



HAL
open science

Adaptive extra-gradient methods for min-max optimization and games

Kimon Antonakopoulos, Veronica E Belmega, Panayotis Mertikopoulos

► **To cite this version:**

Kimon Antonakopoulos, Veronica E Belmega, Panayotis Mertikopoulos. Adaptive extra-gradient methods for min-max optimization and games. ICLR 2021 - 9th International Conference on Learning Representations, May 2021, Virtual, Unknown Region. pp.1-28. hal-03342601

HAL Id: hal-03342601

<https://hal.inria.fr/hal-03342601>

Submitted on 13 Sep 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

ADAPTIVE EXTRA-GRADIENT METHODS FOR MIN-MAX OPTIMIZATION AND GAMES

Kimón Antonakopoulos

Univ. Grenoble Alpes, CNRS, Inria, Grenoble INP
LIG, 38000 Grenoble, France
kimon.antonakopoulos@inria.fr

E. Veronica Belmega

ETIS/ENSEA
Univ. de Cergy-Pontoise-CNRS, France
belmega@ensea.fr

Panayotis Mertikopoulos

Univ. Grenoble Alpes, CNRS, Inria, Grenoble INP, LIG, 38000 Grenoble, France &
Criteo AI Lab
panayotis.mertikopoulos@imag.fr

ABSTRACT

We present a new family of min-max optimization algorithms that automatically exploit the geometry of the gradient data observed at earlier iterations to perform more informative extra-gradient steps in later ones. Thanks to this adaptation mechanism, the proposed method automatically detects whether the problem is smooth or not, without requiring any prior tuning by the optimizer. As a result, the algorithm simultaneously achieves order-optimal convergence rates, i.e., it converges to an ε -optimal solution within $\mathcal{O}(1/\varepsilon)$ iterations in smooth problems, and within $\mathcal{O}(1/\varepsilon^2)$ iterations in non-smooth ones. Importantly, these guarantees do not require any of the standard boundedness or Lipschitz continuity conditions that are typically assumed in the literature; in particular, they apply even to problems with singularities (such as resource allocation problems and the like). This adaptation is achieved through the use of a geometric apparatus based on Finsler metrics and a suitably chosen mirror-prox template that allows us to derive sharp convergence rates for the methods at hand.

1 INTRODUCTION

The surge of recent breakthroughs in generative adversarial networks (GANs) [20], robust reinforcement learning [43], and other adversarial learning models [28] has sparked renewed interest in the theory of min-max optimization problems and games. In this broad setting, it has become empirically clear that, *ceteris paribus*, the simultaneous training of two (or more) antagonistic models faces drastically new challenges relative to the training of a single one. Perhaps the most prominent of these challenges is the appearance of cycles and recurrent (or even chaotic) behavior in min-max games. This has been studied extensively in the context of learning in bilinear games, in both continuous [16, 32, 42] and discrete time [12, 18, 19, 33], and the methods proposed to overcome recurrence typically focus on mitigating the rotational component of min-max games. The method with the richest history in this context is the *extra-gradient* (EG) algorithm of Korpelevich [25] and its variants. The EG algorithm exploits the Lipschitz smoothness of the problem and, if coupled with a Polyak–Ruppert averaging scheme, it achieves an $\mathcal{O}(1/T)$ rate of convergence in smooth, convex-concave min-max problems [36]. This rate is known to be tight [35, 41] but, in order to achieve it, the original method requires the problem’s Lipschitz constant to be known in advance. If the problem is not Lipschitz smooth (or the algorithm is run with a vanishing step-size schedule), the method’s rate of convergence drops to $\mathcal{O}(1/\sqrt{T})$.

Our contributions. Our aim in this paper is to provide an algorithm that automatically adapts to smooth / non-smooth min-max problems and games, and achieves order-optimal rates in both classes without requiring any prior tuning by the optimizer. In this regard, we propose a flexible algorithmic scheme, which we call AdaProx, and which exploits gradient data observed at earlier iterations to perform more informative extra-gradient steps in later ones. Thanks to this mechanism, and to the best of our knowledge, AdaProx is the first algorithm that simultaneously achieves the following:

	EG [24, 25, 36]	GRAAL [30]	GMP [50]	AMP [1, 17]	BL [2]	AdaProx [ours]
PARAMETER-AGNOSTIC	✗	✓	PARTIAL	✓	PARTIAL	✓
RATE INTERPOLATION	✗	✗	✓	✗	✓	✓
UNBOUNDED DOMAIN	✗	✓	✗	✗	✗	✓
SINGULARITIES	✗	✗	✗	✓	✗	✓

Table 1: Overview of related work. For the purposes of this table, “parameter-agnostic” means that the method does not require prior knowledge of the parameters of the problem it was designed to solve (Lipschitz modulus, domain diameter, etc.); “rate interpolation” means that the algorithm’s convergence rate is $\mathcal{O}(1/T)$ or $\mathcal{O}(1/\sqrt{T})$ in smooth/non-smooth problems respectively; “unbounded domain” is self-explanatory; and, finally, “singularities” means that the problem’s defining vector field may blow up at a boundary point of the problem’s domain.

1. An $\mathcal{O}(1/\sqrt{T})$ convergence rate in non-smooth problems and $\mathcal{O}(1/T)$ in smooth ones.
2. Applicability to min-max problems and games where the standard boundedness / Lipschitz continuity conditions required in the literature do not hold.
3. Convergence without prior knowledge of the problem’s parameters (e.g., whether the problem’s defining vector field is smooth or not, its smoothness modulus if it is, etc.).

Our proposed method achieves the above by fusing the following ingredients: *a*) a family of local norms – a *Finsler metric* – capturing any singularities in the problem at hand; *b*) a suitable mirror-prox template; and *c*) an adaptive step-size policy in the spirit of Rakhlin & Sridharan [46]. We also show that, under a suitable coherence assumption, the sequence of iterates generated by the algorithm converges, thus providing an appealing alternative to iterate averaging in cases where the method’s “last iterate” is more appropriate (for instance, if using AdaProx to solve non-monotone problems).

