
Boosting Perturbed Gradient Ascent for Last-Iterate Convergence in Games

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper introduces a payoff perturbation technique, introducing a strong convexity to players' payoff functions in games. This technique is specifically designed 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 containing additive noise. Although perturbation is known to facilitate the convergence of learning algorithms, the magnitude of perturbation requires careful adjustment to ensure last-iterate convergence. Previous studies have proposed a scheme in which the magnitude is determined by the distance from an anchoring or reference strategy, which is periodically re-initialized. In response, this paper proposes 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 algorithms, even in the presence of additive noise.

15 1 Introduction

16 This study considers online learning in monotone games, where the gradient of the payoff function is monotone in the strategy profile space. Monotone games encompassed diverse well-studied games as special instances, such as concave-convex games, zero-sum polymatrix games [Cai and Daskalakis, 2011, Cai et al., 2016], λ -cocoercive games [Lin et al., 2020], and Cournot competition [Bravo et al., 2018]. Due to their wide-ranging applications, there has been growing interest in developing learning algorithms to compute Nash equilibria in monotone games.

22 Typical learning algorithms such as Gradient Ascent [Zinkevich, 2003] and Multiplicative Weights Update [Bailey and Piliouras, 2018] have been extensively studied and shown to converge to equilibria in an average-iterate sense, which is termed *average-iterate convergence*. However, averaging the strategies can be undesirable because it can lead to additional memory or computational costs in the context of training Generative Adversarial Networks [Goodfellow et al., 2014] and preference-based fine-tuning of large language models [Munos et al., 2023, Swamy et al., 2024]. In contrast, *last-iterate convergence*, in which the updated strategy profile itself converges to a Nash equilibrium, has emerged as a stronger notion than average-iterate convergence.

30 Payoff-perturbed algorithms have recently been regaining attention in this context [Sokota et al., 2023, Liu et al., 2023]. Payoff perturbation is a classical technique, e.g., [Facchinei and Pang, 2003] and introduces a strongly convex penalty to the players' payoff functions to stabilize learning, which leads to convergence to approximate equilibria, not only in the *full feedback* setting where the perfect gradient vector of the payoff function can be used to update strategies, but also in the *noisy feedback* setting where the gradient vector is contaminated by noise.

36 However, to ensure convergence toward a Nash equilibrium of the underlying game, the magnitude
 37 of perturbation requires careful adjustment. As a remedy, it is adjusted by the distance from an
 38 anchoring or reference strategy. Koshal et al. [2010] and Tatarenko and Kamgarpour [2019] simply
 39 decay the magnitude in each iteration, and their methods asymptotically converge, since the perturbed
 40 function gradually loses strong convexity. In response to this, recent studies [Perolat et al., 2021, Abe
 41 et al., 2023, 2024] re-initialize the anchoring strategies periodically, or in a predefined interval, so
 42 that they keep the perturbed function strongly convex and achieve non-asymptotic convergence.

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

49 This paper, in particular, focuses on *Adaptively Perturbed Mirror Descent* (APMD) [Abe et al., 2024],
 50 which achieves $\tilde{\mathcal{O}}(1/\sqrt{T})$ ¹ and $\tilde{\mathcal{O}}(1/T^{\frac{1}{10}})$ last-iterate convergence rates in the full/noisy feedback
 51 setting, respectively. The motivation of this study lies in improving the convergence rates of APMD.
 52 We propose an elegant one-line modification of APMD, which effectively accelerates convergence.
 53 In fact, we just add the difference between the current anchoring strategy and the initial anchoring
 54 strategy to the payoff perturbation function in APMD.

55 Our contributions are manifold. Firstly, we propose a novel payoff-perturbed learning algorithm
 56 named *Gradient Ascent with Boosting Payoff Perturbation* (GABP). This method incorporates a
 57 unique perturbation payoff function, enabling it to achieve faster convergence rates than APMD. Sub-
 58 sequently, we prove that GABP exhibits accelerated $\tilde{\mathcal{O}}(1/T)$ and $\tilde{\mathcal{O}}(1/T^{\frac{1}{7}})$ last-iterate convergence
 59 rates to a Nash equilibrium with full and noisy feedback, respectively. We further show that each
 60 player’s individual regret is at most $\mathcal{O}((\ln T)^2)$ in the full feedback setting, provided all players play
 61 according to GABP. Finally, through our experiments, we demonstrate the competitive or superior
 62 performance of GABP over Optimistic Gradient Ascent [Daskalakis et al., 2018, Wei et al., 2021]
 63 and APMD in concave-convex games, irrespective of the presence of noise.

64 2 Preliminaries

65 **Monotone games.** In this study, we focus on a continuous multi-player game, which is denoted
 66 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
 67 $i \in [N]$ chooses a *strategy* π_i from a d_i -dimensional compact convex strategy space \mathcal{X}_i , and we
 68 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
 69 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 ,
 70 and $\pi = (\pi_i)_{i \in [N]} \in \mathcal{X}$ as the *strategy profile*. This paper particularly studies learning in *smooth*
 71 *monotone games*, where the gradient operator $V(\cdot) = (\nabla_{\pi_i} v_i(\cdot))_{i \in [N]}$ of the payoff functions is
 72 monotone: $\forall \pi, \pi' \in \mathcal{X}$,

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

73 and L -Lipschitz for $L > 0$

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

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

75 Many common and well-studied games, such as concave-convex games, zero-sum polymatrix games
 76 [Cai et al., 2016], λ -cocoercive games [Lin et al., 2020], and Cournot competition [Bravo et al.,
 77 2018], are included in the class of monotone games.

78 **Example 2.1 (Concave-Convex Games).** Consider a game defined by $(\{1, 2\}, (\mathcal{X}_1, \mathcal{X}_2), (v, -v))$,
 79 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
 80 minimize v . If v is concave in $x_1 \in \mathcal{X}_1$ and convex in $x_2 \in \mathcal{X}_2$, the game is called a concave-convex
 81 game or minimax optimization problem, and it is not hard to see that this game is a special case of
 82 monotone games.

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

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

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

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

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

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

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

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

93 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 of
 94 $\pi \in \mathcal{X}$. It is easy to see that $\text{GAP}(\pi) \geq 0$ (resp. $r^{\text{tan}}(\pi) \geq 0$) for any $\pi \in \mathcal{X}$, and the equality holds
 95 if and only if π is a Nash equilibrium. Defining $D := \sup_{\pi, \pi' \in \mathcal{X}} \|\pi - \pi'\|$ as the diameter of \mathcal{X} , the
 96 gap function for any given strategy profile $\pi \in \mathcal{X}$ is upper bounded by its tangent residual.

97 **Lemma 2.2** (Lemma 2 of Cai et al. [2022a]). *For any $\pi \in \mathcal{X}$, we have:*

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

98 The gap function and the tangent residual are standard measures of proximity to Nash equilibrium;
 99 e.g., it has been used in Cai and Zheng [2023], Abe et al. [2024].

100 **Problem setting.** This study focuses on the online learning setting in which the following process
 101 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
 102 previously observed feedback; (ii) Each player i receives the (noisy) gradient vector $\widehat{\nabla}_{\pi_i} v_i(\pi^t)$ as
 103 feedback. This study examines two feedback models: *full feedback* and *noisy feedback*. In the full
 104 feedback setting, each player observes the perfect gradient vector $\widehat{\nabla}_{\pi_i} v_i(\pi^t) = \nabla_{\pi_i} v_i(\pi^t)$. In the
 105 noisy feedback setting, each player's gradient feedback $\widehat{\nabla}_{\pi_i} v_i(\pi^t)$ is contaminated by an additive
 106 noise vector ξ_i^t , i.e., $\widehat{\nabla}_{\pi_i} v_i(\pi^t) = \nabla_{\pi_i} v_i(\pi^t) + \xi_i^t$, where $\xi_i^t \in \mathbb{R}^{d_i}$. Throughout the paper, we
 107 assume that ξ_i^t is the zero-mean and bounded-variance noise vector at each iteration t .

108 **Adaptively perturbed Mirror Descent.** To facilitate the convergence in the online learning setting,
 109 recent studies have utilized a *payoff perturbation* technique, where payoff functions are perturbed by
 110 strongly convex functions [Sokota et al., 2023, Liu et al., 2023, Abe et al., 2022]. However, while
 111 the addition of these strongly convex functions leads learning algorithms to converge to a stationary
 112 point, this stationary point may be significantly distant from a Nash equilibrium. Therefore, the
 113 magnitude of perturbation requires careful adjustment. Perolat et al. [2021], Abe et al. [2023, 2024]
 114 have introduced a scheme in which the magnitude is determined by the distance (or divergence
 115 function) from an anchoring strategy σ_i , which is periodically re-initialized. Specifically, Adaptively
 116 Perturbed Mirror Descent (APMD) [Abe et al., 2024] perturbs each player's payoff function by a
 117 strongly convex divergence function $G(\pi_i, \sigma_i) : \mathcal{X}_i \times \mathcal{X}_i \rightarrow [0, \infty)$, where the anchoring strategy σ_i
 118 is periodically replaced by the current strategy π_i^t every predefined iterations T_σ .

