t+1)}} - \pi^{t+1}\|^2 &= \|\pi^{\mu, \sigma^{k(t+1)}} - \sigma^{k(t+1)}\|^2 \\ &= \left(\frac{1}{1+\eta\mu}\right)^{t-(k(t+1)-1)T_\sigma} \|\pi^{\mu, \sigma^{k(t+1)}} - \sigma^{k(t+1)}\|^2. \end{aligned} \quad (19)$$

By combining Eq. (17), Eq. (18), and Eq. (19), we have for any $t \geq 1$:

$$\begin{aligned} \|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 &\leq \left(\frac{1}{1+\eta\mu}\right)^{t-(k(t)-1)T_\sigma} \|\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}\|^2, \\ \|\pi^{\mu, \sigma^{k(t+1)}} - \pi^{t+1}\|^2 &\leq \left(\frac{1}{1+\eta\mu}\right)^{t-(k(t+1)-1)T_\sigma} \|\pi^{\mu, \sigma^{k(t+1)}} - \sigma^{k(t+1)}\|^2. \end{aligned}$$

□

A.3 PROOF OF LEMMA A.2

Proof of Lemma A.2. First, we have for any Nash equilibrium $\pi^* \in \Pi^*$ and $t \geq 1$ such that $k(t) \geq 1$:

$$\begin{aligned} &\frac{(k(t)+1)(k(t)+2)}{2} \|\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}\|^2 + (k(t)+1)(k(t)+2) \langle \hat{\sigma}^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \rangle \\ &= \frac{(k(t)+1)(k(t)+2)}{2} \|\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}\|^2 \\ &\quad + (k(t)+1) \langle (k(t)+1)\sigma^{k(t)+1} + \sigma^1 - (k(t)+2)\pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \rangle \\ &= \frac{(k(t)+1)(k(t)+2)}{2} \|\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}\|^2 + (k(t)+1) \langle \sigma^1 - \sigma^{k(t)+1}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \rangle \\ &\quad + (k(t)+1)(k(t)+2) \langle \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \rangle \\ &= \frac{(k(t)+1)(k(t)+2)}{2} \|\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}\|^2 + (k(t)+1) \langle \sigma^1 - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \rangle \end{aligned}$$

$$\begin{aligned}
& + (k(t) + 1)^2 \left\langle \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle \\
& = \frac{(k(t) + 1)(k(t) + 2)}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 + (k(t) + 1) \left\langle \sigma^1 - \pi^*, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle \\
& \quad + (k(t) + 1) \left\langle \pi^* - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle + (k(t) + 1)^2 \left\langle \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle.
\end{aligned}$$

Here, the first-order optimality condition for $\pi^{\mu, \sigma^{k(t)}}$:

$$\begin{aligned}
& \left\langle V(\pi^{\mu, \sigma^{k(t)}}) - \mu \left(\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right), \pi^{\mu, \sigma^{k(t)}} - \pi^* \right\rangle \geq 0 \\
& \Rightarrow \left\langle \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}, \pi^* - \pi^{\mu, \sigma^{k(t)}} \right\rangle \geq \frac{1}{\mu} \left\langle V(\pi^{\mu, \sigma^{k(t)}}), \pi^* - \pi^{\mu, \sigma^{k(t)}} \right\rangle \geq \frac{1}{\mu} \left\langle V(\pi^*), \pi^* - \pi^{\mu, \sigma^{k(t)}} \right\rangle \geq 0,
\end{aligned}$$

where we use Eq. (1) and the fact that π^* is a Nash equilibrium. Combining these inequalities yields:

$$\begin{aligned}
& \frac{(k(t) + 1)(k(t) + 2)}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 + (k(t) + 1)(k(t) + 2) \left\langle \hat{\sigma}^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle \\
& \geq \frac{(k(t) + 1)(k(t) + 2)}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 + (k(t) + 1) \left\langle \sigma^1 - \pi^*, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle \\
& \quad + (k(t) + 1)^2 \left\langle \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle.
\end{aligned}$$

From Young's inequality, we have for any $\rho_1, \rho_2 > 0$:

$$\begin{aligned}
& \frac{(k(t) + 1)(k(t) + 2)}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 + (k(t) + 1)(k(t) + 2) \left\langle \hat{\sigma}^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle \\
& \geq \frac{(k(t) + 1)(k(t) + 2)}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 - \frac{\rho_1(k(t) + 1)}{2} \left\| \sigma^1 - \pi^* \right\|^2 - \frac{(k(t) + 1)}{2\rho_1} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 \\
& \quad - \frac{\rho_2(k(t) + 1)^2}{2} \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 - \frac{(k(t) + 1)^2}{2\rho_2} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 \\
& = \left(\frac{(k(t) + 1)(k(t) + 2)}{2} - \frac{k(t) + 1}{2\rho_1} - \frac{(k(t) + 1)^2}{2\rho_2} \right) \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 \\
& \quad - \frac{\rho_1(k(t) + 1)}{2} \left\| \sigma^1 - \pi^* \right\|^2 - \frac{\rho_2(k(t) + 1)^2}{2} \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2.
\end{aligned}$$

By setting $\rho_1 = \frac{4}{k(t)+2}$, $\rho_2 = \frac{4(k(t)+1)}{k(t)+2}$, we obtain:

$$\begin{aligned}
& \frac{(k(t) + 1)(k(t) + 2)}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 + (k(t) + 1)(k(t) + 2) \left\langle \hat{\sigma}^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle \\
& \geq \frac{(k(t) + 1)(k(t) + 2)}{4} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 - \frac{2(k(t) + 1)}{k(t) + 2} \left\| \sigma^1 - \pi^* \right\|^2 \\
& \quad - \frac{2(k(t) + 1)^3}{k(t) + 2} \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 \\
& \geq \frac{(k(t) + 1)(k(t) + 2)}{4} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 - 2 \left\| \sigma^1 - \pi^* \right\|^2 - 2(k(t) + 1)^2 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2. \tag{20}
\end{aligned}$$

Here, we introduce the following lemma:

Lemma A.3. For any $t \geq 1$ such that $k(t) \geq 2$, we have:

$$\begin{aligned}
& \frac{(k(t) + 1)(k(t) + 2)}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 + (k(t) + 1)(k(t) + 2) \left\langle \hat{\sigma}^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle \\
& \leq \frac{k(t)(k(t) + 1)}{2} \left\| \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\|^2 + k(t)(k(t) + 1) \left\langle \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\rangle \\
& \quad + (k(t) + 1) \left\langle (k(t) + 1)(\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)+1}) + k(t)(\sigma^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}), \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle.
\end{aligned}$$