Related works. There have been several works improving on the guarantees of the original extra-gradient/mirror-prox template. We review the most relevant of these works below; for convenience, we also tabulate these contributions in Table 1 above. Because many of these works appear in the literature on variational inequalities [15], we also use this language in the sequel. In unconstrained problems with an operator that is locally Lipschitz continuous (but not necessarily globally so), the *golden ratio algorithm* (GRAAL) [30] achieves convergence without requiring prior knowledge of the problem’s Lipschitz parameter. However, GRAAL provides no rate guarantees for non-smooth problems – and hence, a fortiori, no interpolation guarantees either. By contrast, such guarantees are provided in problems with a bounded domain by the *generalized mirror-prox* (GMP) algorithm of [50] under the umbrella of Hölder continuity. Still, nothing is known about the convergence of GRAAL/GMP in problems with singularities (i.e., when the problem’s defining vector field blows up at a boundary point of the problem’s domain). Singularities of this type were treated in a recent series of papers [1, 17, 51] by means of a “Bregman continuity” or “Lipschitz-like” condition. These methods are order-optimal in the smooth case, without requiring any knowledge of the problem’s smoothness modulus. On the other hand, like GRAAL (but unlike GMP), they do not provide any rate interpolation guarantees between smooth and non-smooth problems. Another method that simultaneously achieves an $\mathcal{O}(1/\sqrt{T})$ rate in non-smooth problems and an $\mathcal{O}(1/T)$ rate in smooth ones is the recent algorithm of Bach & Levy [2]. The BL algorithm employs an adaptive, AdaGrad-like step-size policy which allows the method to interpolate between the two regimes – and this, even with noisy gradient feedback. On the negative side, the BL algorithm requires a bounded domain with a (Bregman) diameter that is known in advance; as a result, its theoretical guarantees do not apply to unbounded problems. In addition, the BL algorithm makes crucial use of boundedness and Lipschitz continuity; extending the BL method beyond this standard framework is a highly non-trivial endeavor which formed a big part of this paper’s motivation.

2 PROBLEM SETUP AND BLANKET ASSUMPTIONS

We begin in this section by reviewing some basics for min-max problems and games.

2.1. Min-max / Saddle-point problems. A *min-max game* is a saddle-point problem of the form

$$\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathcal{L}(\theta, \phi) \tag{SP}$$

where Θ, Φ are convex subsets of some ambient real space and $\mathcal{L}: \Theta \times \Phi \rightarrow \mathbb{R}$ is the problem’s *loss function*. In the game-theoretic interpretation of (SP), the player controlling θ seeks to minimize $\mathcal{L}(\theta, \phi)$ for any value of the maximization variable ϕ , while the player controlling ϕ seeks to maximize $\mathcal{L}(\theta, \phi)$ for any value of the minimization variable θ . Accordingly, solving (SP) consists of finding a *Nash equilibrium* (NE), i.e., an action profile $(\theta^*, \phi^*) \in \Theta \times \Phi$ such that

$$\mathcal{L}(\theta^*, \phi) \leq \mathcal{L}(\theta^*, \phi^*) \leq \mathcal{L}(\theta, \phi^*) \quad \text{for all } \theta \in \Theta, \phi \in \Phi. \quad (1)$$

By the minimax theorem of von Neumann [52], Nash equilibria are guaranteed to exist when Θ, Φ are compact and \mathcal{L} is convex-concave (i.e., convex in θ and concave in ϕ). Much of our paper is motivated by the question of calculating a Nash equilibrium (θ^*, ϕ^*) of (SP) in the context of von Neumann’s theorem; we expand on this below.

2.2. Games. Going beyond the min-max setting, a *continuous game in normal form* is defined as follows: First, consider a finite set of players $\mathcal{N} = \{1, \dots, N\}$, each with their own action space $\mathcal{K}_i \in \mathbb{R}^{d_i}$ (assumed convex but possibly not closed). During play, each player selects an action x_i from \mathcal{K}_i with the aim of minimizing a loss determined by the ensemble $x := (x_i; x_{-i}) := (x_1, \dots, x_N)$ of all players’ actions. In more detail, writing $\mathcal{K} := \prod_i \mathcal{K}_i$ for the game’s total action space, we assume that the loss incurred by the i -th player is $\ell_i(x_i; x_{-i})$, where $\ell_i: \mathcal{K} \rightarrow \mathbb{R}$ is the player’s *loss function*. In this context, a Nash equilibrium is any action profile $x^* \in \mathcal{K}$ that is *unilaterally stable*, i.e.,

$$\ell_i(x_i^*; x_{-i}^*) \leq \ell_i(x_i; x_{-i}^*) \quad \text{for all } x_i \in \mathcal{K}_i \text{ and all } i \in \mathcal{N}. \quad (\text{NE})$$

If each \mathcal{K}_i is compact and ℓ_i is convex in x_i , existence of Nash equilibria is guaranteed by the theorem of Debreu [13]. Given that a min-max problem can be seen as a two-player zero-sum game with $\ell_1 = \mathcal{L}, \ell_2 = -\mathcal{L}$, von Neumann’s theorem may in turn be seen as a special case of Debreu’s; in the sequel, we describe a first-order characterization of Nash equilibria that encapsulates both. In most cases of interest, the players’ loss functions are *individually subdifferentiable* on a subset \mathcal{X} of \mathcal{K} with $\text{ri } \mathcal{K} \subseteq \mathcal{X} \subseteq \mathcal{K}$ [21, 47]. This means that there exists a (possibly discontinuous) vector field $V_i: \mathcal{X} \rightarrow \mathbb{R}^{d_i}$ such that

$$\ell_i(x_i'; x_{-i}) \geq \ell_i(x_i; x_{-i}) + \langle V_i(x), x_i' - x_i \rangle \quad (2)$$

for all $x \in \mathcal{X}, x' \in \mathcal{K}$ and all $i \in \mathcal{N}$ [21]. In the simplest case, if ℓ_i is differentiable at x , then $V_i(x)$ can be interpreted as the gradient of ℓ_i with respect to x_i . The *raison d’être* of the more general definition (2) is that it allows us to treat non-smooth loss functions that are common in machine learning (such as L^1 -regularized losses). We make this distinction precise below:

1. If there is no continuous vector field $V_i(x)$ satisfying (2), the game is called *non-smooth*.
2. If there is a continuous vector field $V_i(x)$ satisfying (2), the game is called *smooth*.

Remark. We stress here that the adjective “smooth” refers to the game itself: for instance, if $\ell(x) = |x|$ for $x \in \mathbb{R}$, the game is not smooth and any V satisfying (2) is discontinuous at 0. In this regard, the above boils down to whether the (individual) subdifferential of each ℓ_i admits a continuous selection.