119 Let us define $\sigma_i^{k(t)}$ as the anchoring strategy after $k(t)$ updates. Since σ_i is overwritten every T_σ
 120 iterations, we can write $k(t) = \lfloor (t-1)/T_\sigma \rfloor + 1$ and $\sigma_i^{k(t)} = \pi_i^{T_\sigma(k(t)-1)+1}$. Except for the payoff
 121 perturbation and the update of the anchor strategy, APMD updates each player i 's strategy in the
 122 same way as standard Mirror Descent algorithms:

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

Algorithm 1 GABP for player i .

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

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

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

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

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

5: Update the strategy by

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

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

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

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

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

10: **end if**

11: **end for**

123 where η_t is the learning rate at iteration t , $\mu \in (0, \infty)$ is the *perturbation strength*, and $D_\psi(\pi_i, \pi_i') =$
124 $\psi(\pi_i) - \psi(\pi_i') - \langle \nabla \psi(\pi_i'), \pi_i - \pi_i' \rangle$ as the Bregman divergence associated with a strictly convex
125 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
126 equivalently written as:

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

127 We refer to this version of APMD as Adaptively Perturbed Gradient Ascent (APGA). [Abe et al.](#)
128 [\[2024\]](#) have shown that APGA exhibits the convergence rates of $\tilde{O}(1/\sqrt{T})$ and $\tilde{O}(1/T^{1/10})$ with full
129 and noisy feedback, respectively.

130 3 Gradient ascent with boosting payoff perturbation

131 This section proposes an accelerated version of APGA, Gradient Ascent with Boosting Payoff
132 Perturbation (GABP). The pseudo-code of GABP is outlined in Algorithm 1. In order to obtain faster
133 last-iterate convergence rates compared to APGA, GABP introduces a novel payoff perturbation term
134 in addition to APGA's original payoff perturbation term, $\mu (\pi_i^t - \sigma_i^{k(t)})$. Formally, GABP updates
135 each player's strategy as follows:

$$\pi_i^{t+1} = \arg \max_{x \in \mathcal{X}_i} \left\{ \eta_t \left\langle \widehat{\nabla}_{\pi_i} v_i(\pi^t) - \underbrace{\mu \frac{\sigma_i^{k(t)} - \sigma_i^1}{k(t) + 1}}_{(*)} - \mu (\pi_i^t - \sigma_i^{k(t)}), x \right\rangle - \frac{1}{2} \|x - \pi_i^t\|^2 \right\}. \quad (4)$$

136 The term $(*)$ is our proposed additional perturbation term. It shrinks as $k(t)$, the number of updates
137 of $\sigma_i^{k(t)}$, increases.

138 For a more intuitive explanation of the proposed perturbation term, we present the following update
139 rule, which is equivalent to (4):

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

140 From this formula, it appears that GABP replaces the reference strategy $\sigma_i^{k(t)}$ for the perturbation
141 term in (3) of APGA with $\frac{k(t)\sigma_i^{k(t)} + \sigma_i^1}{k(t) + 1}$. As a result, the anchoring strategy in GABP evolves more
142 gradually than in APGA, leading to further stabilization of the learning dynamics. There is a tradeoff

143 between the shrinking speed of the term (*) and the stabilizing impact on the last-iterate convergence
 144 rate of GABP. The shrinking speed of $1/(k(t) + 1)$ achieves a faster convergence rate, and we believe
 145 that this represents the optimal balance for this trade-off. Although one might think that the term (*)
 146 is closely related to Accelerated Optimistic Gradient (AOG) [Cai and Zheng, 2023], we discuss the
 147 detail in Appendix F to be concise and avoid a complicated explanation.

148 4 Last-iterate convergence rates

149 This section provides the last-iterate convergence rates of GABP in the full/noisy feedback setting,
 150 respectively.

151 4.1 Full feedback setting

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

156 **Theorem 4.1.** *If we use the constant learning rate $\eta_t = \eta \in (0, \frac{\mu}{(L+\mu)^2})$ and the constant perturba-*
 157 *tion 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*
 158 *π^t updated by GABP satisfies for any $t \in \{2, 3, \dots, T+1\}$:*

$$\begin{aligned} \text{GAP}(\pi^t) &\leq \frac{17cD^2 \left(\frac{6 \ln 3(T+1)}{\ln(1+\eta\mu)} + 1 \right)}{t-1} \left(\mu + \frac{1+\eta L}{\eta} \right), \text{ and} \\ r^{\tan}(\pi^t) &\leq \frac{17cD \left(\frac{6 \ln 3(T+1)}{\ln(1+\eta\mu)} + 1 \right)}{t-1} \left(\mu + \frac{1+\eta L}{\eta} \right). \end{aligned}$$

159 This rate is significantly faster than APGA's rate of $\tilde{\mathcal{O}}(1/\sqrt{T})$. Moreover, it is a competitive rate
 160 compared to the previous state-of-the-art rate of $\mathcal{O}(1/T)$ [Yoon and Ryu, 2021, Cai and Zheng, 2023].
 161 Note that the rate in Theorem 4.1 holds for any constant perturbation strength $\mu > 0$.

162 4.1.1 Proof sketch of Theorem 4.1

163 To derive the bound of the gap function $\text{GAP}(\pi^t)$, it is sufficient to derive that of $r^{\tan}(\pi^t)$ due to
 164 Lemma 2.2. This section provides the proof sketch of Theorem 4.1. The complete proof is placed in
 165 Appendix B.

166 **(1) Decomposition of the tangent residual of the last-iterate strategy profile.** From the first-
 167 order optimality condition for π^t , we can see that $V(\pi^{t-1}) - \mu \left(\pi^{t-1} - \frac{k(t-1)\sigma^{k(t-1)} + \sigma^1}{k(t-1)+1} \right) -$
 168 $\frac{1}{\eta} (\pi^t - \pi^{t-1}) \in N_{\mathcal{X}}(\pi^t)$. Therefore, from the triangle inequality and L -smoothness (2) of the
 169 gradient operator, the tangent residual $r^{\tan}(\pi^t)$ can be bounded as:

$$\begin{aligned} r^{\tan}(\pi^t) &= \min_{a \in N_{\mathcal{X}}(\pi^t)} \left\| -V(\pi^t) + a \right\| \\ &\leq \mathcal{O}(\|\pi^t - \pi^{t-1}\|) + \mathcal{O}\left(\left\| \pi^{t-1} - \frac{k(t-1)\sigma^{k(t-1)} + \sigma^1}{k(t-1)+1} \right\|\right) + \mathcal{O}\left(\frac{1}{k(t-1)+1}\right). \end{aligned}$$

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

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

171 where $\hat{\sigma}_i^{k(t)} = \frac{k(t)\sigma_i^{k(t)} + \sigma_i^1}{k(t)+1}$. We will show that π^t converges to the stationary point $\pi^{\mu, \sigma^{k(t)}}$ at an
 172 exponential rate later. By using $\pi^{\mu, \sigma^{k(t)}}$ and applying the triangle inequality to $\|\pi^t - \pi^{t-1}\|$, we de-
 173 compose the term of $\mathcal{O}(\|\pi^t - \pi^{t-1}\|)$ into $\mathcal{O}(\|\pi^t - \pi^{\mu, \sigma^{k(t-1)}}\|)$ and $\mathcal{O}(\|\pi^{\mu, \sigma^{k(t-1)}} - \pi^{t-1}\|)$.

174 Similarly, the term of $\mathcal{O}(\|\pi^{t-1} - \sigma^{k(t-1)}\|)$ is decomposed into $\mathcal{O}(\|\pi^{t-1} - \pi^{\mu, \sigma^{k(t-1)-1}}\|)$ and
 175 $\mathcal{O}(\|\pi^{\mu, \sigma^{k(t-1)-1}} - \sigma^{k(t-1)}\|)$. Then, the tangent residual is bounded as follows:

$$\begin{aligned} r^{\tan}(\pi^t) &\leq \mathcal{O}\left(\left\|\pi^{\mu, \sigma^{k(t-1)}} - \pi^t\right\|\right) + \mathcal{O}\left(\left\|\pi^{\mu, \sigma^{k(t-1)}} - \pi^{t-1}\right\|\right) \\ &\quad + \mathcal{O}\left(\left\|\pi^{\mu, \sigma^{k(t-1)}} - \sigma^{k(t-1)}\right\|\right) + \mathcal{O}\left(\frac{1}{k(t-1)+1}\right). \end{aligned} \quad (5)$$

176 Therefore, it is enough to derive the convergence rate on $\|\pi^{\mu, \sigma^{k(t)-1}} - \pi^t\|$ and $\|\pi^{\mu, \sigma^{k(t-1)}} - \sigma^{k(t-1)}\|$.