972 By combining Eq. (20) and Lemma A.3, we get:
 973

$$\begin{aligned}
 & \frac{(k(t)+1)(k(t)+2)}{4} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 \\
 & \leq \frac{(k(t)+1)(k(t)+2)}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 + (k(t)+1)(k(t)+2) \left\langle \hat{\sigma}^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle \\
 & \quad + 2 \|\sigma^1 - \pi^*\|^2 + 2(k(t)+1)^2 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 \\
 & \leq 3 \left\| \pi^{\mu, \sigma^1} - \hat{\sigma}^1 \right\|^2 + 6 \left\langle \hat{\sigma}^2 - \pi^{\mu, \sigma^1}, \pi^{\mu, \sigma^1} - \hat{\sigma}^1 \right\rangle + 2 \|\sigma^1 - \pi^*\|^2 + 2(k(t)+1)^2 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 \\
 & \quad + \sum_{l=2}^{k(t)} (l+1) \left\langle (l+1)(\pi^{\mu, \sigma^l} - \sigma^{l+1}) + l(\sigma^l - \pi^{\mu, \sigma^{l-1}}), \hat{\sigma}^l - \pi^{\mu, \sigma^l} \right\rangle \\
 & = 3 \left\| \pi^{\mu, \sigma^1} - \hat{\sigma}^1 \right\|^2 + 2 \left\langle 2\sigma^2 + \sigma^1 - 3\pi^{\mu, \sigma^1}, \pi^{\mu, \sigma^1} - \hat{\sigma}^1 \right\rangle + 2 \|\sigma^1 - \pi^*\|^2 \\
 & \quad + 2(k(t)+1)^2 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 + \sum_{l=2}^{k(t)} (l+1) \left\langle (l+1)(\pi^{\mu, \sigma^l} - \sigma^{l+1}) + l(\sigma^l - \pi^{\mu, \sigma^{l-1}}), \hat{\sigma}^l - \pi^{\mu, \sigma^l} \right\rangle \\
 & = 3 \left\| \pi^{\mu, \sigma^1} - \hat{\sigma}^1 \right\|^2 + 2 \left\langle \sigma^1 - \pi^{\mu, \sigma^1}, \pi^{\mu, \sigma^1} - \hat{\sigma}^1 \right\rangle + 4 \left\langle \sigma^2 - \pi^{\mu, \sigma^1}, \pi^{\mu, \sigma^1} - \hat{\sigma}^1 \right\rangle \\
 & \quad + 2 \|\sigma^1 - \pi^*\|^2 + 2(k(t)+1)^2 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 \\
 & \quad + \sum_{l=2}^{k(t)} (l+1) \left\langle (l+1)(\pi^{\mu, \sigma^l} - \sigma^{l+1}) + l(\sigma^l - \pi^{\mu, \sigma^{l-1}}), \hat{\sigma}^l - \pi^{\mu, \sigma^l} \right\rangle \\
 & = \left\| \pi^{\mu, \sigma^1} - \hat{\sigma}^1 \right\|^2 + 4 \left\langle \sigma^2 - \pi^{\mu, \sigma^1}, \pi^{\mu, \sigma^1} - \hat{\sigma}^1 \right\rangle + 2 \|\sigma^1 - \pi^*\|^2 + 2(k(t)+1)^2 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 \\
 & \quad + \sum_{l=2}^{k(t)} (l+1) \left\langle (l+1)(\pi^{\mu, \sigma^l} - \sigma^{l+1}) + l(\sigma^l - \pi^{\mu, \sigma^{l-1}}), \hat{\sigma}^l - \pi^{\mu, \sigma^l} \right\rangle \\
 & = \left\| \pi^{\mu, \sigma^1} - \hat{\sigma}^1 \right\|^2 + 2 \|\sigma^1 - \pi^*\|^2 + 2(k(t)+1)^2 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 \\
 & \quad + \sum_{l=1}^{k(t)} (l+1)^2 \left\langle \pi^{\mu, \sigma^l} - \sigma^{l+1}, \hat{\sigma}^l - \pi^{\mu, \sigma^l} \right\rangle + \sum_{l=2}^{k(t)} l(l+1) \left\langle \sigma^l - \pi^{\mu, \sigma^{l-1}}, \hat{\sigma}^l - \pi^{\mu, \sigma^l} \right\rangle \\
 & \leq 3D^2 + 2(k(t)+1)^2 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 + 2D(k(t)+1)^2 \sum_{l=1}^{k(t)} \left\| \pi^{\mu, \sigma^l} - \sigma^{l+1} \right\|.
 \end{aligned}$$

1012 Therefore, we have for any $t \geq 1$ such that $k(t) \geq 2$:

$$\left\| \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\|^2 \leq \frac{12D^2}{(k(t)+1)^2} + 8 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 + 8D \sum_{l=1}^{k(t)} \left\| \pi^{\mu, \sigma^l} - \sigma^{l+1} \right\|.$$

1019 By the definition of $\hat{\sigma}^{k(t)}$,

$$\begin{aligned}
 & \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 + \frac{\left\| \sigma^{k(t)} - \sigma^1 \right\|^2}{(k(t)+1)^2} + \frac{2}{k(t)+1} \left\langle \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}, \sigma^{k(t)} - \sigma^1 \right\rangle \\
 & \leq \frac{12D^2}{(k(t)+1)^2} + 8 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 + 8D \sum_{l=1}^{k(t)} \left\| \pi^{\mu, \sigma^l} - \sigma^{l+1} \right\|.
 \end{aligned}$$

1026 Therefore, from Cauchy-Schwarz inequality, we have:
 1027

$$\begin{aligned}
 & \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 \\
 & \leq \frac{2}{k(t) + 1} \left\langle \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}, \sigma^1 - \sigma^{k(t)} \right\rangle + \frac{12D^2}{(k(t) + 1)^2} \\
 & \quad + 8 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 + 8D \sum_{l=1}^{k(t)} \left\| \pi^{\mu, \sigma^l} - \sigma^{l+1} \right\|^2 \\
 & \leq \frac{2D}{k(t) + 1} \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 + \frac{12D^2}{(k(t) + 1)^2} + 8 \left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 + 8D \sum_{l=1}^{k(t)} \left\| \pi^{\mu, \sigma^l} - \sigma^{l+1} \right\|^2.
 \end{aligned} \tag{21}$$

1039 Furthermore, from Lemma A.1, we have for any $k \geq 1$:

$$\left\| \pi^{\mu, \sigma^k} - \sigma^{k+1} \right\|^2 \leq \left(\frac{1}{1 + \eta\mu} \right)^{T_\sigma} \left\| \pi^{\mu, \sigma^k} - \sigma^k \right\|^2. \tag{22}$$

1044 Combining Eq. (21) nad Eq. (22), we have for any $t \geq 1$ such that $k(t) \geq 2$:

$$\begin{aligned}
 \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 & \leq \frac{2D}{k(t) + 1} \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 + \frac{12D^2}{(k(t) + 1)^2} \\
 & \quad + 8 \left(\frac{1}{1 + \eta\mu} \right)^{T_\sigma} \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 + 8D^2 k(t) \left(\frac{1}{1 + \eta\mu} \right)^{\frac{T_\sigma}{2}}.
 \end{aligned}$$

1052 Therefore, since $T_\sigma \geq \max(1, \frac{6 \ln 3(T+1)}{\ln(1+\eta\mu)}) \Rightarrow \left(\frac{1}{1 + \eta\mu} \right)^{T_\sigma} \leq \frac{(k(t)+1)^3}{(1+\eta\mu)^{T_\sigma}} \leq \frac{1}{16}$, we have for $k(t) \geq 2$:

$$\frac{1}{2} \left(\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 - \frac{2D}{k(t) + 1} \right)^2 \leq \frac{2D^2}{(k(t) + 1)^2} + \frac{12D^2}{(k(t) + 1)^2} + \frac{D^2}{2(k(t) + 1)^2} \leq \frac{16D^2}{(k(t) + 1)^2},$$

1056 and then:

$$\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| \leq \frac{2D}{k(t) + 1} + \frac{4\sqrt{2}D}{k(t) + 1} \leq \frac{8D}{k(t) + 1}.$$

1061 On the other hand, for $k(t) = 1$, we have:

$$\left\| \pi^{\mu, \sigma^1} - \sigma^1 \right\| \leq D \leq \frac{8D}{1 + 1}.$$

1065 In summary, for any $t \geq 1$, we have:

$$\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| \leq \frac{8D}{k(t) + 1}.$$

1069 \square

1071 A.4 PROOF OF LEMMA A.3

1073 *Proof of Lemma A.3.* From the first-order optimality condition for $\pi^{\mu, \sigma^{k(t)}}$, we have:

$$\left\langle V(\pi^{\mu, \sigma^{k(t)}}) - \mu(\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}), \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\rangle \geq 0.$$

1077 Similarly, from the first-order optimality condition for $\pi^{\mu, \sigma^{k(t)-1}}$, we have:

$$\left\langle V(\pi^{\mu, \sigma^{k(t)-1}}) - \mu(\pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1}), \pi^{\mu, \sigma^{k(t)-1}} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \geq 0.$$