2.3. Resource allocation and equilibrium problems. The notion of a Nash equilibrium captures the unilateral minimization of the players’ individual loss functions. In many practical cases of interest, a notion of equilibrium is still relevant, even though it is not necessarily attached to the minimization of individual loss functions. Such problems are known as “equilibrium problems” [15, 26]; to avoid unnecessary generalities, we focus here on a relevant problem that arises in distributed computing architectures (such as GPU clusters and the like). To state the problem, consider a distributed computing grid consisting of N parallel processors that serve demands arriving at a rate of ρ per unit of time (measured e.g., in flop/s). If the maximum processing rate of the i -th node is μ_i (without overclocking), and jobs are buffered and served on a first-come, first-served (FCFS) basis, the mean time required to process a unit demand at the i -th node is given by the Kleinrock M/M/1 response function $\tau_i(x_i) = 1/(\mu_i - x_i)$, where x_i denotes the node’s *load* [5]. Accordingly, the set of feasible loads that can be processed by the grid is $\mathcal{X} := \{(x_1, \dots, x_N) : 0 \leq x_i < \mu_i, x_1 + \dots + x_N = \rho\}$. In this context, a load profile $x^* \in \mathcal{X}$ is said to be *balanced* if no infinitesimal process can be better served by buffering it at a different node [40]; formally, this amounts to the so-called *Wardrop equilibrium condition*

$$\tau_i(x_i^*) \leq \tau_j(x_j^*) \quad \text{for all } i, j \in \mathcal{N} \text{ with } x_i^* > 0. \quad (\text{WE})$$

We note here a crucial difference between (WE) and (NE): if we view the grid’s computing nodes as “players”, the constraint $\sum_i x_i = \rho$ means that there is no allowable unilateral deviation $(x_i^*; x_{-i}^*) \mapsto (x_i; x_{-i}^*)$ with $x_i \neq x_i^*$. As a result, (NE) is meaningless as a requirement for this equilibrium problem.

As we discuss below, this resource allocation problem will require the full capacity of our framework.

2.4. Variational inequalities. Importantly, all of the above problems can be restated as a *variational inequality* of the form

$$\text{Find } x^* \in \mathcal{X} \text{ such that } \langle V(x^*), x - x^* \rangle \geq 0 \text{ for all } x \in \mathcal{X}. \quad (\text{VI})$$

In the above, \mathcal{X} is a convex subset of \mathbb{R}^d (not necessarily closed) that represents the problem’s *domain*. The problem’s *defining vector field* $V: \mathcal{X} \rightarrow \mathbb{R}^d$ is then given as follows: In min-max problems and games, V is any field satisfying (2); otherwise, in equilibrium problems of the form (WE), the components of V are $V_i = \tau_i$ (we leave the details of this verification to the reader). This equivalent formulation is quite common in the literature on min-max / equilibrium problems [14, 15, 26, 31], and it is often referred to as the “vector field formulation” [3, 8, 23]. Its usefulness lies in that it allows us to abstract away from the underlying game-theoretic complications (multiple indices, individual subdifferentials, etc.) and provides a unifying framework for a wide range of problems in machine learning, signal processing, operations research, and many other fields [15, 48]. For this reason, our analysis will focus almost exclusively on solving (VI), and we will treat V and $\mathcal{X} \subseteq \mathbb{R}^d$, $d = \sum_i d_i$, as the problem’s primitive data.

2.5. Merit functions and monotonicity. A widely used assumption in the literature on equilibrium problems and variational inequalities is the *monotonicity condition*

$$\langle V(x) - V(x'), x - x' \rangle \geq 0 \text{ for all } x, x' \in \mathcal{X}. \quad (\text{Mon})$$

In single-player games, monotonicity is equivalent to convexity of the optimizer’s loss function; in min-max games, it is equivalent to \mathcal{L} being convex-concave [26]; etc. In the absence of monotonicity, approximating an equilibrium is PPAD-hard [11], so we will state most of our results under (Mon).

Now, to assess the quality of a candidate solution $\hat{x} \in \mathcal{X}$, we will employ the *restricted merit function*

$$\text{Gap}_{\mathcal{C}}(\hat{x}) = \sup_{x \in \mathcal{C}} \langle V(x), \hat{x} - x \rangle, \quad (3)$$

where the “*test domain*” \mathcal{C} is a nonempty convex subset of \mathcal{X} [15, 24, 38]. The motivation for this is provided by the following proposition:

Proposition 1. *Let \mathcal{C} be a nonempty convex subset of \mathcal{X} . Then: a) $\text{Gap}_{\mathcal{C}}(\hat{x}) \geq 0$ whenever $\hat{x} \in \mathcal{C}$; and b) if $\text{Gap}_{\mathcal{C}}(\hat{x}) = 0$ and \mathcal{C} contains a neighborhood of \hat{x} , then \hat{x} is a solution of (VI).*

Proposition 1 generalizes an earlier characterization by Nesterov [38] and justifies the use of $\text{Gap}_{\mathcal{C}}(x)$ as a merit function for (VI); to streamline our presentation, we defer the proof to the paper’s supplement. Moreover, to avoid trivialities, we will also assume that the solution set \mathcal{X}^* of (VI) is nonempty and we will reserve the notation x^* for solutions of (VI). Together with monotonicity, this will be our only blanket assumption.

3 THE EXTRA-GRADIENT ALGORITHM AND ITS LIMITS

Perhaps the most widely used solution method for games and variational inequalities (VIs) is the *extra-gradient* (EG) algorithm of Korpelevich [25] and its variants [29, 45, 46]. This algorithm has a rich history in optimization, and it has recently attracted considerable interest in the fields of machine learning and AI, see e.g., [8, 12, 18, 22, 23, 33, 34] and references therein.

In its simplest form, for problems with closed domains, the algorithm proceeds recursively as

$$X_{t+1/2} = \Pi(X_t - \gamma_t V_t), \quad X_{t+1} = \Pi(X_t - \gamma_t V_{t+1/2}), \quad (\text{EG})$$

where $\Pi(x) = \arg \min_{x' \in \mathcal{X}} \|x' - x\|$ is the Euclidean projection on \mathcal{X} , $V_t := V(X_t)$ for $t = 1, 3/2, \dots$, and $\gamma_t > 0$, is the method’s step-size. Then, running (EG) for T iterations, the algorithm returns the “ergodic average”

$$\bar{X}_T = \frac{\sum_{t=1}^T \gamma_t X_{t+1/2}}{\sum_{t=1}^T \gamma_t}. \quad (4)$$

In this setting, the main guarantees for (EG) date back to [36] and can be summarized as follows:

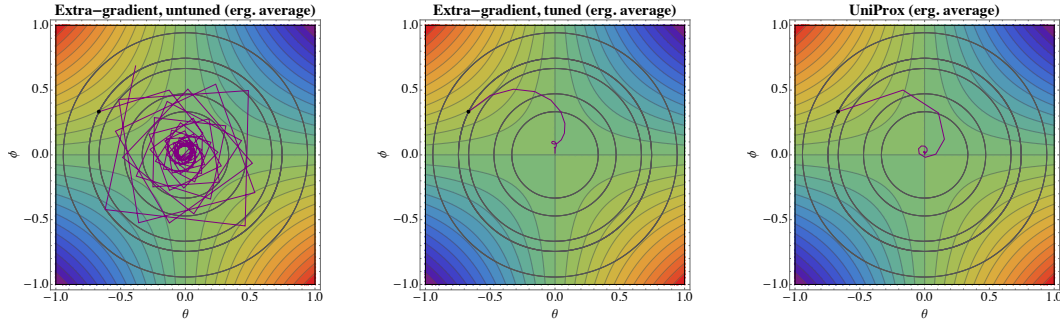


Figure 1: The behavior of (EG) in the bilinear min-max problem $\mathcal{L}(\theta, \phi) = \theta\phi$ with $\theta, \phi \in [-1, 1]$. Given the clipping at $[-1, 1]$, this problem is smooth with $L = 1$; instead, in the unconstrained case, both (BD) and (LC) fail. Still, even in the constrained case, running (EG) with a step-size only slightly above the $1/L$ bound ($L = 1, \gamma = 1.04$) results in a dramatic convergence failure (left plot). Tuning the step-size of (EG) resolves this problem (center), but a constant step-size makes the algorithm unnecessarily conservative towards the end. The proposed AdaProx algorithm automatically exploits previous gradient data to perform more informative extra-gradient steps in later ones, thus achieving faster convergence without tuning.