177 **(2) Convergence rate of π^t to the stationary point $\pi^{\mu, \sigma^{k(t)}}$.** Using the strong convexity of the
 178 perturbation payoff function, $\frac{\mu}{2}\|x - \hat{\sigma}_i^{k(t)}\|^2$, we show that π^t converges to $\pi^{\mu, \sigma^{k(t)}}$ exponentially
 179 fast (in Lemma B.1). That is, we have for any $t \geq 1$:

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

180 Since the first and second terms of the right-hand side of (5) are bounded by the distance between the
 181 stationary point and the anchoring strategy by using (6), we have:

$$r^{\tan}(\pi^t) \leq \mathcal{O}\left(\left\|\pi^{\mu, \sigma^{k(t-1)}} - \sigma^{k(t-1)}\right\|\right) + \mathcal{O}\left(\frac{1}{k(t-1)+1}\right). \quad (7)$$

182 **(3) Potential function for bounding the distance between $\pi^{\mu, \sigma^{k(t)-1}}$ and $\sigma^{k(t)-1}$.** To derive the
 183 upper bound on $\left\|\pi^{\mu, \sigma^{k(t)-1}} - \sigma^{k(t-1)}\right\|$, we define the following potential function $P^{k(t)}$:

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

184 By some algebra, we can see that $P^{k(t)}$ is approximately non-increasing (in Lemma B.3). That is, we
 185 have for any $t \geq 1$ such that $k(t) \geq 2$:

$$P^{k(t)+1} \leq P^{k(t)} + (k(t)+1)^2 \cdot \mathcal{O}\left(\left\|\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)+1}\right\| + \left\|\pi^{\mu, \sigma^{k(t)-1}} - \sigma^{k(t)}\right\|\right). \quad (8)$$

186 Using (6) again, it is easy to show that $\left\|\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)+1}\right\| + \left\|\pi^{\mu, \sigma^{k(t)-1}} - \sigma^{k(t)}\right\| \leq \mathcal{O}\left(\frac{1}{(k(t)+1)^3}\right)$
 187 for a sufficiently large T_σ . Therefore, under the assumption that $T_\sigma \geq \Omega(\ln T)$, by telescoping of (8)
 188 and some algebra, we can derive the following upper bound on $\left\|\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}\right\|$ (in Lemma B.2):

$$\left\|\pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)}\right\| \leq \mathcal{O}\left(\frac{1}{k(t)+1}\right). \quad (9)$$

189 **(4) Putting it all together: last-iterate convergence rate of π^t .** By combining (7) and (9), we get
 190 $r^{\tan}(\pi^t) \leq \mathcal{O}\left(\frac{1}{k(t-1)+1}\right)$. Therefore, since $k(t) = \lfloor \frac{t-1}{T_\sigma} \rfloor + 1$, it holds that $r^{\tan}(\pi^t) \leq \mathcal{O}\left(\frac{T_\sigma}{t+T_\sigma-2}\right)$.
 191 Finally, taking $T_\sigma = \Theta(\ln T)$, we have:

$$r^{\tan}(\pi^t) \leq \mathcal{O}\left(\frac{\ln T}{t-1}\right).$$

192 The upper bound on the gap function is immediately obtained since we have Lemma 2.2. \square

193 4.2 Noisy feedback setting

194 Next, we establish the last-iterate convergence rate in the *noisy feedback* setting, where each
 195 player i observes a noisy gradient vector contaminated by an additive noise vector $\xi_i^t \in \mathbb{R}^{d_i}$:
 196 $\widehat{\nabla}_{\pi_i} v_i(\pi^t) + \xi_i^t$. We assume that the noisy vector ξ_i^t is zero-mean and its variance is bounded.
 197 Formally, defining the sigma-algebra generated by the history of the observations as $\mathcal{F}_t :=$
 198 $\sigma\left(\left(\widehat{\nabla}_{\pi_i} v_i(\pi^1)\right)_{i \in [N]}, \dots, \left(\widehat{\nabla}_{\pi_i} v_i(\pi^{t-1})\right)_{i \in [N]}\right)$, $\forall t \geq 1$, the noisy vector ξ_i^t is assumed to satisfy
 199 the following conditions:

200 **Assumption 4.2.** For all $t \geq 1$ and $i \in [N]$, the noise vector ξ_i^t satisfies the following properties: (a)
 201 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
 202 $C > 0$.

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

207 **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
 208 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 $c \geq 1$.
 209 If we use the constant perturbation strength $\mu > 0$ and the decreasing learning rate sequence
 210 $\eta_t = \frac{1}{\kappa(t - T_\sigma(k(t) - 1)) + 2\theta}$, then the strategy π^{T+1} satisfies:

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

211 Note that the non-increasing property, as described in (8), of the potential function holds regardless
 212 of the presence of noise. This implies that a proof technique similar to the one used with the potential
 213 function in the full feedback setting can also be applied in the noisy feedback setting. The detailed
 214 proof can be found in Appendix C.

215 5 Individual regret bound

216 In this section, we present an upper bound on an individual regret for each player. Specifically,
 217 we examine two performance measures in our study: the *external regret* and the *dynamic regret*
 218 [Zinkevich, 2003]. The external regret is a conventional measure in online learning. In online learning
 219 in games, the external regret for player i is defined as the gap between the player's realized cumulative
 220 payoff and the cumulative payoff of the best fixed strategy in hindsight:

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

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

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

222 We show in Theorem 5.1 that the individual regret is at most $\mathcal{O}((\ln T)^2)$ if each player $i \in [N]$ plays
 223 according to GABP in the full feedback setting:

224 **Theorem 5.1.** In the same setup of Theorem 4.1, we have for any player $i \in [N]$ and $T \geq 3$:

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

225 This regret bound is significantly superior to the $\mathcal{O}(\sqrt{T})$ regret bound of Optimistic Gradient Ascent,
 226 and it is slightly inferior to the $\mathcal{O}(\ln T)$ regret bound of AOG [Cai and Zheng, 2023]. The proof is
 227 given in Appendix D.

228 6 Experiments

229 In this section, we present the empirical results of our GABP, comparing its performance with
 230 Adaptively Perturbed Gradient Ascent (APGA) [Abe et al., 2024] and Optimistic Gradient Ascent
 231 (OGA) [Daskalakis et al., 2018, Wei et al., 2021]. We conduct experiments on two classes of concave-
 232 convex games. The first experiment is carried out on random payoff games, which are two-player
 233 zero-sum normal-form games with payoff matrices of size d . In this game, each player's strategy

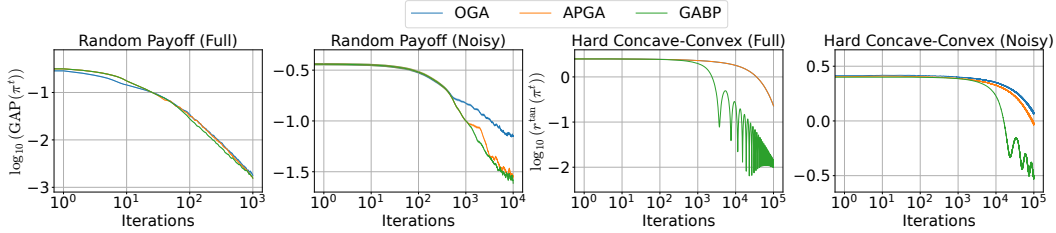


Figure 1: Performance of π^t for GABP, APGA, and OGA with full and noisy feedback in the random payoff and hard concave-convex games, respectively. The shaded area represents the standard errors. Note that we report the gap function for the random payoff game, while the tangent residual is reported for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

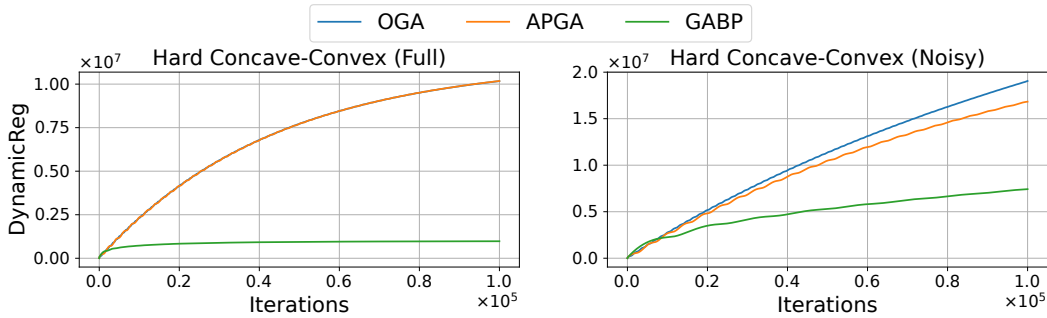


Figure 2: Dynamic regret for GABP, APGA, and OGA with full and noisy feedback.