1080 Summing up these inequalities, we get for any $t \geq 1$ such that $k(t) \geq 2$:

$$\begin{aligned}
1082 \quad & 0 \leq \left\langle V(\pi^{\mu, \sigma^{k(t)}}) - V(\pi^{\mu, \sigma^{k(t)-1}}), \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\rangle - \mu \left\langle \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}, \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\rangle \\
1083 \quad & + \mu \left\langle \hat{\sigma}^{k(t)-1} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
1084 \quad & \leq -\mu \left\langle \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}, \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\rangle + \mu \left\langle \hat{\sigma}^{k(t)-1} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)-1}} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
1085 \quad & = -\mu \left\langle \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} + \sigma^{k(t)} - \hat{\sigma}^{k(t)}, \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} + \sigma^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}} \right\rangle \\
1086 \quad & + \mu \left\langle \hat{\sigma}^{k(t)-1} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)-1}} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
1087 \quad & = -\mu \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 - \mu \left\langle \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}, \sigma^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}} \right\rangle \\
1088 \quad & - \mu \left\langle \sigma^{k(t)} - \hat{\sigma}^{k(t)}, \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\rangle + \mu \left\langle \hat{\sigma}^{k(t)-1} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)-1}} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
1089 \quad & = -\mu \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 - \mu \left\langle \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}, \sigma^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}} \right\rangle \\
1090 \quad & + \mu \left\langle \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)-1}} - \sigma^{k(t)} \right\rangle + \mu \left\langle \hat{\sigma}^{k(t)-1} - \hat{\sigma}^{k(t)}, \pi^{\mu, \sigma^{k(t)-1}} - \pi^{\mu, \sigma^{k(t)}} \right\rangle. \\
1091
\end{aligned}$$

1099 Here, for any vectors a, b, c , it holds that:

$$\begin{aligned}
1100 \quad & \langle a - b, b - c \rangle = \frac{1}{2} \|a - c\|^2 - \frac{1}{2} \|b - c\|^2 - \frac{1}{2} \|a - b\|^2, \\
1101 \quad & \langle a - b, c - d \rangle = \frac{1}{2} \|a - b\|^2 + \frac{1}{2} \|c - d\|^2 - \frac{1}{2} \|d - c + a - b\|^2. \\
1102 \\
1103 \\
1104
\end{aligned}$$

1105 Thus, we have:

$$\begin{aligned}
1106 \quad & 0 \leq -\mu \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 \\
1107 \quad & + \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)-1}} - \sigma^{k(t)} \right\|^2 + \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 \\
1108 \quad & + \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)-1}} - \sigma^{k(t)} \right\|^2 - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 \\
1109 \quad & + \frac{\mu}{2} \left\| \hat{\sigma}^{k(t)-1} - \hat{\sigma}^{k(t)} \right\|^2 + \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 \\
1110 \quad & - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} + \hat{\sigma}^{k(t)-1} + \hat{\sigma}^{k(t)} \right\|^2 \\
1111 \quad & = -\frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 + \frac{\mu}{2} \left\| \hat{\sigma}^{k(t)} - \hat{\sigma}^{k(t)-1} \right\|^2 \\
1112 \quad & - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} + \hat{\sigma}^{k(t)-1} + \hat{\sigma}^{k(t)} \right\|^2 \\
1113 \quad & \leq -\frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 + \frac{\mu}{2} \left\| \hat{\sigma}^{k(t)} - \hat{\sigma}^{k(t)-1} \right\|^2 \\
1114 \quad & = -\frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 + \frac{\mu}{2} \left\| \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}} + \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\|^2 \\
1115 \quad & = \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\|^2 + \frac{\mu}{2} \left\| \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 \\
1116 \quad & + \mu \left\langle \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\rangle. \\
1117 \\
1118 \\
1119 \\
1120 \\
1121 \\
1122 \\
1123 \\
1124 \\
1125 \\
1126 \\
1127 \\
1128
\end{aligned} \tag{23}$$

1129 Here, from the definition of $\hat{\sigma}^{k(t)}$, we have:

$$\begin{aligned}
1130 \quad & \frac{1}{2} \left\| \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 - \frac{1}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 \\
1131 \quad & = \frac{1}{2} \left\| \frac{k(t)\sigma^{k(t)} + \sigma^1}{k(t) + 1} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 - \frac{1}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 \\
1132 \\
1133
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \left\langle \frac{k(t)\sigma^{k(t)} + \sigma^1}{k(t) + 1} - \pi^{\mu, \sigma^{k(t)-1}} + \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}}, \frac{k(t)\sigma^{k(t)} + \sigma^1}{k(t) + 1} - \pi^{\mu, \sigma^{k(t)-1}} - \pi^{\mu, \sigma^{k(t)}} + \pi^{\mu, \sigma^{k(t)-1}} \right\rangle \\
&= \frac{1}{2} \left\langle \frac{\sigma^1 + (k(t) + 1)\pi^{\mu, \sigma^{k(t)}} - 2(k(t) + 1)\pi^{\mu, \sigma^{k(t)-1}} + k(t)\sigma^{k(t)}}{k(t) + 1}, \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
&= \frac{1}{2k(t)} \left\langle 2(k(t) + 1)\sigma^{k(t)+1} + 2\sigma^1 - 2(k(t) + 2)\pi^{\mu, \sigma^{k(t)}}, \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
&\quad + \frac{1}{2k(t)} \left\langle -\frac{k(t) + 2}{k(t) + 1}\sigma^1 + (3k(t) + 4)\pi^{\mu, \sigma^{k(t)}} - 2(k(t) + 1)\sigma^{k(t)+1}\hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
&\quad + \frac{1}{2k(t)} \left\langle -2k(t)\pi^{\mu, \sigma^{k(t)-1}} + \frac{k(t)^2}{k(t) + 1}\sigma^{k(t)}, \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
&= \frac{k(t) + 2}{k(t)} \left\langle \hat{\sigma}^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
&\quad + \frac{1}{2k(t)} \left\langle -\frac{k(t) + 2}{k(t) + 1}\sigma^1 - \frac{k(t)(k(t) + 2)}{k(t) + 1}\sigma^{k(t)} + (k(t) + 2)\pi^{\mu, \sigma^{k(t)}}, \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
&\quad + \frac{1}{2k(t)} \left\langle 2(k(t) + 1)(\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)+1}) + 2k(t)(\sigma^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}), \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
&= -\frac{k(t) + 2}{k(t)} \left\langle \hat{\sigma}^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle \\
&\quad - \frac{k(t) + 2}{2k(t)} \left\langle \frac{k(t)\sigma^{k(t)} + \sigma^1}{k(t) + 1} - \pi^{\mu, \sigma^{k(t)}}, \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
&\quad + \frac{1}{k(t)} \left\langle (k(t) + 1)(\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)+1}) + k(t)(\sigma^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}), \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle \\
&= -\frac{k(t) + 2}{k(t)} \left\langle \hat{\sigma}^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle - \frac{k(t) + 2}{2k(t)} \|\hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}}\|^2 \\
&\quad + \frac{1}{k(t)} \left\langle (k(t) + 1)(\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)+1}) + k(t)(\sigma^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}), \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle. \tag{24}
\end{aligned}$$

Combining Eq. (23) and Eq. (24) yields for any $t \geq 1$ such that $k(t) \geq 2$:

$$\begin{aligned}
&\frac{k(t) + 2}{2k(t)} \|\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}\|^2 + \frac{k(t) + 2}{k(t)} \left\langle \hat{\sigma}^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle \\
&\leq \frac{1}{2} \|\pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1}\|^2 + \left\langle \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\rangle \\
&\quad + \frac{1}{k(t)} \left\langle (k(t) + 1)(\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)+1}) + k(t)(\sigma^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}), \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle.
\end{aligned}$$

Multiplying both sides by $k(t)(k(t) + 1)$, we have:

$$\begin{aligned}
&\frac{(k(t) + 1)(k(t) + 2)}{2} \|\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}\|^2 + (k(t) + 1)(k(t) + 2) \left\langle \hat{\sigma}^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}}, \pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)} \right\rangle \\
&\leq \frac{k(t)(k(t) + 1)}{2} \|\pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1}\|^2 + k(t)(k(t) + 1) \left\langle \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\rangle \\
&\quad + (k(t) + 1) \left\langle (k(t) + 1)(\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)+1}) + k(t)(\sigma^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}), \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)}} \right\rangle.
\end{aligned}$$