1. For non-smooth problems (discontinuous V): Assume V is bounded, i.e., there exists some $M > 0$ such that

$$\|V(x)\| \leq M \quad \text{for all } x \in \mathcal{X}. \quad (\text{BD})$$

Then, if (EG) is run with a step-size of the form $\gamma_t \propto 1/\sqrt{t}$, we have

$$\text{Gap}_C(\bar{X}_T) = \mathcal{O}(1/\sqrt{T}). \quad (5)$$

2. For smooth problems (continuous V): Assume V is L -Lipschitz continuous, i.e.,

$$\|V(x) - V(x')\| \leq L\|x - x'\| \quad \text{for all } x, x' \in \mathcal{X}. \quad (\text{LC})$$

Then, if (EG) is run with a constant step-size $\gamma < 1/L$, we have

$$\text{Gap}_C(\bar{X}_T) = \mathcal{O}(1/T). \quad (6)$$

Remark. In the above, $\|\cdot\|$ is tacitly assumed to be the standard Euclidean norm. Non-Euclidean considerations will play a crucial role in the sequel, but they are not necessary for the moment.

Importantly, the distinction between smooth and non-smooth problems cannot be lifted: the bounds (5) and (6) are tight in their respective problem classes and they cannot be improved without further assumptions [35, 41]. Moreover, we should also note the following:

1. The algorithm changes drastically from the non-smooth to the smooth case: non-smoothness requires $\gamma_t \propto 1/\sqrt{t}$, but such a step-size cannot achieve a fast $\mathcal{O}(1/T)$ rate.
2. If (EG) is run with a constant step-size, L must be known in advance; otherwise, running (EG) with an ill-adapted step-size ($\gamma > 1/L$) could lead to non-convergence.

We illustrate this failure of (EG) in Fig. 1. As we discussed in the introduction, our aim in the sequel will be to provide a single, *adaptive* algorithm that simultaneously achieves the following: *a*) an order-optimal $\mathcal{O}(1/\sqrt{T})$ convergence rate in non-smooth problems and $\mathcal{O}(1/T)$ in smooth ones; *b*) convergence in problems where the boundedness / Lipschitz continuity conditions (BD) / (LC) no longer hold; and *c*) achieves all this without prior knowledge of the problem’s parameters.

4 RATE INTERPOLATION: THE EUCLIDEAN CASE

As a prelude to our main result, we provide in this section an adaptive version of (EG) that achieves the “best of both worlds” in the Euclidean setting of Section 3, i.e., an $\mathcal{O}(1/\sqrt{T})$ convergence rate in problems satisfying (BD), and an $\mathcal{O}(1/T)$ rate in problems satisfying (LC). Our starting point is the observation that, if the sequence X_t produced by (EG) converges to a solution of (VI), the difference

$$\delta_t := \|V_{t+1/2} - V_t\| = \|V(X_{t+1/2}) - V(X_t)\| \quad (7)$$

must itself become vanishingly small if V is (Lipschitz) continuous. On the contrary, if V is *discontinuous*, this difference may remain bounded away from zero (consider for example the L^1 loss $\ell(x) = |x|$ near 0). Based on this observation, we consider the adaptive step-size policy:

$$\gamma_{t+1} = 1/\sqrt{1 + \sum_{s=1}^t \delta_s^2}. \quad (8)$$

The intuition behind (8) is as follows: If V is not smooth and $\liminf_{t \rightarrow \infty} \delta_t > 0$, then γ_t will vanish at a $\Theta(1/\sqrt{t})$ rate, which is the optimal step-size schedule for problems satisfying (BD) but not (LC). Instead, if V satisfies (LC) and X_t converges to a solution x^* of (VI), it is plausible to expect that the infinite series $\sum_t \delta_t^2$ is summable, in which case the step-size γ_t will not vanish as $t \rightarrow \infty$. Furthermore, since δ_t is defined in terms of successive gradient differences, it automatically exploits the variation of the gradient data observed up to time t , so it can be expected to adjust to the “local” Lipschitz constant of V around a solution x^* of (VI).

Our step-size policy and motivation are similar in spirit to the “predictable sequence” approach of [46]. For now, we only state (without proof) our main result for problems satisfying (BD) or (LC).

Theorem 1. *Suppose V satisfies (Mon), let \mathcal{C} be a compact neighborhood of a solution of (VI), and let $H = \sup_{x \in \mathcal{C}} \|X_1 - x\|^2$. If (EG) is run with the adaptive step-size policy (8), we have:*

$$a) \text{ If } V \text{ satisfies (BD): } \text{Gap}_{\mathcal{C}}(\bar{X}_T) = \mathcal{O}\left(\frac{H + 4M^3 + \log(1 + 4M^2T)}{\sqrt{T}}\right). \quad (9a)$$

$$b) \text{ If } V \text{ satisfies (LC): } \text{Gap}_{\mathcal{C}}(\bar{X}_T) = \mathcal{O}(H/T). \quad (9b)$$

Theorem 1 (which is proved in the sequel as a special case of **Theorem 2**) should be compared to the corresponding results of **Bach & Levy** [2]. In the non-smooth case, [2] provides a bound of the form $\tilde{\mathcal{O}}(\alpha MD/\sqrt{T})$ with $D^2 = \frac{1}{2} \max_{x \in \mathcal{X}} \|x\|^2 - \frac{1}{2} \min_{x \in \mathcal{X}} \|x\|^2$ (recall that [2] only treats problems with a bounded domain), and $\alpha = \max\{M/M_0, M_0/M\}$ where M_0 is an initial estimate of M . The worst-case value of α is $\mathcal{O}(M)$ when good estimates are not readily available; in this regard, (9a) essentially replaces the $\mathcal{O}(D)$ constant of **Bach & Levy** [2] by $\mathcal{O}(M)$. Since $D = \infty$ in problems with an unbounded domain, **Theorem 1** provides a significant improvement in this regard.