234 space is represented by the d -dimensional probability simplex, i.e., $\mathcal{X}_1 = \mathcal{X}_2 = \Delta^d$. All entries
 235 of the payoff matrix are drawn independently from a uniform distribution over the interval $[-1, 1]$.
 236 We set $d = 50$ and the initial strategies are set to $\pi_1^1 = \pi_2^1 = \frac{1}{d}\mathbf{1}$. The second instance is a *hard*
 237 *concave-convex game* [Ouyang and Xu, 2021], formulated as the following max-min optimization
 238 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
 239 the setup in Cai and Zheng [2023], we choose $\mathcal{X}_1 = \mathcal{X}_2 = [-200, 200]^d$ with $d = 100$. The precise
 240 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 E.2. All algorithms
 241 are executed with initial strategies $\pi_1^1 = \pi_2^1 = \frac{1}{n}\mathbf{1}$. The detailed hyperparameters of the algorithms,
 242 tuned for best performance, are shown in Table 1 in Appendix E.3.

243 The numerical results of the random payoff game and the hard concave-convex game are shown in
 244 Figure 1. Both the full feedback and noisy feedback experiments in the random payoff game were
 245 conducted with 50 different random seeds, which corresponds to using 50 different payoff matrices.
 246 For experiments on the hard concave-convex game with noisy feedback, we use 10 different random
 247 seeds. We assume that the noise vector ξ_i^t is generated from the multivariate Gaussian distribution
 248 $\mathcal{N}(0, 0.1^2\mathbf{I})$ in an i.i.d. manner for both games. We observe that GABP exhibits competitive or faster
 249 performance over APGA and OGA in all experiments.

250 Figure 2 illustrates the dynamic regret in the hard concave-convex game. GABP exhibits lower
 251 regret than APGA and OGA in both settings, demonstrating its efficiency and robustness. Note that
 252 APGA and OGA exhibit almost identical trajectories with full feedback, with their plots overlapping
 253 completely.

254 7 Related literature

255 No-regret learning algorithms have been extensively studied with the intent of achieving key objectives
 256 such as average-iterate convergence or last-iterate convergence. Recently, learning algorithms
 257 introducing optimism [Rakhlin and Sridharan, 2013a,b], such as optimistic Follow the Regularized
 258 Leader [Shalev-Shwartz and Singer, 2006] and optimistic Mirror Descent [Zhou et al., 2017, Hsieh
 259 et al., 2021], have been introduced to admit last-iterate convergence in a broad spectrum of game

260 settings. These optimistic algorithms with full feedback have been shown to achieve last-iterate
 261 convergence in various classes of games, including bilinear games [Daskalakis et al., 2018, Daskalakis
 262 and Panageas, 2019, Liang and Stokes, 2019, de Montbrun and Renault, 2022], cocoercive games
 263 [Lin et al., 2020], and saddle point problems [Daskalakis and Panageas, 2018, Mertikopoulos et al.,
 264 2019, Golowich et al., 2020b, Wei et al., 2021, Lei et al., 2021, Yoon and Ryu, 2021, Lee and Kim,
 265 2021, Cevher et al., 2023]. Recent studies have provided finite convergence rates for monotone games
 266 [Golowich et al., 2020a, Cai et al., 2022a,b, Gorbunov et al., 2022, Cai and Zheng, 2023].

267 Compared to the full feedback setting, there are significant challenges in learning with noisy feedback.
 268 For example, a learning algorithm must estimate the gradient from feedback that is contaminated by
 269 noise. Despite the challenge, a vast literature has successfully achieved last-iterate convergence with
 270 noisy feedback in specific classes of games, including potential games [Cohen et al., 2017], strongly
 271 monotone games [Giannou et al., 2021b,a], and two-player zero-sum games [Abe et al., 2023]. These
 272 results have often leveraged unique structures of their payoff functions, such as strict (or strong)
 273 monotonicity [Bravo et al., 2018, Kannan and Shanbhag, 2019, Hsieh et al., 2019, Anagnostides
 274 and Panageas, 2022] and strict variational stability [Mertikopoulos et al., 2019, Azizian et al., 2021,
 275 Mertikopoulos and Zhou, 2019, Mertikopoulos et al., 2022]. Without these restrictions, convergence
 276 is mainly demonstrated in an asymptotic manner, with no quantification of the rate [Koshal et al.,
 277 2010, 2013, Yousefian et al., 2017, Tatarenko and Kamgarpour, 2019, Hsieh et al., 2020, 2022, Abe
 278 et al., 2023]. Consequently, an exceedingly large number of iterations might be necessary to reach an
 279 equilibrium.

280 There have been several studies focusing on payoff-regularized learning, where each player’s payoff
 281 or utility function is perturbed or regularized via strongly convex functions [Cen et al., 2021, 2023,
 282 Pattathil et al., 2023]. Previous studies have successfully achieved convergence to stationary points,
 283 which are approximate equilibria. For instance, Sokota et al. [2023] have demonstrated that their
 284 perturbed mirror descent algorithm converges to a quantal response equilibrium [McKelvey and
 285 Palfrey, 1995, 1998]. Similar results have been obtained with the Boltzmann Q-learning dynam-
 286 ics [Tuyls et al., 2006] and penalty-regularized dynamics [Coucheney et al., 2015] in continuous-time
 287 settings [Leslie and Collins, 2005, Abe et al., 2022, Hussain et al., 2023]. To ensure convergence
 288 toward a Nash equilibrium of the underlying game, the magnitude of perturbation requires careful
 289 adjustment. Several learning algorithms have been proposed to gradually reduce the perturbation
 290 strength μ in response to this [Bernasconi et al., 2022, Liu et al., 2023, Cai et al., 2023]. These
 291 include well-studied methods such as iterative Tikhonov regularization [Facchinei and Pang, 2003,
 292 Koshal et al., 2010, Tatarenko and Kamgarpour, 2019]. Alternatively, Perolat et al. [2021] and Abe
 293 et al. [2023] have employed a payoff perturbation scheme, where the magnitude of perturbation is
 294 determined by the distance from an anchoring strategy, which is periodically re-initialized by the
 295 current strategy. Recently, Abe et al. [2024] have established $\tilde{O}(1/\sqrt{T})$ and $\tilde{O}(1/T^{1/10})$ last-iterate
 296 convergence rates for the payoff perturbation scheme in the full/noisy feedback setting, respectively.
 297 Our algorithm achieves faster $\tilde{O}(1/T)$ and $\tilde{O}(1/T^{1/2})$ last-iterate convergence rates by modifying the
 298 periodically re-initializing anchoring strategy scheme so that the anchoring strategy evolves more
 299 gradually.

300 8 Conclusion

301 This study proposes a novel payoff-perturbed algorithm, Gradient Ascent with Boosting Payoff
 302 Perturbation, which achieves $\tilde{O}(1/T)$ and $\tilde{O}(1/T^{1/2})$ last-iterate convergence rates in monotone
 303 games with full/noisy feedback, respectively. Extending our results in settings where each player
 304 only observes bandit feedback is an intriguing and challenging future direction.

305 References

- 306 Kenshi Abe, Mitsuki Sakamoto, and Atsushi Iwasaki. Mutation-driven follow the regularized leader
 307 for last-iterate convergence in zero-sum games. In *UAI*, pages 1–10, 2022.
- 308 Kenshi Abe, Kaito Ariu, Mitsuki Sakamoto, Kentaro Toyoshima, and Atsushi Iwasaki. Last-iterate
 309 convergence with full and noisy feedback in two-player zero-sum games. In *AISTATS*, pages
 310 7999–8028, 2023.