□

1188 **B PROOFS FOR THEOREM 4.3**
 1189

1190 **B.1 PROOF OF THEOREM 4.3**
 1191

1192 *Proof of Theorem 4.3.* Let us define $K := \frac{T}{T_\sigma}$. We can decompose the gap function for π^{T+1} as
 1193 follows:

$$\begin{aligned}
 & \text{GAP}(\pi^{T+1}) \\
 &= \max_{x \in \mathcal{X}} \langle V(\pi^{T+1}), x - \pi^{T+1} \rangle \\
 &= \max_{x \in \mathcal{X}} \left(\langle V(\pi^{\mu, \sigma^K}), x - \pi^{\mu, \sigma^K} \rangle - \langle V(\pi^{\mu, \sigma^K}), x - \pi^{\mu, \sigma^K} \rangle + \langle V(\pi^{T+1}), x - \pi^{T+1} \rangle \right) \\
 &= \max_{x \in \mathcal{X}} \left(\langle V(\pi^{\mu, \sigma^K}), x - \pi^{\mu, \sigma^K} \rangle - \langle V(\pi^{\mu, \sigma^K}) - V(\pi^{T+1}), x - \pi^{T+1} \rangle + \langle V(\pi^{\mu, \sigma^K}), \pi^{\mu, \sigma^K} - \pi^{T+1} \rangle \right) \\
 &\leq \max_{x \in \mathcal{X}} \left(\langle V(\pi^{\mu, \sigma^K}), x - \pi^{\mu, \sigma^K} \rangle + D \|V(\pi^{\mu, \sigma^K}) - V(\pi^{T+1})\| + \zeta \|\pi^{\mu, \sigma^K} - \pi^{T+1}\| \right) \\
 &\leq \text{GAP}(\pi^{\mu, \sigma^K}) + (LD + \zeta) \|\pi^{\mu, \sigma^K} - \pi^{T+1}\| \\
 &\leq D \cdot \min_{c \in N_{\mathcal{X}}(\pi^{\mu, \sigma^K})} \| -V(\pi^{\mu, \sigma^K}) + c \| + (LD + \zeta) \|\pi^{\mu, \sigma^K} - \pi^{T+1}\|,
 \end{aligned}$$

1207 where the last inequality follows from Lemma 2.3. From the first-order optimality condition for
 1208 π^{μ, σ^K} , we have for any $x \in \mathcal{X}$:

$$\left\langle V(\pi^{\mu, \sigma^K}) - \mu \left(\pi^{\mu, \sigma^K} - \frac{K\sigma^K + \sigma^1}{K+1} \right), \pi^{\mu, \sigma^K} - x \right\rangle \geq 0,$$

1212 and then $V(\pi^{\mu, \sigma^K}) - \mu \left(\pi^{\mu, \sigma^K} - \frac{K\sigma^K + \sigma^1}{K+1} \right) \in N_{\mathcal{X}}(\pi^{\mu, \sigma^K})$. Thus, the gap function for π^{T+1} can
 1213 be bounded by:

$$\begin{aligned}
 \text{GAP}(\pi^{T+1}) &\leq \mu D \cdot \left\| \pi^{\mu, \sigma^K} - \frac{K\sigma^K + \sigma^1}{K+1} \right\| + (LD + \zeta) \|\pi^{\mu, \sigma^K} - \pi^{T+1}\| \\
 &= \mu D \cdot \left\| \frac{\sigma^K - \sigma^1}{K+1} + \pi^{\mu, \sigma^K} - \sigma^K \right\| + (LD + \zeta) \|\pi^{\mu, \sigma^K} - \pi^{T+1}\| \\
 &\leq \mu D \cdot \left(\frac{D}{K+1} + \|\pi^{\mu, \sigma^K} - \sigma^K\| \right) + (LD + \zeta) \|\pi^{\mu, \sigma^K} - \pi^{T+1}\|.
 \end{aligned}$$

1222 Taking its expectation yields:

$$\begin{aligned}
 \mathbb{E} [\text{GAP}(\pi^{T+1})] &\leq \frac{\mu D^2}{K+1} + \mu D \cdot \mathbb{E} [\|\pi^{\mu, \sigma^K} - \sigma^K\|] + (LD + \zeta) \cdot \mathbb{E} [\|\pi^{\mu, \sigma^K} - \pi^{T+1}\|] \\
 &\leq \frac{\mu D^2}{K+1} + \mu D \cdot \mathbb{E} [\|\pi^{\mu, \sigma^K} - \sigma^K\|] + (LD + \zeta) \cdot \sqrt{\mathbb{E} [\|\pi^{\mu, \sigma^K} - \pi^{T+1}\|^2]}.
 \end{aligned} \tag{25}$$

1230 Here, we derive the following upper bound on $\mathbb{E} [\|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2]$:

1232 **Lemma B.1.** Let $\kappa = \frac{\mu}{2}, \theta = \frac{3\mu^2 + 8L^2}{2\mu}$. Suppose that Assumption 4.2 holds. If we set $\eta_t =$
 1233 $\frac{1}{\kappa(t-T_\sigma(k(t)-1))+2\theta}$, we have for any $t \geq 1$:

$$\mathbb{E} [\|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2] \leq \frac{2\theta}{\kappa(t-(k(t)-1)T_\sigma)+2\theta} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa(t-(k(t)-1)T_\sigma)}{2\theta} + 1 \right) \right).$$

1238 Setting $t = T = KT_\sigma$, we can write $k(t) = \lfloor \frac{KT_\sigma-1}{T_\sigma} \rfloor + 1 = K$. Therefore, from Lemma B.1, we
 1239 have:

$$\mathbb{E} [\|\pi^{\mu, \sigma^K} - \pi^{T+1}\|^2] \leq \frac{2\theta}{\kappa T_\sigma + 2\theta} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T_\sigma}{2\theta} + 1 \right) \right). \tag{26}$$

On the other hand, in terms of $\mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| \right]$, we introduce the following lemma:

Lemma B.2. If we set $\eta_t = \frac{1}{\kappa(t-T_\sigma(k(t)-1))+2\theta}$ and $T_\sigma \geq \max(1, T^{\frac{6}{7}})$, we have for any $t \geq 1$:

$$\mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| \right] \leq \frac{6 \left(\sqrt{\kappa} + \sqrt{\theta} + \sqrt{D\theta} + \sqrt{D} \right)}{k(t)} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right).$$

By setting $t = KT_\sigma$ in this lemma, we get:

$$\mathbb{E} \left[\left\| \pi^{\mu, \sigma^K} - \sigma^K \right\| \right] \leq \frac{6 \left(\sqrt{\kappa} + \sqrt{\theta} + \sqrt{D\theta} + \sqrt{D} \right)}{K} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right). \quad (27)$$

Combining Eq. (25), Eq. (26), and Eq. (27), we have:

$$\begin{aligned} & \mathbb{E} [\text{GAP}(\sigma^{K+1})] \\ & \leq \frac{\mu D^2}{K+1} + \mu D \cdot \frac{6 \left(\sqrt{\kappa} + \sqrt{\theta} + \sqrt{D\theta} + \sqrt{D} \right)}{K} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right) \\ & \quad + (LD + \zeta) \cdot \sqrt{\frac{2\theta}{\kappa T_\sigma + 2\theta} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T_\sigma}{2\theta} + 1 \right) \right)} \\ & \leq \mu D^2 \frac{T_\sigma}{T} + \mu D \cdot \frac{6T_\sigma \left(\sqrt{\kappa} + \sqrt{\theta} + \sqrt{D\theta} + \sqrt{D} \right)}{T} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right) \\ & \quad + (LD + \zeta) \cdot \sqrt{\frac{2\theta}{\kappa T_\sigma} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)}, \end{aligned}$$

where the second inequality follows from $K = \frac{T}{T_\sigma}$. Finally, since $T_\sigma = c \cdot \max(1, T^{\frac{6}{7}})$, we have for any $T \geq T_\sigma$:

$$\begin{aligned} & \mathbb{E} [\text{GAP}(\sigma^{K+1})] \\ & \leq \frac{c\mu D^2}{T^{\frac{1}{7}}} + \frac{6c\mu D \left(\sqrt{\kappa} + \sqrt{\theta} + \sqrt{D\theta} + \sqrt{D} \right)}{T^{\frac{1}{7}}} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right) \\ & \quad + \frac{(LD + \zeta)}{T^{\frac{3}{7}}} \sqrt{\frac{2\theta}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} \\ & \leq \frac{6c\mu D \left(\sqrt{\kappa} + \sqrt{\theta} + \sqrt{D\theta} + \sqrt{D} + D \right)}{T^{\frac{1}{7}}} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right) \\ & \quad + \frac{(LD + \zeta)\sqrt{2\theta}}{T^{\frac{1}{7}}} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right) \\ & \leq \frac{9c(\mu D + LD + \zeta) \left(\sqrt{\kappa} + \sqrt{\theta} + \sqrt{D\theta} + \sqrt{D} + D \right)}{T^{\frac{1}{7}}} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right). \end{aligned}$$