In terms of L , the smooth guarantee of [2] is $\tilde{\mathcal{O}}(\alpha^2 LD^2/T)$, so the multiplicative constant in the bound also becomes infinite in problems with an unbounded domain. In our case, D^2 is replaced by H (which is also finite) times an additional multiplicative constant which is increasing in M and L (but is otherwise asymptotic, so it is not included in the statement of **Theorem 1**). This removes an additional limitation in the results of [2]; in the next sections we drop even the Euclidean regularity requirements (BD)/(LC), and we provide a rate interpolation result that does not require either condition.

5 FINSLER REGULARITY

To motivate our analysis outside the setting of (BD)/(LC), consider the vector field

$$V_i(x) = (\mu_i - x_i)^{-1} + \lambda \mathbb{1}\{x_i > 0\}, \quad i = 1, \dots, N, \quad (10)$$

which corresponds to the distributed computing problem of **Section 2.3** plus a regularization term designed to limit the activation of computing nodes at low loads. Clearly, we have $\|V(x)\| \rightarrow \infty$ whenever $x_i \rightarrow 0^+$, so (BD) and (LC) both fail (the latter even if $\lambda = 0$). On the other hand, if we consider the “local” norm $\|v\|_{x,*} = \sum_{i=1}^d (\mu_i - x_i) |v_i|$, we have $\|V(x)\|_{x,*} \leq d + \lambda \sum_{i=1}^d \mu_i$, so V is *bounded relative to* $\|\cdot\|_{x,*}$. This observation motivates the use of a *local* – as opposed to *global* – norm, which we define formally as follows:

Definition 1. A *Finsler metric* on a convex subset \mathcal{X} of \mathbb{R}^d is a continuous function $F: \mathcal{X} \times \mathbb{R}^d \rightarrow \mathbb{R}_+$ which satisfies the following properties for all $x \in \mathcal{X}$ and all $z, z' \in \mathbb{R}^d$:

1. *Subadditivity:* $F(x; z + z') \leq F(x; z) + F(x; z')$.
2. *Absolute homogeneity:* $F(x; \lambda z) = |\lambda|F(x; z)$ for all $\lambda \in \mathbb{R}$.
3. *Positive-definiteness:* $F(x; z) \geq 0$ with equality if and only if $z = 0$.

Given a Finsler metric on \mathcal{X} , the induced *primal/dual local norms* on \mathcal{X} are respectively defined as

$$\|z\|_x = F(x; z) \quad \text{and} \quad \|v\|_{x,*} = \max\{\langle v, z \rangle : F(x; z) = 1\} \quad (11)$$

for all $x \in \mathcal{X}$ and all $z, v \in \mathbb{R}^d$. We will also say that a Finsler metric on \mathcal{X} is *regular* when $\|v\|_{x',*}/\|v\|_{x,*} = 1 + \mathcal{O}(\|x' - x\|_x)$ for all $x, x' \in \mathcal{X}$, $v \in \mathbb{R}^d$. Finally, for simplicity, we will also assume in the sequel that $\|\cdot\|_x \geq \nu\|\cdot\|$ for some $\nu > 0$ and all $x \in \mathcal{X}$ (this last assumption is for convenience only, as the norm could be redefined to $\|\cdot\|_x \leftarrow \|\cdot\|_x + \nu\|\cdot\|$ without affecting our theoretical analysis).

When \mathcal{X} is equipped with a regular Finsler metric as above, we will say that it is a *Finsler space*.

Example 5.1. Let $F(x; z) = \|z\|$ where $\|\cdot\|$ denotes the reference norm of $\mathcal{X} = \mathbb{R}^d$. Then the properties of [Definition 1](#) are satisfied trivially. ◀

Example 5.2. For a more interesting example of a Finsler structure, consider the set $\mathcal{X} = (0, 1]^d$ and the metric $\|z\|_x = \max_i |z_i|/x_i$, $z \in \mathbb{R}^d$, $x \in \mathcal{X}$. In this case $\|v\|_{x,*} = \sum_{i=1}^d x_i |v_i|$ for all $v \in \mathbb{R}^d$, and the only property of [Definition 1](#) that remains to be proved is that of regularity. To that end, we have

$$\|v\|_{x',*} - \|v\|_{x,*} \leq \sum_{i=1}^d |v_i| \cdot |x'_i - x_i| = \sum_{i=1}^d x_i |v_i| \cdot |x'_i - x_i|/x_i \leq \|v\|_{x,*} \cdot \|x' - x\|_x. \quad (12)$$

Hence, by dividing by $\|v\|_{x,*}$, we readily get $\|v\|_{x',*}/\|v\|_{x,*} \leq 1 + \|x - x'\|_x$ i.e., $\|\cdot\|_x$ is regular in the sense of [Definition 1](#). As we discuss in the sequel, this metric plays an important role for distributed computing problems of the form presented in [Section 2.3](#). ◀

With all this in hand, we will say that a vector field $V: \mathcal{X} \rightarrow \mathbb{R}^d$ is

1. *Metrically bounded* if there exists some $M > 0$ such that

$$\|V(x)\|_{x,*} \leq M \quad \text{for all } x \in \mathcal{X}. \quad (\text{MB})$$

2. *Metrically smooth* if there exists some $L > 0$ such that

$$\|V(x') - V(x)\|_{x,*} \leq L\|x' - x\|_{x'} \quad \text{for all } x', x \in \mathcal{X}. \quad (\text{MS})$$

The notion of metric boundedness/smoothness extends that of ordinary boundedness/Lipschitz continuity to a Finsler context; note also that, even though neither side of [\(MS\)](#) is unilaterally symmetric under the change $x \leftrightarrow x'$, the condition [\(MS\)](#) as a whole *is*. Our next example shows that this extension is *proper*, i.e., [\(BD\)](#)/[\(LC\)](#) may both fail while [\(MB\)](#)/[\(MS\)](#) both hold:

Example 5.3. Consider the change of variables $x_i \rightsquigarrow 1 - x_i/\mu_i$ in the resource allocation problem of [Section 2.3](#). Then, writing $V_i(x) = -(1/x_i) - \lambda \mathbb{1}_{\{x_i < 1\}}$ for the transformed field [\(10\)](#) under this change of variables, we readily get $V_i(x) \rightarrow -\infty$ as $x_i \rightarrow 0^+$; as a result, both [\(BD\)](#) and [\(LC\)](#) fail to hold for *any* global norm on \mathbb{R}^d . Instead, under the *local* norm $\|z\|_x = \max_i |z_i|/x_i$, we have:

1. For all $\lambda \geq 0$, V satisfies [\(MB\)](#) with $M = d(1 + \lambda)$: $\|V(x)\|_{x,*} \leq \sum_{i=1}^d x_i \cdot (1/x_i + \lambda) = d(1 + \lambda)$.
2. For $\lambda = 0$, V satisfies [\(MS\)](#) with $L = d$: indeed, for all $x, x' \in \mathcal{X}$, we have

$$\|V(x') - V(x)\|_{x,*} = \sum_{i=1}^d x_i \left| \frac{1}{x'_i} - \frac{1}{x_i} \right| = \sum_{i=1}^d \frac{|x'_i - x_i|}{x'_i} \leq d \max_i \frac{|x'_i - x_i|}{x'_i} = d\|x' - x\|_{x'}. \quad (13)$$