- 311 Kenshi Abe, Kaito Ariu, Mitsuki Sakamoto, and Atsushi Iwasaki. Adaptively perturbed mirror
312 descent for learning in games. In *ICML*, 2024.
- 313 Ioannis Anagnostides and Ioannis Panageas. Frequency-domain representation of first-order methods:
314 A simple and robust framework of analysis. In *SOSA*, pages 131–160, 2022.
- 315 Waïss Azizian, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. The last-iterate
316 convergence rate of optimistic mirror descent in stochastic variational inequalities. In *COLT*, pages
317 326–358, 2021.
- 318 James P Bailey and Georgios Piliouras. Multiplicative weights update in zero-sum games. In
319 *Economics and Computation*, pages 321–338, 2018.
- 320 Martino Bernasconi, Alberto Marchesi, and Francesco Trovò. Last-iterate convergence to trembling-
321 hand perfect equilibria. *arXiv preprint arXiv:2208.08238*, 2022.
- 322 Mario Bravo, David Leslie, and Panayotis Mertikopoulos. Bandit learning in concave N-person
323 games. In *NeurIPS*, pages 5666–5676, 2018.
- 324 Yang Cai and Constantinos Daskalakis. On minmax theorems for multiplayer games. In *SODA*,
325 pages 217–234, 2011.
- 326 Yang Cai and Weiqiang Zheng. Doubly optimal no-regret learning in monotone games. In *ICML*,
327 pages 3507–3524, 2023.
- 328 Yang Cai, Ozan Candogan, Constantinos Daskalakis, and Christos Papadimitriou. Zero-sum polyma-
329 trix games: A generalization of minmax. *Mathematics of Operations Research*, 41(2):648–655,
330 2016.
- 331 Yang Cai, Argyris Oikonomou, and Weiqiang Zheng. Finite-time last-iterate convergence for learning
332 in multi-player games. In *NeurIPS*, pages 33904–33919, 2022a.
- 333 Yang Cai, Argyris Oikonomou, and Weiqiang Zheng. Tight last-iterate convergence of the extragradi-
334 ent method for constrained monotone variational inequalities. *arXiv preprint arXiv:2204.09228*,
335 2022b.
- 336 Yang Cai, Haipeng Luo, Chen-Yu Wei, and Weiqiang Zheng. Uncoupled and convergent learning in
337 two-player zero-sum markov games with bandit feedback. In *NeurIPS*, pages 36364–36406, 2023.
- 338 Shicong Cen, Yuting Wei, and Yuejie Chi. Fast policy extragradient methods for competitive games
339 with entropy regularization. In *NeurIPS*, pages 27952–27964, 2021.
- 340 Shicong Cen, Yuejie Chi, Simon S Du, and Lin Xiao. Faster last-iterate convergence of policy
341 optimization in zero-sum Markov games. In *ICLR*, 2023.
- 342 Volkan Cevher, Georgios Piliouras, Ryann Sim, and Stratis Skoulakis. Min-max optimization made
343 simple: Approximating the proximal point method via contraction maps. In *Symposium on*
344 *Simplicity in Algorithms (SOSA)*, pages 192–206, 2023.
- 345 Johanne Cohen, Amélie Héliou, and Panayotis Mertikopoulos. Learning with bandit feedback in
346 potential games. In *NeurIPS*, pages 6372–6381, 2017.
- 347 Pierre Coucheney, Bruno Gaujal, and Panayotis Mertikopoulos. Penalty-regulated dynamics and
348 robust learning procedures in games. *Mathematics of Operations Research*, 40(3):611–633, 2015.
- 349 Constantinos Daskalakis and Ioannis Panageas. The limit points of (optimistic) gradient descent in
350 min-max optimization. In *NeurIPS*, pages 9256–9266, 2018.
- 351 Constantinos Daskalakis and Ioannis Panageas. Last-iterate convergence: Zero-sum games and
352 constrained min-max optimization. In *ITCS*, pages 27:1–27:18, 2019.
- 353 Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training GANs with
354 optimism. In *ICLR*, 2018.

- 355 Étienne de Montbrun and Jérôme Renault. Convergence of optimistic gradient descent ascent in
356 bilinear games. *arXiv preprint arXiv:2208.03085*, 2022.
- 357 Gerard Debreu. A social equilibrium existence theorem. *Proceedings of the National Academy of*
358 *Sciences*, 38(10):886–893, 1952.
- 359 Francisco Facchinei and Jong-Shi Pang. *Finite-dimensional variational inequalities and complemen-*
360 *tarity problems*. Springer, 2003.
- 361 Angeliki Giannou, Emmanouil Vasileios Vlatakis-Gkaragkounis, and Panayotis Mertikopoulos.
362 Survival of the strictest: Stable and unstable equilibria under regularized learning with partial
363 information. In *COLT*, pages 2147–2148, 2021a.
- 364 Angeliki Giannou, Emmanouil-Vasileios Vlatakis-Gkaragkounis, and Panayotis Mertikopoulos. On
365 the rate of convergence of regularized learning in games: From bandits and uncertainty to optimism
366 and beyond. In *NeurIPS*, pages 22655–22666, 2021b.
- 367 Noah Golowich, Sarath Pattathil, and Constantinos Daskalakis. Tight last-iterate convergence rates
368 for no-regret learning in multi-player games. In *NeurIPS*, pages 20766–20778, 2020a.
- 369 Noah Golowich, Sarath Pattathil, Constantinos Daskalakis, and Asuman Ozdaglar. Last iterate is
370 slower than averaged iterate in smooth convex-concave saddle point problems. In *COLT*, pages
371 1758–1784, 2020b.
- 372 Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
373 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, pages 2672–2680,
374 2014.
- 375 Eduard Gorbunov, Adrien Taylor, and Gauthier Gidel. Last-iterate convergence of optimistic gradient
376 method for monotone variational inequalities. In *NeurIPS*, pages 21858–21870, 2022.
- 377 Yu-Guan Hsieh, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. On the convergence
378 of single-call stochastic extra-gradient methods. In *NeurIPS*, pages 6938–6948, 2019.
- 379 Yu-Guan Hsieh, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. Explore aggressively,
380 update conservatively: Stochastic extragradient methods with variable stepsize scaling. In *NeurIPS*,
381 pages 16223–16234, 2020.
- 382 Yu-Guan Hsieh, Kimon Antonakopoulos, and Panayotis Mertikopoulos. Adaptive learning in
383 continuous games: Optimal regret bounds and convergence to Nash equilibrium. In *COLT*, pages
384 2388–2422, 2021.
- 385 Yu-Guan Hsieh, Kimon Antonakopoulos, Volkan Cevher, and Panayotis Mertikopoulos. No-regret
386 learning in games with noisy feedback: Faster rates and adaptivity via learning rate separation. In
387 *NeurIPS*, pages 6544–6556, 2022.
- 388 Aamal Abbas Hussain, Francesco Belardinelli, and Georgios Piliouras. Asymptotic convergence and
389 performance of multi-agent Q-learning dynamics. *arXiv preprint arXiv:2301.09619*, 2023.
- 390 Aswin Kannan and Uday V. Shanbhag. Optimal stochastic extragradient schemes for pseudomonotone
391 stochastic variational inequality problems and their variants. *Computational Optimization and*
392 *Applications*, 74(3):779–820, 2019.
- 393 Jayash Koshal, Angelia Nedić, and Uday V Shanbhag. Single timescale regularized stochastic
394 approximation schemes for monotone nash games under uncertainty. In *CDC*, pages 231–236.
395 IEEE, 2010.
- 396 Jayash Koshal, Angelia Nedic, and Uday V. Shanbhag. Regularized iterative stochastic approximation
397 methods for stochastic variational inequality problems. *IEEE Transactions on Automatic Control*,
398 58(3):594–609, 2013.
- 399 Suchool Lee and Donghwan Kim. Fast extra gradient methods for smooth structured nonconvex-
400 nonconcave minimax problems. In *NeurIPS*, pages 22588–22600, 2021.