1296 Since $T = T_\sigma K$, we have finally:
 1297

$$\begin{aligned}
 & \mathbb{E} [\text{GAP}(\pi^{T+1})] \\
 & \leq \frac{9c(\mu D + LD + \zeta) (\sqrt{\kappa} + \sqrt{\theta} + \sqrt{D\theta} + \sqrt{D} + D)}{T^{\frac{1}{7}}} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right) \\
 & = \frac{9c(D(\mu + L) + \zeta) (\sqrt{\kappa} + (\sqrt{D} + 1)(\sqrt{D} + \sqrt{\theta}))}{T^{\frac{1}{7}}} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right) \\
 & \leq \frac{18c(D(\mu + L) + \zeta) (\sqrt{\kappa} + \sqrt{(D+1)(D+\theta)})}{T^{\frac{1}{7}}} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right) \\
 & \leq \frac{26c(D(\mu + L) + \zeta) \sqrt{(D+1)(D+\theta) + \kappa}}{T^{\frac{1}{7}}} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right).
 \end{aligned}$$

□

B.2 PROOF OF LEMMA B.1

Proof of Lemma B.1. From the first-order optimality condition for π^{t+1} , we have for $t \geq 1$:

$$\langle \eta_t (\hat{V}(\pi^t) - \mu(\pi^t - \hat{\sigma}^{k(t)})) - \pi^{t+1} + \pi^t, \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle \geq 0. \quad (28)$$

Combining Eq. (28), Eq. (12), and Eq. (14), we have:

$$\begin{aligned}
 & \frac{1}{2} \|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 - \frac{1}{2} \|\pi^{\mu, \sigma^{k(t)}} - \pi^t\|^2 + \frac{1}{2} \|\pi^{t+1} - \pi^t\|^2 \\
 & \leq \eta_t \langle \hat{V}(\pi^t) - \mu(\pi^t - \hat{\sigma}^{k(t)}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle \\
 & = \eta_t \langle V(\pi^{t+1}) - \mu(\pi^{t+1} - \hat{\sigma}^{k(t)}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle + \eta_t \langle \hat{V}(\pi^t) - V(\pi^{t+1}) - \mu(\pi^t - \pi^{t+1}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle \\
 & \leq \eta_t \langle V(\pi^{\mu, \sigma^{k(t)}}) - \mu(\pi^{t+1} - \hat{\sigma}^{k(t)}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle + \eta_t \langle \hat{V}(\pi^t) - V(\pi^{t+1}) - \mu(\pi^t - \pi^{t+1}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle \\
 & = \eta_t \langle V(\pi^{\mu, \sigma^{k(t)}}) - \mu(\pi^{\mu, \sigma^{k(t)}} - \hat{\sigma}^{k(t)}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle - \eta_t \mu \|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 \\
 & \quad + \eta_t \langle V(\pi^t) - V(\pi^{t+1}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle - \eta_t \mu \langle \pi^t - \pi^{t+1}, \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle + \eta_t \langle \xi^t, \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle \\
 & \leq -\eta_t \mu \|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 + \eta_t \mu \langle \pi^{t+1} - \pi^t, \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle \\
 & \quad + \eta_t \langle V(\pi^t) - V(\pi^{t+1}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle + \eta_t \langle \xi^t, \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle \\
 & = -\eta_t \mu \|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 + \frac{\eta_t \mu}{2} \|\pi^{t+1} - \pi^t\|^2 + \frac{\eta_t \mu}{2} \|\pi^{t+1} - \pi^{\mu, \sigma^{k(t)}}\|^2 - \frac{\eta_t \mu}{2} \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 \\
 & \quad + \eta_t \langle V(\pi^t) - V(\pi^{t+1}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle + \eta_t \langle \xi^t, \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle \\
 & = -\frac{\eta_t \mu}{2} \|\pi^{t+1} - \pi^{\mu, \sigma^{k(t)}}\|^2 - \frac{\eta_t \mu}{2} \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 + \frac{\eta_t \mu}{2} \|\pi^{t+1} - \pi^t\|^2 \\
 & \quad + \eta_t \langle V(\pi^t) - V(\pi^{t+1}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle + \eta_t \langle \xi^t, \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle,
 \end{aligned} \quad (29)$$

where the third inequality follows from Eq. (1). From Cauchy-Schwarz inequality and Young's inequality, the fourth term on the right-hand side of this inequality can be bounded by:

$$\begin{aligned}
 & \langle V(\pi^t) - V(\pi^{t+1}), \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle \\
 & \leq \|V(\pi^t) - V(\pi^{t+1})\| \cdot \|\pi^{t+1} - \pi^{\mu, \sigma^{k(t)}}\|
 \end{aligned}$$

$$\begin{aligned}
&\leq L \|\pi^t - \pi^{t+1}\| \cdot \|\pi^{t+1} - \pi^{\mu, \sigma^{k(t)}}\| \\
&\leq \frac{2L^2}{\mu} \|\pi^t - \pi^{t+1}\|^2 + \frac{\mu}{8} \|\pi^{t+1} - \pi^{\mu, \sigma^{k(t)}}\|^2 \\
&\leq \frac{2L^2}{\mu} \|\pi^t - \pi^{t+1}\|^2 + \frac{\mu}{4} \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 + \frac{\mu}{4} \|\pi^{t+1} - \pi^t\|^2 \\
&= \left(\frac{4L^2}{\mu} + \frac{\mu}{2} \right) \frac{\|\pi^t - \pi^{t+1}\|^2}{2} + \frac{\mu}{2} \frac{\|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2}{2}. \tag{30}
\end{aligned}$$

By combining Eq. (29) and Eq. (30), we have:

$$\begin{aligned}
\|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 &\leq -\eta_t \mu \|\pi^{t+1} - \pi^{\mu, \sigma^{k(t)}}\|^2 + \left(1 - \frac{\eta_t \mu}{2}\right) \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 \\
&\quad - \left(1 - \eta_t \left(\frac{3\mu}{2} + \frac{4L^2}{\mu}\right)\right) \|\pi^{t+1} - \pi^t\|^2 + 2\eta_t \langle \xi^t, \pi^{t+1} - \pi^{\mu, \sigma^{k(t)}} \rangle \\
&\leq \left(1 - \frac{\eta_t \mu}{2}\right) \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 - \left(1 - \eta_t \left(\frac{3\mu}{2} + \frac{4L^2}{\mu}\right)\right) \|\pi^{t+1} - \pi^t\|^2 \\
&\quad + 2\eta_t \langle \xi^t, \pi^t - \pi^{\mu, \sigma^{k(t)}} \rangle + 2\eta_t \langle \xi^t, \pi^{t+1} - \pi^t \rangle \\
&= (1 - \eta_t \kappa) \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 - (1 - \eta_t \theta) \|\pi^{t+1} - \pi^t\|^2 \\
&\quad + 2\eta_t \langle \xi^t, \pi^t - \pi^{\mu, \sigma^{k(t)}} \rangle + 2\eta_t \langle \xi^t, \pi^{t+1} - \pi^t \rangle.
\end{aligned}$$