6 THE ADAPROX ALGORITHM AND ITS GUARANTEES

The method. We are now in a position to define a family of algorithms that is capable of interpolating between the optimal smooth/non-smooth convergence rates for solving [\(VI\)](#) without requiring either [\(BD\)](#) or [\(LC\)](#). To do so, the key steps in our approach will be to (i) equip \mathcal{X} with a suitable Finsler structure (as in [Section 5](#)); and (ii) replace the Euclidean projection in [\(EG\)](#) with a suitable ‘‘Bregman proximal’’ step that is compatible with the chosen Finsler structure on \mathcal{X} . We begin with the latter (assuming that \mathcal{X} is equipped with an arbitrary Finsler structure):

Definition 2. We say that $h: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ is a *Bregman-Finsler function* on \mathcal{X} if:

1. h is convex, lower semi-continuous (l.s.c.), $\text{cl}(\text{dom } h) = \text{cl}(\mathcal{X})$, and $\text{dom } \partial h = \mathcal{X}$.
2. The subdifferential of h admits a *continuous selection* $\nabla h(x) \in \partial h(x)$ for all $x \in \mathcal{X}$.
3. h is *strongly convex*, i.e., there exists some $K > 0$ such that

$$h(x') \geq h(x) + \langle \nabla h(x), x' - x \rangle + \frac{K}{2} \|x' - x\|_x^2 \quad (14)$$

for all $x \in \mathcal{X}$ and all $x' \in \text{dom } h$.

The *Bregman divergence* induced by h is defined for all $x \in \mathcal{X}$, $x' \in \text{dom } h$ as

$$D(x', x) = h(x') - h(x) - \langle \nabla h(x), x' - x \rangle \quad (15)$$

and the associated *prox-mapping* is defined for all $x \in \mathcal{X}$ and $y \in \mathbb{R}^d$ as

$$P_x(y) = \arg \min_{x' \in \mathcal{X}} \{ \langle y, x - x' \rangle + D(x', x) \}. \quad (16)$$

[Definition 2](#) is fairly technical, so some clarifications are in order. First, to connect this definition with the Euclidean setup of [Section 4](#), the prox-mapping (16) should be seen as the Bregman equivalent of a Euclidean projection step, i.e., $\Pi(x+y) \leftrightarrow P_x(y)$. Second, a key difference between [Definition 2](#) and other definitions of Bregman functions in the literature [[4](#), [6](#), [7](#), [9](#), [24](#), [37](#), [38](#), [49](#)] is that h is assumed strongly convex relative to a *local* norm – not a global norm. This “locality” will play a crucial role in allowing the proposed methods to adapt to the geometry of the problem. For concreteness, we provide below an example that expands further on [Examples 5.2](#) and [5.3](#):

Example 6.1. Consider the local norm $\|z\|_x = \max_i |z_i|/x_i$ on $\mathcal{X} = (0, 1]^d$ and let $h(x) = \sum_{i=1}^d 1/x_i$ on $(0, 1]^d$. We then have

$$D(x', x) = \sum_{i=1}^d \left[\frac{1}{x'_i} - \frac{1}{x_i} + \frac{x'_i - x_i}{x_i^2} \right] = \sum_{i=1}^d \frac{(x'_i - x_i)^2}{x_i^2 x'_i} \geq \sum_{i=1}^d (1 - x'_i/x_i)^2 \geq \|x' - x\|_x^2 \quad (17)$$

i.e., h is 1-strongly convex relative to $\|\cdot\|_x$ on \mathcal{X} . \blacktriangleleft

With all this in place, the extra-gradient method can be adapted to our current setting as follows:

$$\begin{aligned} X_{t+1/2} &= P_{X_t}(-\gamma_t V_t) & \delta_t &= \|V_{t+1/2} - V_t\|_{X_{t+1/2},*} \\ X_{t+1} &= P_{X_t}(-\gamma_t V_{t+1/2}) & \gamma_{t+1} &= 1/\sqrt{1 + \sum_{s=1}^t \delta_s^2} \end{aligned} \quad (\text{AdaProx})$$

with $V_t = V(X_t)$, $t = 1, 3/2, \dots$, as in [Section 3](#). In words, this method builds on the template of [\(EG\)](#) by (i) replacing the Euclidean projection with a mirror step; (ii) replacing the global norm in [\(8\)](#) with a dual Finsler norm evaluated at the algorithm’s leading state $X_{t+1/2}$.

Convergence speed. With all this in hand, our main result for AdaProx can be stated as follows:

Theorem 2. *Suppose V satisfies [\(Mon\)](#), let \mathcal{C} be a compact neighborhood of a solution of [\(VI\)](#), and set $H = \sup_{x \in \mathcal{C}} D(x, X_1)$. Then, the AdaProx algorithm enjoys the guarantees:*

$$a) \text{ If } V \text{ satisfies } (\text{MB}): \quad \text{Gap}_{\mathcal{C}}(\bar{X}_T) = \mathcal{O}\left(\frac{H + M^3(1 + 1/K)^2 + \log(1 + 4M^2(1 + 2/K)^2 T)}{\sqrt{T}}\right). \quad (18a)$$

$$b) \text{ If } V \text{ satisfies } (\text{MS}): \quad \text{Gap}_{\mathcal{C}}(\bar{X}_T) = \mathcal{O}(H/T). \quad (18b)$$

For the constants that appear in [Eq. \(18\)](#), we refer the reader to the discussion following [Theorem 1](#). Moreover, we defer the proof of [Theorem 2](#) to the paper’s supplement. We only mention here that its key element is the determination of the asymptotic behavior of the adaptive step-size policy γ_t in the non-smooth and smooth regimes, i.e., under [\(MB\)](#) and [\(MS\)](#) respectively. At a very high level, [\(MB\)](#) guarantees that the difference sequence δ_t is bounded, which implies in turn that $\sum_{t=1}^T \gamma_t = \Omega(\sqrt{T})$ and eventually yields the bound [\(18a\)](#) for the algorithm’s ergodic average \bar{X}_T . On the other hand, if [\(MS\)](#) kicks in, we have the following finer result:

Lemma 1. *Assume V satisfies [\(MS\)](#). Then, a) γ_t decreases monotonically to a strictly positive limit $\gamma_\infty = \lim_{t \rightarrow \infty} \gamma_t > 0$; and b) the sequence δ_t is square summable: in particular, $\sum_{t=1}^{\infty} \delta_t^2 = 1/\gamma_\infty^2 - 1$.*

By means of this lemma (which we prove in the paper’s supplement), it follows that $\sum_{t=1}^T \gamma_t \geq \gamma_\infty T = \Omega(T)$; hence it ultimately follows that AdaProx enjoys an $\mathcal{O}(1/T)$ rate of convergence under [\(MS\)](#).