- 401 Qi Lei, Sai Ganesh Nagarajan, Ioannis Panageas, et al. Last iterate convergence in no-regret learning:
402 constrained min-max optimization for convex-concave landscapes. In *AISTATS*, pages 1441–1449,
403 2021.
- 404 David S Leslie and Edmund J Collins. Individual q-learning in normal form games. *SIAM Journal*
405 *on Control and Optimization*, 44(2):495–514, 2005.
- 406 Tengyuan Liang and James Stokes. Interaction matters: A note on non-asymptotic local convergence
407 of generative adversarial networks. In *AISTATS*, pages 907–915, 2019.
- 408 Tianyi Lin, Zhengyuan Zhou, Panayotis Mertikopoulos, and Michael Jordan. Finite-time last-iterate
409 convergence for multi-agent learning in games. In *ICML*, pages 6161–6171, 2020.
- 410 Mingyang Liu, Asuman Ozdaglar, Tiancheng Yu, and Kaiqing Zhang. The power of regularization in
411 solving extensive-form games. In *ICLR*, 2023.
- 412 Richard D McKelvey and Thomas R Palfrey. Quantal response equilibria for normal form games.
413 *Games and economic behavior*, 10(1):6–38, 1995.
- 414 Richard D McKelvey and Thomas R Palfrey. Quantal response equilibria for extensive form games.
415 *Experimental economics*, 1:9–41, 1998.
- 416 Panayotis Mertikopoulos and Zhengyuan Zhou. Learning in games with continuous action sets and
417 unknown payoff functions. *Mathematical Programming*, 173(1):465–507, 2019.
- 418 Panayotis Mertikopoulos, Bruno Lecouat, Houssam Zenati, Chuan-Sheng Foo, Vijay Chandrasekhar,
419 and Georgios Piliouras. Optimistic mirror descent in saddle-point problems: Going the extra
420 (gradient) mile. In *ICLR*, 2019.
- 421 Panayotis Mertikopoulos, Ya-Ping Hsieh, and Volkan Cevher. Learning in games from a stochastic
422 approximation viewpoint. *arXiv preprint arXiv:2206.03922*, 2022.
- 423 Rémi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland,
424 Zhaohan Daniel Guo, Yunhao Tang, Matthieu Geist, Thomas Mesnard, Andrea Michi, et al. Nash
425 learning from human feedback. *arXiv preprint arXiv:2312.00886*, 2023.
- 426 John Nash. Non-cooperative games. *Annals of mathematics*, pages 286–295, 1951.
- 427 Angelia Nedić and Soomin Lee. On stochastic subgradient mirror-descent algorithm with weighted
428 averaging. *SIAM Journal on Optimization*, 24(1):84–107, 2014.
- 429 A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to
430 stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- 431 Yuyuan Ouyang and Yangyang Xu. Lower complexity bounds of first-order methods for convex-
432 concave bilinear saddle-point problems. *Mathematical Programming*, 185(1):1–35, 2021.
- 433 Sarath Pattathil, Kaiqing Zhang, and Asuman Ozdaglar. Symmetric (optimistic) natural policy
434 gradient for multi-agent learning with parameter convergence. In *AISTATS*, pages 5641–5685,
435 2023.
- 436 Julien Perolat, Remi Munos, Jean-Baptiste Lespiau, Shayegan Omidshafiei, Mark Rowland, Pedro
437 Ortega, Neil Burch, Thomas Anthony, David Balduzzi, Bart De Vylder, et al. From Poincaré
438 recurrence to convergence in imperfect information games: Finding equilibrium via regularization.
439 In *ICML*, pages 8525–8535, 2021.
- 440 Alexander Rakhlin and Karthik Sridharan. Online learning with predictable sequences. In *COLT*,
441 pages 993–1019, 2013a.
- 442 Sasha Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable sequences.
443 In *NeurIPS*, pages 3066–3074, 2013b.
- 444 Shai Shalev-Shwartz and Yoram Singer. Convex repeated games and fenchel duality. *Advances in*
445 *neural information processing systems*, 19, 2006.

- 446 Samuel Sokota, Ryan D’Orazio, J Zico Kolter, Nicolas Loizou, Marc Lanctot, Ioannis Mitliagkas,
447 Noam Brown, and Christian Kroer. A unified approach to reinforcement learning, quantal response
448 equilibria, and two-player zero-sum games. In *ICLR*, 2023.
- 449 Gokul Swamy, Christoph Dann, Rahul Kidambi, Zhiwei Steven Wu, and Alekh Agarwal. A minimaxi-
450 malist approach to reinforcement learning from human feedback. *arXiv preprint arXiv:2401.04056*,
451 2024.
- 452 Tatiana Tatarenko and Maryam Kamgarpour. Learning Nash equilibria in monotone games. In *CDC*,
453 pages 3104–3109. IEEE, 2019.
- 454 Karl Tuyls, Pieter Jan Hoen, and Bram Vanschoenwinkel. An evolutionary dynamical analysis
455 of multi-agent learning in iterated games. *Autonomous Agents and Multi-Agent Systems*, 12(1):
456 115–153, 2006.
- 457 Chen-Yu Wei, Chung-Wei Lee, Mengxiao Zhang, and Haipeng Luo. Linear last-iterate convergence
458 in constrained saddle-point optimization. In *ICLR*, 2021.
- 459 TaeHo Yoon and Ernest K Ryu. Accelerated algorithms for smooth convex-concave minimax
460 problems with $\mathcal{O}(1/k^2)$ rate on squared gradient norm. In *ICML*, pages 12098–12109, 2021.
- 461 Farzad Yousefian, Angelia Nedić, and Uday V Shanbhag. On smoothing, regularization, and averaging
462 in stochastic approximation methods for stochastic variational inequality problems. *Mathematical*
463 *Programming*, 165:391–431, 2017.
- 464 Zhengyuan Zhou, Panayotis Mertikopoulos, Aris L Moustakas, Nicholas Bambos, and Peter Glynn.
465 Mirror descent learning in continuous games. In *CDC*, pages 5776–5783. IEEE, 2017.
- 466 Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In
467 *ICML*, pages 928–936, 2003.

468 A Broader impact

469 Our study can bring about a positive impact on society by contributing to the advancement of the
470 Game AI industry. However, as far as we can envision, there are no conceivable negative social
471 impacts.

472 B Proofs for Theorem 4.1

473 B.1 Proof of Theorem 4.1

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

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

475 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
476 residual for π^t can be bounded as:

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

477 Letting us define

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

478 then we get by triangle inequality:

$$\begin{aligned}
r^{\tan}(\pi^t) &\leq \left\| -V(\pi^t) + V(\pi^{t-1}) - \frac{\mu}{k(t-1)+1}(\sigma^{k(t-1)} - \sigma^1) \right. \\
&\quad \left. - \mu(\pi^{\mu, \sigma^{k(t-1)}} - \pi^{\mu, \sigma^{k(t-1)}} + \pi^{t-1} - \sigma^{k(t-1)}) - \frac{1}{\eta}(\pi^t - \pi^{t-1}) \right\| \\
&\leq \left\| -V(\pi^t) + V(\pi^{t-1}) \right\| + \frac{\mu}{k(t-1)+1} \left\| \sigma^{k(t-1)} - \sigma^1 \right\| \\
&\quad + \mu \left\| \pi^{\mu, \sigma^{k(t-1)}} - \sigma^{k(t-1)} \right\| + \mu \left\| \pi^{\mu, \sigma^{k(t-1)}} - \pi^{t-1} \right\| + \frac{1}{\eta} \left\| \pi^t - \pi^{t-1} \right\| \\
&\leq \frac{1+\eta L}{\eta} \left\| \pi^t - \pi^{t-1} \right\| + \frac{\mu D}{k(t-1)+1} \\
&\quad + \mu \left\| \pi^{\mu, \sigma^{k(t-1)}} - \sigma^{k(t-1)} \right\| + \mu \left\| \pi^{\mu, \sigma^{k(t-1)}} - \pi^{t-1} \right\| \\
&\leq \frac{1+\eta L}{\eta} \left\| \pi^{\mu, \sigma^{k(t-1)}} - \pi^t \right\| + \frac{\mu D}{k(t-1)+1} + \mu \left\| \pi^{\mu, \sigma^{k(t-1)}} - \sigma^{k(t-1)} \right\| \\
&\quad + \left(\mu + \frac{1+\eta L}{\eta} \right) \left\| \pi^{\mu, \sigma^{k(t-1)}} - \pi^{t-1} \right\|. \tag{10}
\end{aligned}$$

479 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
480 lemma:

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

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

482 Combining (10) and Lemma B.1, we have:

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

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

484 **Lemma B.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 $t \geq 1$:*

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

485 By combining (11) and Lemma B.2, we get:

$$\begin{aligned}
r^{\tan}(\pi^t) &\leq \frac{16D}{k(t-1)+1} \left(\mu + \frac{1+\eta L}{\eta} \right) + \frac{\mu D}{k(t-1)+1} \\
&\leq \frac{17D}{k(t-1)+1} \left(\mu + \frac{1+\eta L}{\eta} \right).
\end{aligned}$$

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

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

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

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

488 □

489 **B.2 Proof of Lemma B.1**

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

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

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

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

492 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} , we
493 have for any $t \geq 1$:

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

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

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

495 Combining (12), (13), and (14) yields:

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

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

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

498 where the second inequality follow from (2), and the last inequality follows from the assumption that
 499 $\eta \leq \frac{\mu}{(L+\mu)^2}$. By combining (15) and (16), we get:

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

500 Thus,

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

508

□

509 **B.3 Proof of Lemma B.2**

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

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

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

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

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

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

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

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

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

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

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

515 Here, we introduce the following lemma:

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

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

517 By combining (20) and Lemma B.3, we get:

$$\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 \left\| \sigma^1 - \pi^* \right\|^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 \left\| \sigma^1 - \pi^* \right\|^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} - \sigma^1 \right\|^2 + 2 \left\langle 2\sigma^2 + \sigma^1 - 3\pi^{\mu, \sigma^1}, \pi^{\mu, \sigma^1} - \sigma^1 \right\rangle + 2 \left\| \sigma^1 - \pi^* \right\|^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} - \sigma^1 \right\|^2 + 2 \left\langle \sigma^1 - \pi^{\mu, \sigma^1}, \pi^{\mu, \sigma^1} - \sigma^1 \right\rangle + 4 \left\langle \sigma^2 - \pi^{\mu, \sigma^1}, \pi^{\mu, \sigma^1} - \sigma^1 \right\rangle \\
& \quad + 2 \left\| \sigma^1 - \pi^* \right\|^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} - \sigma^1 \right\|^2 + 4 \left\langle \sigma^2 - \pi^{\mu, \sigma^1}, \pi^{\mu, \sigma^1} - \sigma^1 \right\rangle + 2 \left\| \sigma^1 - \pi^* \right\|^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} - \sigma^1 \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 \\
& \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}$$

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

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

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

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

520 Therefore, from Cauchy-Schwarz inequality, we have:

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

521 Furthermore, from Lemma B.1, we have for any $k \geq 1$:

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

522 Combining (21) and (22), we have for any $t \geq 1$ such that $k(t) \geq 2$:

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

523 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 $k(t) \geq 2$:

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

524 and then:

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

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

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

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

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

527

□

528 **B.4 Proof of Lemma B.3**

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

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

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

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

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

$$\begin{aligned} 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 \\ &\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 \\ &\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 \\ &= -\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 \\ &\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 \\ &= -\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 \\ &\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 \\ &= -\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 \\ &\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. \end{aligned}$$

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

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

533 Thus, we have:

$$\begin{aligned} 0 &\leq -\mu \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 \\ &\quad + \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)-1}} - \sigma^{k(t)} \right\|^2 + \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 \\ &\quad + \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \sigma^{k(t)} \right\|^2 - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)-1}} - \sigma^{k(t)} \right\|^2 - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 \\ &\quad + \frac{\mu}{2} \left\| \hat{\sigma}^{k(t)-1} - \hat{\sigma}^{k(t)} \right\|^2 + \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 \\ &\quad - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} + \hat{\sigma}^{k(t)-1} + \hat{\sigma}^{k(t)} \right\|^2 \\ &= -\frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 + \frac{\mu}{2} \left\| \hat{\sigma}^{k(t)} - \hat{\sigma}^{k(t)-1} \right\|^2 \\ &\quad - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} + \hat{\sigma}^{k(t)-1} + \hat{\sigma}^{k(t)} \right\|^2 \\ &\leq -\frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 + \frac{\mu}{2} \left\| \hat{\sigma}^{k(t)} - \hat{\sigma}^{k(t)-1} \right\|^2 \\ &= -\frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 + \frac{\mu}{2} \left\| \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}} + \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\|^2 \end{aligned}$$

$$\begin{aligned}
&= \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\|^2 + \frac{\mu}{2} \left\| \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 - \frac{\mu}{2} \left\| \pi^{\mu, \sigma^{k(t)}} - \pi^{\mu, \sigma^{k(t)-1}} \right\|^2 \\
&\quad + \mu \left\langle \hat{\sigma}^{k(t)} - \pi^{\mu, \sigma^{k(t)-1}}, \pi^{\mu, \sigma^{k(t)-1}} - \hat{\sigma}^{k(t)-1} \right\rangle. \tag{23}
\end{aligned}$$

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

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

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

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

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

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

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

537

□

538 C Proofs for Theorem 4.3

539 C.1 Proof of Theorem 4.3

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

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

542 where the last inequality follows from Lemma 2.2. From the first-order optimality condition for
543 π^{μ, σ^K} , we have for any $x \in \mathcal{X}$:

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

544 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
545 be bounded by:

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

546 Taking its expectation yields:

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

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

548 **Lemma C.1.** Let $\kappa = \frac{\mu}{2}, \theta = \frac{3\mu^2 + 8L^2}{2\mu}$. Suppose that Assumption 4.2 holds. If we set $\eta_t =$
549 $\frac{1}{\kappa(t - T_\sigma(k(t)-1)) + 2\theta}$, we have for any $t \geq 1$:

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

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

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

552 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:

553 **Lemma C.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 $t \geq 1$:*

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

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

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

555 Combining (25), (26), and (27), we have:

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

556 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
 557 any $T \geq T_\sigma$:

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

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

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

559

□

560 C.2 Proof of Lemma C.1

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

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

562 Combining (28), (12), and (14), we have:

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

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

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

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

565 By combining (29) and (30), we have:

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

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

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

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

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

568 Rearranging and taking the expectations, we get:

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

569 Since $k(s) = k(t)$ for any $s \in [(k(t) - 1)T_\sigma + 1, T]$, telescoping the sum yields:

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

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

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

571 Therefore,

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

572 Here, we have:

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

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

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

574

□

575 C.3 Proof of Lemma C.2

576 *Proof of Lemma C.2.* First, from Lemma C.1, we have for any $k \geq 1$:

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

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

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

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

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

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

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

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

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

581 Then, we have:

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

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

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

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

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

584 □

585 **D Proof of Theorem 5.1**

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

$$\begin{aligned} \text{DynamicReg}_i(T) &= \sum_{t=1}^T \left(\max_{x \in \mathcal{X}_i} v_i(x, \pi_{-i}^t) - v_i(\pi^t) \right) \\ &\leq \mathcal{O}(1) + \sum_{t=3}^T \sum_{i=1}^N \left(\max_{x \in \mathcal{X}_i} v_i(x, \pi_{-i}^t) - v_i(\pi^t) \right). \end{aligned}$$

587 Here, we introduce the following lemma:

588 **Lemma D.1** (Lemma 2 of [Cai et al. \[2022a\]](#)). For any $\pi \in \mathcal{X}$, we have:

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

589 Therefore, we have:

$$\text{DynamicReg}_i(T) \leq \mathcal{O}(1) + \sum_{t=3}^T \text{GAP}(\pi^t).$$

590 Thus, from [Theorem 4.1](#):

$$\begin{aligned} \text{DynamicReg}_i(T) &\leq \mathcal{O}(1) + \sum_{t=3}^T \mathcal{O}\left(\frac{\ln T}{t}\right) \\ &\leq \mathcal{O}((\ln T)^2). \end{aligned}$$

591

□

592 E Experimental details

593 E.1 Information on the computer resources

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

595 E.2 Hard concave-convex game

596 Following the setup in [Ouyang and Xu \[2021\]](#), [Cai and Zheng \[2023\]](#), we choose

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

597 and $H = 2A^\top A$.

598 E.3 Hyperparameters

599 For each game, we carefully tuned the hyperparameters for each algorithm to ensure optimal perfor-
600 mance. The specific parameters for each game and setting are summarized in [Table 1](#).

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

Table 1: Hyperparameters

601 **F Relationship with accelerated optimistic gradient algorithm**

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

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

604 This can be equivalently written as:

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

605 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
 606 strategy π_i^1 as the proximal point in gradient ascent. However, our GABP diverges from AOG in that it
 607 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
 608 term.

609 **NeurIPS Paper Checklist**

610 **1. Claims**

611 Question: Do the main claims made in the abstract and introduction accurately reflect the
612 paper’s contributions and scope?

613 Answer: [\[Yes\]](#)

614 Justification: We have clearly stated the contributions and scope of this study.

615 Guidelines:

- 616 • The answer NA means that the abstract and introduction do not include the claims
617 made in the paper.
- 618 • The abstract and/or introduction should clearly state the claims made, including the
619 contributions made in the paper and important assumptions and limitations. A No or
620 NA answer to this question will not be perceived well by the reviewers.
- 621 • The claims made should match theoretical and experimental results, and reflect how
622 much the results can be expected to generalize to other settings.
- 623 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
624 are not attained by the paper.