By taking the expectation conditioned on \mathcal{F}_t for both sides and using Assumption 4.2 (a) and (b),

$$\begin{aligned}
&\mathbb{E} \left[\|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 \mid \mathcal{F}_t \right] \\
&\leq (1 - \eta_t \kappa) \mathbb{E} \left[\|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 \mid \mathcal{F}_t \right] - (1 - \eta_t \theta) \mathbb{E} \left[\|\pi^{t+1} - \pi^t\|^2 \mid \mathcal{F}_t \right] \\
&\quad + 2\eta_t \left\langle \mathbb{E} [\xi^t \mid \mathcal{F}_t], \pi^t - \pi^{\mu, \sigma^{k(t)}} \right\rangle + 2\eta_t \mathbb{E} [\langle \xi^t, \pi^{t+1} - \pi^t \rangle \mid \mathcal{F}_t] \\
&= (1 - \eta_t \kappa) \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 - (1 - \eta_t \theta) \mathbb{E} [\|\pi^{t+1} - \pi^t\|^2 \mid \mathcal{F}_t] + 2\eta_t \mathbb{E} [\langle \xi^t, \pi^{t+1} - \pi^t \rangle \mid \mathcal{F}_t] \\
&\leq (1 - \eta_t \kappa) \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 - (1 - \eta_t \theta) \mathbb{E} [\|\pi^{t+1} - \pi^t\|^2 \mid \mathcal{F}_t] \\
&\quad + \frac{\eta_t^2}{1 - \eta_t \theta} \mathbb{E} [\|\xi^t\|^2 \mid \mathcal{F}_t] + (1 - \eta_t \theta) \mathbb{E} [\|\pi^{t+1} - \pi^t\|^2 \mid \mathcal{F}_t] \\
&\leq (1 - \eta_t \kappa) \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 + \frac{\eta_t^2}{1 - \eta_t \theta} \mathbb{E} [\|\xi^t\|^2 \mid \mathcal{F}_t] \\
&\leq (1 - \eta_t \kappa) \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 + 2\eta_t^2 \mathbb{E} [\|\xi^t\|^2 \mid \mathcal{F}_t] \\
&\leq (1 - \eta_t \kappa) \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 + 2\eta_t^2 C^2.
\end{aligned}$$

Therefore, under the setting where $\eta_t = \frac{1}{\kappa(t - T_\sigma(k(t) - 1)) + 2\theta}$, we have for any $t \geq 1$:

$$\mathbb{E} \left[\|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 \mid \mathcal{F}_t \right] \leq \left(1 - \frac{1}{t - T_\sigma(k(t) - 1) + 2\theta/\kappa}\right) \|\pi^t - \pi^{\mu, \sigma^{k(t)}}\|^2 + 2\eta_t^2 C^2.$$

Rearranging and taking the expectations, we get:

$$\begin{aligned}
&(t - T_\sigma(k(t) - 1) + 2\theta/\kappa) \mathbb{E} \left[\|\pi^{\mu, \sigma^{k(t)}} - \pi^{t+1}\|^2 \right] \\
&\leq (t - 1 - T_\sigma(k(t) - 1) + 2\theta/\kappa) \mathbb{E} \left[\|\pi^{\mu, \sigma^{k(t)}} - \pi^t\|^2 \right] + \frac{2C^2}{\kappa(t - T_\sigma(k(t) - 1) + 2\theta)}.
\end{aligned}$$

1404 Since $k(s) = k(t)$ for any $s \in [(k(t) - 1)T_\sigma + 1, T]$, telescoping the sum yields:
1405
1406
$$(t - T_\sigma(k(t) - 1) + 2\theta/\kappa) \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{t+1} \right\|^2 \right]$$

1407
1408
$$\leq (s - 1 - T_\sigma(k(t) - 1) + 2\theta/\kappa) \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \pi^s \right\|^2 \right] + \sum_{m=s}^t \frac{2C^2}{\kappa(\kappa(m - T_\sigma(k(t) - 1)) + 2\theta)}.$$

1409
1410
1411

1412 Defining $s = (k(t) - 1)T_\sigma + 1$,

$$\begin{aligned} 1413 & (t - T_\sigma(k(t) - 1) + 2\theta/\kappa) \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{t+1} \right\|^2 \right] \\ 1414 & \leq \frac{2\theta}{\kappa} \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{(k(t)-1)T_\sigma+1} \right\|^2 \right] + \frac{2C^2}{\kappa} \sum_{m=(k(t)-1)T_\sigma+1}^t \frac{1}{\kappa(m - T_\sigma(k(t) - 1)) + 2\theta}. \end{aligned}$$

1415
1416 Therefore,

$$\begin{aligned} 1417 & \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{t+1} \right\|^2 \right] \leq \frac{2\theta}{\kappa(t - T_\sigma(k(t) - 1)) + 2\theta} \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{(k(t)-1)T_\sigma+1} \right\|^2 \right] \\ 1418 & + \frac{2C^2}{\kappa(t - T_\sigma(k(t) - 1)) + 2\theta} \sum_{m=1}^{t-(k(t)-1)T_\sigma} \frac{1}{\kappa m + 2\theta}. \end{aligned} \quad (31)$$

1419
1420 Here, we have:

$$\begin{aligned} 1421 & \sum_{m=1}^{t-(k(t)-1)T_\sigma} \frac{1}{\kappa m + 2\theta} \leq \int_0^{t-(k(t)-1)T_\sigma} \frac{1}{\kappa x + 2\theta} dx = \frac{1}{\kappa} \ln \left(\frac{\kappa(t - (k(t) - 1)T_\sigma)}{2\theta} + 1 \right). \end{aligned} \quad (32)$$

1422 Combining Eq. (31), Eq. (32), and the fact that $\pi^{(k(t)-1)T_\sigma+1} = \sigma^{k(t)}$, we have:

$$\begin{aligned} 1423 & \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{t+1} \right\|^2 \right] \\ 1424 & \leq \frac{2\theta}{\kappa(t - (k(t) - 1)T_\sigma) + 2\theta} \left(\mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 \right] + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa(t - (k(t) - 1)T_\sigma)}{2\theta} + 1 \right) \right) \\ 1425 & \leq \frac{2\theta}{\kappa(t - (k(t) - 1)T_\sigma) + 2\theta} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa(t - (k(t) - 1)T_\sigma)}{2\theta} + 1 \right) \right). \end{aligned}$$

1426
1427
1428 \square

1429 B.3 PROOF OF LEMMA B.2

1430 Proof of Lemma B.2. First, from Lemma B.1, we have for any $k \geq 1$:

$$\mathbb{E} \left[\left\| \pi^{\mu, \sigma^k} - \sigma^{k+1} \right\|^2 \right] \leq \frac{2\theta}{\kappa T_\sigma + 2\theta} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T_\sigma}{2\theta} + 1 \right) \right).$$

1431 Moreover, by taking the expectation of Eq. (21), we have for any $t \geq 1$ such that $k(t) \geq 2$:

$$\begin{aligned} 1432 & \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 \right] \leq \frac{2D}{k(t) + 1} \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| \right] + \frac{12D^2}{(k(t) + 1)^2} \\ 1433 & + 8\mathbb{E} \left[\left\| \sigma^{k(t)+1} - \pi^{\mu, \sigma^{k(t)}} \right\|^2 \right] + 8D \sum_{l=1}^{k(t)} \mathbb{E} \left[\left\| \pi^{\mu, \sigma^l} - \sigma^{l+1} \right\| \right]. \end{aligned}$$

Combining these inequalities, we get for any $t \geq 1$ such that $k(t) \geq 2$:

$$\begin{aligned} \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 \right] &\leq \frac{2D}{k(t)+1} \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| \right] + \frac{12D^2}{(k(t)+1)^2} \\ &+ \frac{16\theta}{\kappa T_\sigma} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T_\sigma}{2\theta} + 1 \right) \right) + 8Dk(t) \sqrt{\frac{2\theta}{\kappa T_\sigma} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T_\sigma}{2\theta} + 1 \right) \right)}. \end{aligned}$$

Since $T_\sigma \geq \max(1, T^{\frac{6}{7}}) \Rightarrow \frac{k(t)^3}{\sqrt{T_\sigma}} \leq 1$, we have:

$$\begin{aligned} \mathbb{E} \left[\left(\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| - \frac{D}{k(t)+1} \right)^2 \right] &\leq \frac{13D^2}{k(t)^2} + \frac{16\theta}{\kappa k(t)^2} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right) \\ &+ \frac{8D}{k(t)^2} \sqrt{\frac{2\theta}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)}. \end{aligned}$$