Trajectory convergence. In complement to [Theorem 2](#), we also provide a trajectory convergence result that governs the *actual* iterates of the AdaProx algorithm:

Theorem 3. *Suppose that $\langle V(x), x - x^* \rangle < 0$ whenever x^* is a solution of [\(VI\)](#) and x is not. If, in addition, V satisfies [\(MB\)](#) or [\(MS\)](#), the iterates X_t of AdaProx converge to a solution of [\(VI\)](#).*

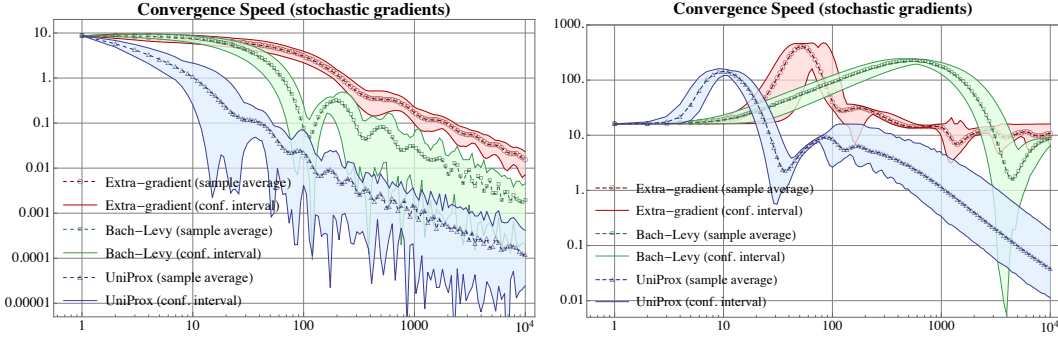


Figure 2: Numerical comparison between the extra-gradient (EG), Bach-Levy (BL) and UniProx algorithms (red circles, green squares and blue triangles respectively). The figure on the left shows the methods’ convergence in a 100×100 bilinear game; the one on the right shows the methods’ convergence in a non-convex/non-concave covariance learning problem. In both cases, the parameters of the EG and BL algorithms have been tuned with a grid search (UniProx has no parameters to tune). All curves have been averaged over $S = 100$ sample runs, and the 95% confidence interval is indicated by the shaded area.

The importance of this result is that, in many practical applications (especially in non-monotone problems), it is more common to harvest the “last iterate” of the method (X_t) rather than its ergodic average (\bar{X}_T); as such, [Theorem 3](#) provides a certain justification for this design choice.

The proof of [Theorem 3](#) relies on non-standard arguments, so we relegate it to the supplement. Structurally, the first step is to show that X_t visits any neighborhood of a solution point $x^* \in \mathcal{X}^*$ infinitely often (this is where the coherence assumption $\langle V(x), x - x^* \rangle$ is used). The second is to use this trapping property in conjunction with a suitable “energy inequality” to establish convergence via the use of a quasi-Fejér technique as in [\[10\]](#); this part is detailed in a separate appendix.

7 NUMERICAL EXPERIMENTS

We conclude in this section with a numerical illustration of the convergence properties of UniProx in two different settings: *a*) bilinear min-max games; and *b*) a simple Wasserstein GAN in the spirit of Daskalakis et al. [\[12\]](#) with the aim of learning an unknown covariance matrix.

Bilinear min-max games. For our first set of experiments, we consider a min-max game of the form $\mathcal{L}(\theta, \phi) = (\theta - \theta^*)^\top A(\phi - \phi^*)$ with $\theta, \phi \in \mathbb{R}^{100}$ and $A \in \mathbb{R}^{100} \times \mathbb{R}^{100}$ (drawn i.i.d. component-wise from a standard Gaussian). To test the convergence of UniProx beyond the “full gradient” framework, we ran the algorithm with stochastic gradient signals of the form $V_t = V(X_t) + U_t$ where U_t is drawn i.i.d. from a centered Gaussian distribution with unit covariance matrix. We then plotted in [Fig. 2](#) the squared gradient norm $\|V(\bar{X}_T)\|^2$ of the method’s ergodic average \bar{X}_T after T iterations (so values closer to zero are better). For benchmarking purposes, we also ran the extra-gradient (EG) and Bach-Levy (BL) algorithms [\[2\]](#) with the same random seed for the simulated gradient noise. The step-size parameter of the EG algorithm was chosen as $\gamma_t = 0.025/\sqrt{t}$, whereas the BL algorithm was run with diameter and gradient bound estimation parameters $D_0 = .5$ and $M_0 = 2.5$ respectively (both determined after a hyper-parameter search since the only *theoretically* allowable values are $D_0 = M_0 = \infty$; interestingly, very large values for D_0 and M_0 did not yield good results). The experiment was repeated $S = 100$ times, and UniProx gave consistently faster rates.

Covariance matrix learning. Going a step further, consider the covariance learning game

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{x \sim \mathcal{N}(0, \Sigma)}[x^\top \theta x] - \mathbb{E}_{z \sim \mathcal{N}(0, I)}[z^\top \theta^\top \phi z], \quad \theta, \phi \in \mathbb{R}^d \times \mathbb{R}^d. \quad (19)$$

The goal here is to generate data drawn from a centered Gaussian distribution with unknown covariance Σ ; in particular, this model follows the Wasserstein GAN formulation of Daskalakis et al. [\[12\]](#) with generator and discriminator respectively given by $G(z) = \theta z$ and $D(x) = x^\top \phi x$ (no clipping). For the experiments, we took $d = 100$, a mini-batch of $m = 128$ samples per update, and we ran the EG, BL and UniProx algorithms as above, tracing the square norm of V as a measure of convergence. Since the problem is non-monotone, there are several disjoint equilibrium components so the algorithms’ behavior is considerably more erratic; however, after this initial warm-up phase, UniProx again gave the faster convergence rates.

ACKNOWLEDGMENTS

This research was partially supported by the COST Action CA16228 “European Network for Game Theory” (GAMENET) and the French National Research Agency (ANR) in the framework of the grants ORACLESS (ANR-16-CE33-0004-01) and ELIOT (ANR-18-CE40-0030 and FAPESP 2018/12579-7), the “Investissements d’avenir” program (ANR-15-IDEX-02), the LabEx PERSYVAL (ANR-11-LABX-0025-01), and MIAI@Grenoble Alpes (ANR-19-P3IA-0003).