625 **2. Limitations**

626 Question: Does the paper discuss the limitations of the work performed by the authors?

627 Answer: [\[Yes\]](#)

628 Justification: In “Introduction” and “Conclusion”, we have reiterated the limitation of this
629 study.

630 Guidelines:

- 631 • The answer NA means that the paper has no limitation while the answer No means that
632 the paper has limitations, but those are not discussed in the paper.
- 633 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 634 • The paper should point out any strong assumptions and how robust the results are to
635 violations of these assumptions (e.g., independence assumptions, noiseless settings,
636 model well-specification, asymptotic approximations only holding locally). The authors
637 should reflect on how these assumptions might be violated in practice and what the
638 implications would be.
- 639 • The authors should reflect on the scope of the claims made, e.g., if the approach was
640 only tested on a few datasets or with a few runs. In general, empirical results often
641 depend on implicit assumptions, which should be articulated.
- 642 • The authors should reflect on the factors that influence the performance of the approach.
643 For example, a facial recognition algorithm may perform poorly when image resolution
644 is low or images are taken in low lighting. Or a speech-to-text system might not be
645 used reliably to provide closed captions for online lectures because it fails to handle
646 technical jargon.
- 647 • The authors should discuss the computational efficiency of the proposed algorithms
648 and how they scale with dataset size.
- 649 • If applicable, the authors should discuss possible limitations of their approach to
650 address problems of privacy and fairness.
- 651 • While the authors might fear that complete honesty about limitations might be used by
652 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
653 limitations that aren’t acknowledged in the paper. The authors should use their best
654 judgment and recognize that individual actions in favor of transparency play an impor-
655 tant role in developing norms that preserve the integrity of the community. Reviewers
656 will be specifically instructed to not penalize honesty concerning limitations.

657 **3. Theory Assumptions and Proofs**

658 Question: For each theoretical result, does the paper provide the full set of assumptions and
659 a complete (and correct) proof?

660 Answer: [\[Yes\]](#)

661 Justification: Please see the theoretical results and their proofs in the Appendix.

662 Guidelines:

- 663 • The answer NA means that the paper does not include theoretical results.
- 664 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
665 referenced.
- 666 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 667 • The proofs can either appear in the main paper or the supplemental material, but if
668 they appear in the supplemental material, the authors are encouraged to provide a short
669 proof sketch to provide intuition.
- 670 • Inversely, any informal proof provided in the core of the paper should be complemented
671 by formal proofs provided in appendix or supplemental material.
- 672 • Theorems and Lemmas that the proof relies upon should be properly referenced.

673 4. Experimental Result Reproducibility

674 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
675 perimental results of the paper to the extent that it affects the main claims and/or conclusions
676 of the paper (regardless of whether the code and data are provided or not)?

677 Answer: [Yes]

678 Justification: We have provided descriptions of experimental setups in the experiments
679 section.

680 Guidelines:

- 681 • The answer NA means that the paper does not include experiments.
- 682 • If the paper includes experiments, a No answer to this question will not be perceived
683 well by the reviewers: Making the paper reproducible is important, regardless of
684 whether the code and data are provided or not.
- 685 • If the contribution is a dataset and/or model, the authors should describe the steps taken
686 to make their results reproducible or verifiable.
- 687 • Depending on the contribution, reproducibility can be accomplished in various ways.
688 For example, if the contribution is a novel architecture, describing the architecture fully
689 might suffice, or if the contribution is a specific model and empirical evaluation, it may
690 be necessary to either make it possible for others to replicate the model with the same
691 dataset, or provide access to the model. In general, releasing code and data is often
692 one good way to accomplish this, but reproducibility can also be provided via detailed
693 instructions for how to replicate the results, access to a hosted model (e.g., in the case
694 of a large language model), releasing of a model checkpoint, or other means that are
695 appropriate to the research performed.
- 696 • While NeurIPS does not require releasing code, the conference does require all submis-
697 sions to provide some reasonable avenue for reproducibility, which may depend on the
698 nature of the contribution. For example
 - 699 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
700 to reproduce that algorithm.
 - 701 (b) If the contribution is primarily a new model architecture, the paper should describe
702 the architecture clearly and fully.
 - 703 (c) If the contribution is a new model (e.g., a large language model), then there should
704 either be a way to access this model for reproducing the results or a way to reproduce
705 the model (e.g., with an open-source dataset or instructions for how to construct
706 the dataset).
 - 707 (d) We recognize that reproducibility may be tricky in some cases, in which case
708 authors are welcome to describe the particular way they provide for reproducibility.
709 In the case of closed-source models, it may be that access to the model is limited in
710 some way (e.g., to registered users), but it should be possible for other researchers
711 to have some path to reproducing or verifying the results.

712 5. Open access to data and code

713 Question: Does the paper provide open access to the data and code, with sufficient instruc-
714 tions to faithfully reproduce the main experimental results, as described in supplemental
715 material?

716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766

Answer: [Yes]

Justification: We have included the experimental code in the supplementary material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We have provided descriptions of experimental setups in the experiments section.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Please see Figures.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 767
- 768
- 769
- 770
- 771
- 772
- 773
- 774
- 775
- 776
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
 - It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

777 8. Experiments Compute Resources

778 Question: For each experiment, does the paper provide sufficient information on the com-
779 puter resources (type of compute workers, memory, time of execution) needed to reproduce
780 the experiments?

781 Answer: [Yes]

782 Justification: We have shown the computer resources for this study in Appendix E.

783 Guidelines:

- 784
- 785
- 786
- 787
- 788
- 789
- 790
- 791
- The answer NA means that the paper does not include experiments.
 - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
 - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

792 9. Code Of Ethics

793 Question: Does the research conducted in the paper conform, in every respect, with the
794 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

795 Answer: [Yes]

796 Justification: We have carefully reviewed and followed the NeurIPS Code of Ethics.

797 Guidelines:

- 798
- 799
- 800
- 801
- 802
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

803 10. Broader Impacts

804 Question: Does the paper discuss both potential positive societal impacts and negative
805 societal impacts of the work performed?

806 Answer: [Yes]

807 Justification: We have described the potential societal impacts of our study in Appendix A.

808 Guidelines:

- 809
- 810
- 811
- 812
- 813
- 814
- 815
- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- 816
- 817
- 818
- 819
- 820
- 821
- 822
- 823
- 824
- 825
- 826
- 827
- 828
- 829
- 830
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

831 11. Safeguards

832 Question: Does the paper describe safeguards that have been put in place for responsible
833 release of data or models that have a high risk for misuse (e.g., pretrained language models,
834 image generators, or scraped datasets)?

835 Answer: [NA]

836 Justification: There are no such risks associated with the paper.

837 Guidelines:

- 838
- 839
- 840
- 841
- 842
- 843
- 844
- 845
- 846
- 847
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

848 12. Licenses for existing assets

849 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
850 the paper, properly credited and are the license and terms of use explicitly mentioned and
851 properly respected?

852 Answer: [NA]

853 Justification: This study does not use existing assets.

854 Guidelines:

- 855
- 856
- 857
- 858
- 859
- 860
- 861
- 862
- 863
- 864
- 865
- 866
- 867
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

868 • If this information is not available online, the authors are encouraged to reach out to
869 the asset’s creators.

870 13. New Assets

871 Question: Are new assets introduced in the paper well documented and is the documentation
872 provided alongside the assets?

873 Answer: [NA]

874 Justification: This paper does not release new assets.

875 Guidelines:

- 876 • The answer NA means that the paper does not release new assets.
- 877 • Researchers should communicate the details of the dataset/code/model as part of their
878 submissions via structured templates. This includes details about training, license,
879 limitations, etc.
- 880 • The paper should discuss whether and how consent was obtained from people whose
881 asset is used.
- 882 • At submission time, remember to anonymize your assets (if applicable). You can either
883 create an anonymized URL or include an anonymized zip file.

884 14. Crowdsourcing and Research with Human Subjects

885 Question: For crowdsourcing experiments and research with human subjects, does the paper
886 include the full text of instructions given to participants and screenshots, if applicable, as
887 well as details about compensation (if any)?

888 Answer: [NA]

889 Justification: This paper does not involve crowdsourcing nor research with human subjects.

890 Guidelines:

- 891 • The answer NA means that the paper does not involve crowdsourcing nor research with
892 human subjects.
- 893 • Including this information in the supplemental material is fine, but if the main contribu-
894 tion of the paper involves human subjects, then as much detail as possible should be
895 included in the main paper.
- 896 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
897 or other labor should be paid at least the minimum wage in the country of the data
898 collector.

899 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 900 Subjects

901 Question: Does the paper describe potential risks incurred by study participants, whether
902 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
903 approvals (or an equivalent approval/review based on the requirements of your country or
904 institution) were obtained?

905 Answer: [NA]

906 Justification: This paper does not involve crowdsourcing nor research with human subjects.

907 Guidelines:

- 908 • The answer NA means that the paper does not involve crowdsourcing nor research with
909 human subjects.
- 910 • Depending on the country in which research is conducted, IRB approval (or equivalent)
911 may be required for any human subjects research. If you obtained IRB approval, you
912 should clearly state this in the paper.
- 913 • We recognize that the procedures for this may vary significantly between institutions
914 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
915 guidelines for their institution.
- 916 • For initial submissions, do not include any information that would break anonymity (if
917 applicable), such as the institution conducting the review.