Since $\mathbb{E}[X]^2 \leq \mathbb{E}[X^2]$ for any random variable X , we get:

$$\begin{aligned} \frac{13D^2}{k(t)^2} + \frac{16\theta}{\kappa k(t)^2} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right) + \frac{8D}{k(t)^2} \sqrt{\frac{2\theta}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} \\ \geq \mathbb{E} \left[\left(\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| - \frac{D}{k(t)+1} \right)^2 \right] \\ \geq \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| - \frac{D}{k(t)+1} \right]^2 \\ = \left(\mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| \right] - \frac{D}{k(t)+1} \right)^2. \end{aligned}$$

Then, we have:

$$\begin{aligned} \mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| \right] &\leq \frac{D}{k(t)} + \frac{4D}{k(t)} + \frac{4\sqrt{\theta}}{\sqrt{\kappa}k(t)} \sqrt{D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right)} + \frac{3\sqrt{D}}{k(t)} \left(\frac{2\theta}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right) \right)^{\frac{1}{4}} \\ &\leq \frac{5(\sqrt{\kappa} + \sqrt{\theta})}{k(t)\sqrt{\kappa}} \sqrt{D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right)} + \frac{6\sqrt{D}(\sqrt{\theta} + 1)}{k(t)} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right) \\ &\leq \frac{6(\sqrt{\kappa} + \sqrt{\theta} + \sqrt{D\theta} + \sqrt{D})}{k(t)} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right). \end{aligned}$$

Furthermore, for $k(t) = 1$, we have:

$$\mathbb{E} \left[\left\| \pi^{\mu, \sigma^1} - \sigma^1 \right\| \right] \leq D \leq \frac{6(\sqrt{\kappa} + \sqrt{\theta} + \sqrt{D\theta} + \sqrt{D})}{1} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right).$$

Therefore, we have for any $t \geq 1$:

$$\mathbb{E} \left[\left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\| \right] \leq \frac{6(\sqrt{\kappa} + \sqrt{\theta} + \sqrt{D\theta} + \sqrt{D})}{k(t)} \left(\sqrt{\frac{1}{\kappa} \left(D^2 + \frac{C^2}{\kappa\theta} \ln \left(\frac{\kappa T}{2\theta} + 1 \right) \right)} + 1 \right).$$

□

1512 **C PROOF OF THEOREM 5.1**
 1513

1514 *Proof of Theorem 5.1.* By the definition of dynamic regret, we have:
 1515

$$\begin{aligned} \text{DynamicReg}_i(T) &= \sum_{t=1}^T \left(\max_{x \in \mathcal{X}_i} v_i(x, \pi_{-i}^t) - v_i(\pi^t) \right) \\ &\leq \mathcal{O}(1) + \sum_{t=2}^T \sum_{i=1}^N \left(\max_{x \in \mathcal{X}_i} v_i(x, \pi_{-i}^t) - v_i(\pi^t) \right). \end{aligned}$$

1522 Here, we introduce the following lemma:
 1523

1525 **Lemma C.1** (Lemma 2 of Cai et al. (2022a)). *For any $\pi \in \mathcal{X}$, we have:*
 1526

$$\sum_{i=1}^N \left(\max_{\tilde{\pi}_i \in \mathcal{X}_i} v_i(\tilde{\pi}_i, \pi_{-i}) - v_i(\pi) \right) \leq \text{GAP}(\pi) \leq D \cdot \max_{\tilde{\pi} \in \mathcal{X}} \langle V(\pi), \tilde{\pi} - \pi \rangle.$$

1530 Therefore, we have:
 1531

$$\text{DynamicReg}_i(T) \leq \mathcal{O}(1) + \sum_{t=2}^T \text{GAP}(\pi^t).$$

1535 Thus, from Theorem 4.1:
 1536

$$\begin{aligned} \text{DynamicReg}_i(T) &\leq \mathcal{O}(1) + \sum_{t=2}^T \mathcal{O}\left(\frac{\ln T}{t}\right) \\ &\leq \mathcal{O}((\ln T)^2). \end{aligned}$$

1542 \square
 1543

1544 **D EXPERIMENTAL DETAILS**
 1545

1547 **D.1 INFORMATION ON THE COMPUTER RESOURCES**
 1548

1549 The experiments were conducted on macOS Sonoma 14.4.1 with Apple M2 Max and 32GB RAM.
 1550

1551 **D.2 HARD CONCAVE-CONVEX GAME**
 1552

1553 Following the setup in Ouyang & Xu (2021); Cai & Zheng (2023), we choose
 1554

$$A = \frac{1}{4} \begin{bmatrix} & & -1 & 1 \\ & \cdots & \cdots & \\ -1 & 1 & 1 & \\ 1 & & & \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad b = \frac{1}{4} \begin{bmatrix} 1 \\ 1 \\ \cdots \\ 1 \end{bmatrix} \in \mathbb{R}^n, \quad h = \frac{1}{4} \begin{bmatrix} 0 \\ 0 \\ \cdots \\ 0 \\ 1 \end{bmatrix} \in \mathbb{R}^n,$$

1560 and $H = 2A^\top A$.
 1561

1562 **D.3 HYPERPARAMETERS**
 1563

1564 For each game, we carefully tuned the hyperparameters for each algorithm to ensure optimal performance. The specific parameters for each game and setting are summarized in Table 1.
 1565

| Game | Algorithm | η | T_σ | μ |
|--------------------------------------|-----------|--------|------------|-------|
| Random Payoff (Full Feedback) | OG | 0.05 | - | - |
| | AOG | 0.05 | - | - |
| | APGA | 0.05 | 20 | 1.0 |
| | GABP | 0.05 | 10 | 1.0 |
| Random Payoff (Noisy Feedback) | OG | 0.001 | - | - |
| | AOG | 0.001 | - | - |
| | APGA | 0.001 | 2000 | 1.0 |
| | GABP | 0.001 | 1000 | 1.0 |
| Hard Concave-Convex (Full Feedback) | OG | 1.0 | - | - |
| | AOG | 1.0 | - | - |
| | APGA | 1.0 | 20 | 0.1 |
| | GABP | 1.0 | 20 | 0.1 |
| Hard Concave-Convex (Noisy Feedback) | OG | 0.5 | - | - |
| | AOG | 0.5 | - | - |
| | APGA | 0.5 | 50 | 0.1 |
| | GABP | 0.1 | 100 | 0.1 |

Table 1: Hyperparameters

E COMPARISON WITH EXISTING LEARNING ALGORITHMS

E.1 RELATIONSHIP WITH ACCELERATED OPTIMISTIC GRADIENT ALGORITHM

Our GABP bears some relation to Accelerated Optimistic Gradient (AOG) ([Cai & Zheng, 2023](#)), which updates the strategy by:

$$\begin{aligned}\pi_i^{t+\frac{1}{2}} &= \arg \max_{x \in \mathcal{X}_i} \left\{ \left\langle \eta \widehat{\nabla}_{\pi_i} v_i(\pi^{t-\frac{1}{2}}) + \frac{\pi_i^1 - \pi_i^t}{t+1}, x \right\rangle - \frac{1}{2} \|x - \pi_i^t\|^2 \right\}, \\ \pi_i^{t+1} &= \arg \max_{x \in \mathcal{X}_i} \left\{ \left\langle \eta \widehat{\nabla}_{\pi_i} v_i(\pi^{t+\frac{1}{2}}) + \frac{\pi_i^1 - \pi_i^t}{t+1}, x \right\rangle - \frac{1}{2} \|x - \pi_i^t\|^2 \right\}.\end{aligned}$$

This can be equivalently written as:

$$\begin{aligned}\pi_i^{t+\frac{1}{2}} &= \arg \max_{x \in \mathcal{X}_i} \left\{ \eta \left\langle \widehat{\nabla}_{\pi_i} v_i(\pi^{t-\frac{1}{2}}), x \right\rangle - \frac{1}{2} \left\| x - \frac{t\pi_i^t + \pi_i^1}{t+1} \right\|^2 \right\}, \\ \pi_i^{t+1} &= \arg \max_{x \in \mathcal{X}_i} \left\{ \eta \left\langle \widehat{\nabla}_{\pi_i} v_i(\pi^{t+\frac{1}{2}}), x \right\rangle - \frac{1}{2} \left\| x - \frac{t\pi_i^t + \pi_i^1}{t+1} \right\|^2 \right\}.\end{aligned}$$

This means that AOG employs a convex combination $\frac{t\pi_i^t + \pi_i^1}{t+1}$ of the current strategy π_i^t and initial strategy π_i^1 as the proximal point in gradient ascent. However, our GABP diverges from AOG in that it uses a convex combination $\frac{k(t)\sigma_i^{k(t)} + \sigma_i^1}{k(t)+1}$ of $\sigma_i^{k(t)}$ and σ_i^1 as the reference strategy for the perturbation term.

In terms of proof for a last-iterate convergence result, [Cai & Zheng \(2023\)](#) have employed the following potential function for the full feedback setting:

$$P^t := \frac{t(t+1)}{2} \left(\eta^2 \left\| -\widehat{V}(\pi^t) + c^t \right\|^2 + \eta^2 \left\| \widehat{V}(\pi^t) - \widehat{V}(\pi^{t-\frac{1}{2}}) \right\|^2 \right) + t\eta \langle -\widehat{V}(\pi^t) + c^t, \pi^t - \pi^1 \rangle,$$

where $c^t = \frac{\pi^{t-1} + \eta \widehat{V}(\pi^{t-\frac{1}{2}}) + \frac{1}{t}(\pi^1 - \pi^{t-1}) - \pi^t}{\eta}$ and $\widehat{V}(\cdot) = (\widehat{\nabla}_{\pi_i} v_i(\cdot))_{i \in [N]}$. Compared to our potential function $P^{k(t)}$, their potential function includes the term of $\eta^2 \left\| \widehat{V}(\pi^t) - \widehat{V}(\pi^{t-\frac{1}{2}}) \right\|^2$. In the noisy feedback setting, the value of this term could remain high even if $\pi^t = \pi^{t-\frac{1}{2}}$. Therefore, this term complicates providing a last-iterate convergence result for the noisy feedback setting. In contrast,

1620 our potential function $P^{k(t)}$ does not contain the term depending on $\hat{\nabla}v_i(\pi^t)$:
 1621

$$\begin{aligned} 1622 \quad P^{k(t)} &:= \frac{k(t)(k(t)+1)}{2} \left\| \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\|^2 \\ 1623 \\ 1624 \quad &+ k(t)(k(t)+1) \left\langle \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\rangle. \\ 1625 \end{aligned}$$

1626 This allows us to provide the last-iterate convergence rates both for full and noisy feedback settings.
 1627

1628 E.2 COMPARISON OF LAST-ITERATE CONVERGENCE RESULTS

1630 In this section, we provide Table 2 for comparison with last-iterate convergence results of existing
 1631 representative learning algorithms in monotone games.

| Algorithm | Full Feedback | Noisy Feedback |
|---|-----------------------------------|---|
| Extragradient (Cai et al., 2022a;b) | $\mathcal{O}(1/\sqrt{T})$ | N/A |
| Optimistic Gradient (Golowich et al., 2020b; Gorbunov et al., 2022; Cai et al., 2022a) | $\mathcal{O}(1/\sqrt{T})$ | N/A |
| Extra Anchored Gradient (Yoon & Ryu, 2021) | $\mathcal{O}(1/T)$ | N/A |
| Accelerated Optimistic Gradient (Cai & Zheng, 2023) | $\mathcal{O}(1/T)$ | N/A |
| Iterative Tikhonov Regularization (Koshal et al., 2010; Tatarenko & Kamgarpour, 2019) | N/A | Asymptotic* |
| Adaptively Perturbed Gradient Ascent (Abe et al., 2024) | $\tilde{\mathcal{O}}(1/\sqrt{T})$ | $\tilde{\mathcal{O}}(1/T^{\frac{1}{10}})$ |
| Ours | $\tilde{\mathcal{O}}(1/T)$ | $\tilde{\mathcal{O}}(1/T^{\frac{1}{7}})$ |

1647 Table 2: Existing results on last-iterate convergence in monotone games. (*) means the result holds
 1648 only under bandit feedback.
 1649

1652 F ADDITIONAL EXPERIMENTAL RESULTS

1655 In this section, we experimentally compare the gap function for GABP with APGA, OG, and AOG
 1656 in the full/noisy feedback setting. The experimental settings are identical to those in Section 6.
 1657 Figure 3 illustrates the logarithm of the gap function for π^t . We can observe the same trends as in
 1658 Section 6.
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

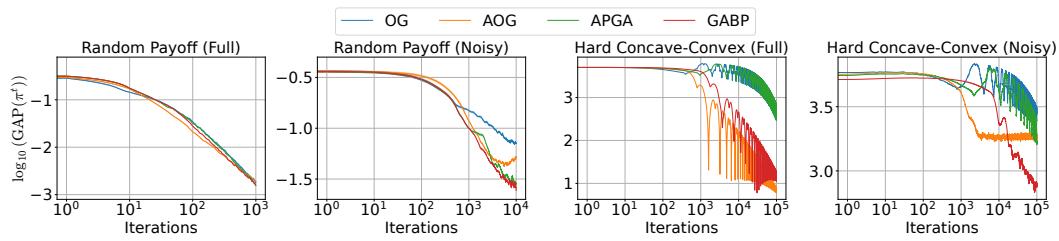


Figure 3: The gap function for π^t for GABP, APGA, OG, and AOG with full and noisy feedback in the random payoff and hard concave-convex games, respectively. The shaded area represents the standard errors.

G NOTATIONS

In this section, we summarize the notations we use in Table 3.

Table 3: Notations

| Symbol | Description |
|---------------------------------|---|
| N | Number of players |
| \mathcal{X}_i | Strategy space for player i |
| \mathcal{X} | Joint strategy space: $\mathcal{X} = \prod_{i=1}^N \mathcal{X}_i$ |
| v_i | Payoff function for player i |
| π_i | Strategy for player i |
| π | Strategy profile: $\pi = (\pi_i)_{i \in [N]}$ |
| π^* | Nash equilibrium |
| Π^* | Set of Nash equilibria |
| $\text{GAP}(\pi)$ | Gap function of π : $\text{GAP}(\pi) = \max_{\tilde{\pi} \in \mathcal{X}} \sum_{i=1}^N \langle \nabla_{\pi_i} v_i(\pi), \tilde{\pi}_i - \pi_i \rangle$ |
| $r^{\tan}(\pi)$ | Tangent residual of π : $r^{\tan}(\pi) = \min_{a \in N_{\mathcal{X}}(\pi)} \ -V(\pi) + a \ $ |
| $\nabla_{\pi_i} v_i(\pi)$ | Gradient vector of v_i with respect to π_i |
| $\hat{\nabla}_{\pi_i} v_i(\pi)$ | Noisy gradient vector of v_i with respect to π_i : $\hat{\nabla}_{\pi_i} v_i(\pi) = \nabla_{\pi_i} v_i(\pi) + \xi_i^t$ |
| ξ_i^t | Noise vector for player i at iteration t |
| $V(\cdot)$ | Gradient operator of the payoff functions: $V(\cdot) = (\nabla_{\pi_i} v_i(\cdot))_{i \in [N]}$ |
| T | Total number of iterations |
| η_t | Learning rate at iteration t |
| μ | Perturbation strength |
| T_σ | Update interval for the anchoring strategy |
| π^t | Strategy profile at iteration t |
| $k(t)$ | Number of updates of the anchoring strategy up to iteration t |
| K | Total number of the updates of the anchoring strategy |
| $\sigma^{k(t)}$ | Anchoring strategy profile at iteration t |
| $\hat{\sigma}^{k(t)}$ | convex combination of $\sigma^{k(t)}$ and σ^1 : $\hat{\sigma}^{k(t)} = \frac{k(t)\sigma^{k(t)} + \sigma^1}{k(t) + 1}$ |
| $\pi^{\mu, \sigma^{k(t)}}$ | Stationary point satisfies: $\forall i \in [N], \pi_i^{\mu, \sigma^{k(t)}} = \arg \max_{x \in \mathcal{X}_i} \left\{ v_i(x, \pi_{-i}^{\mu, \sigma^{k(t)}}) - \frac{\mu}{2} \ x - \hat{\sigma}^{k(t)}\ ^2 \right\}$ |
| L | Smoothness parameter of $(v_i)_{i \in [N]}$ |