A PROPERTIES OF THE RESTRICTED GAP FUNCTION

In this appendix, we discuss the basic properties of the restricted merit function $\text{Gap}_{\mathcal{C}}$ introduced in (3). For completeness, we provide the proof of Proposition 1, which itself is an extension of a similar result by Nesterov [38]:

Proof of Proposition 1. Let $x^* \in \mathcal{X}$ be a solution of (VI) so $\langle V(x^*), x - x^* \rangle \geq 0$ for all $x \in \mathcal{X}$. Then, by monotonicity, we get:

$$\begin{aligned} \langle V(x), x^* - x \rangle &\leq \langle V(x) - V(x^*), x^* - x \rangle + \langle V(x^*), x^* - x \rangle \\ &= -\langle V(x^*) - V(x), x^* - x \rangle - \langle V(x^*), x - x^* \rangle \leq 0, \end{aligned} \quad (\text{A.1})$$

so $\text{Gap}_{\mathcal{C}}(x^*) \leq 0$. On the other hand, if $x^* \in \mathcal{C}$, we also get $\text{Gap}(x^*) \geq \langle V(x^*), x^* - x^* \rangle = 0$, so we conclude that $\text{Gap}_{\mathcal{C}}(x^*) = 0$.

For the converse statement, assume that $\text{Gap}_{\mathcal{C}}(\hat{x}) = 0$ for some $\hat{x} \in \mathcal{C}$ and suppose that \mathcal{C} contains a neighborhood of \hat{x} in \mathcal{X} . First, we claim that the following inequality holds:

$$\langle V(x), x - \hat{x} \rangle \geq 0 \quad \text{for all } x \in \mathcal{C}. \quad (\text{A.2})$$

Indeed, assume to the contrary that there exists some $x_1 \in \mathcal{C}$ such that

$$\langle V(x_1), x_1 - \hat{x} \rangle < 0. \quad (\text{A.3})$$

This would then give

$$0 = \text{Gap}_{\mathcal{C}}(\hat{x}) \geq \langle V(x_1), \hat{x} - x_1 \rangle > 0, \quad (\text{A.4})$$

which is a contradiction. Now, we further claim that \hat{x} is a solution of (VI), i.e.,:

$$\langle V(\hat{x}), x - \hat{x} \rangle \geq 0 \quad \text{for all } x \in \mathcal{X}. \quad (\text{A.5})$$

If we suppose that there exists some $z_1 \in \mathcal{X}$ such that $\langle V(\hat{x}), z_1 - \hat{x} \rangle < 0$, then, by the continuity of V , there exists a neighborhood U' of \hat{x} in \mathcal{X} such that

$$\langle V(x), z_1 - x \rangle < 0 \quad \text{for all } x \in U'. \quad (\text{A.6})$$

Hence, assuming without loss of generality that $U' \subset U \subset \mathcal{C}$ (the latter assumption due to the assumption that \mathcal{C} contains a neighborhood of \hat{x}), and taking $\lambda > 0$ sufficiently small so that $x = \hat{x} + \lambda(z_1 - \hat{x}) \in U'$, we get that $\langle V(x), x - \hat{x} \rangle = \lambda \langle V(x), z_1 - \hat{x} \rangle < 0$, in contradiction to (A.2). We conclude that \hat{x} is a solution of (VI), as claimed. \square

B PROPERTIES OF BREGMAN FUNCTIONS AND PROXIMAL MAPPINGS

In this appendix, we present some basic facts about Bregman-Finsler functions and proximal mappings. Similar results exist in the literature in different contexts (see e.g., [24, 38, 39] and references therein), but given that many of our results rely on the use of *local* – as opposed to *global* – norms, we provide here complete statements and proofs. We then have the following basic lemma connecting the above notions:

Lemma B.1. *Let h be a Bregman-Finsler function on \mathcal{X} . Then, for all $p \in \text{dom } h$, $x \in \text{dom } \partial h$ and all $y \in \partial h(x)$, we have*

$$\langle \nabla h(x), x - p \rangle \leq \langle y, x - p \rangle. \quad (\text{B.1})$$

Proof. By a simple continuity argument, it is sufficient to show that the inequality holds for the relative interior $\text{ri } \mathcal{X}$ of \mathcal{X} . In order to show this, pick a base point $p \in \text{ri } \mathcal{X}$, and let

$$\phi(t) = h(x + t(p - x)) - [h(x) + \langle y, t(p - x) \rangle] \quad \text{for all } t \in [0, 1]. \quad (\text{B.2})$$

Since, h is strongly convex and $y \in \partial h(x)$ due to the first equivalence, it follows that $\phi(t) \geq 0$ with equality if and only if $t = 0$. Since, $\psi(t) = \langle \nabla h(x + t(p - x)) - y, p - x \rangle$ is a continuous selection of subgradients of ϕ and both ϕ and ψ are continuous over $[0, 1]$, it follows that ϕ is continuously differentiable with $\phi' = \psi$ on $[0, 1]$. Hence, with ϕ convex and $\phi(t) \geq 0 = \phi(0)$ for all $t \in [0, 1]$, we conclude that $\phi'(0) = \langle \nabla h(x) - y, p - x \rangle \geq 0$ and thus we obtain the result. \square

The basic ingredient for establishing connections in the Bregman framework is a generalization of the rule of cosines which is known in the literature as the ‘‘three-point identity’’ [9] and will be the main tool for deriving the main estimations for our analysis. Being more precise, we have the following lemma:

Lemma B.2. *Let h be a Bregman-Finsler function on \mathcal{X} . Then, for all $p \in \mathcal{X}$ and all $x, x' \in \mathcal{X}$, we have:*

$$D(p, x') = D(p, x) + D(x, x') + \langle \nabla h(x') - \nabla h(x), x - p \rangle \quad (\text{B.3})$$

The proof of this lemma follows as in the classic Bregman case [9] so we omit it and proceed to derive some key bounds for the Bregman divergence before and after a mirror step:

Proposition B.1. *Let h be a Bregman-Finsler function with strong convexity modulus $K > 0$. Fix some $p \in \mathcal{X}$ and let $x^+ = P_x(v)$ for some $x \in \mathcal{X}^o$ and $v \in \mathbb{R}^d$. We then have:*

$$D(p, x^+) \leq D(p, x) - D(x^+, x) + \langle v, x^+ - p \rangle \quad (\text{B.4})$$

Proof. By the three-point identity established in Lemma B.2, we get:

$$D(p, x) = D(p, x^+) + D(x^+, x) + \langle \nabla h(x) - \nabla h(x^+), x^+ - p \rangle \quad (\text{B.5})$$

By rearranging the terms we get:

$$D(p, x^+) = D(p, x) - D(x^+, x) + \langle \nabla h(x^+) - \nabla h(x), x^+ - p \rangle \quad (\text{B.6})$$

Due to (B.1) and the fact that $x^+ = P_x(v)$ so $\nabla h(x) + v \in \partial h(x^+)$, we get the result. \square

Thanks to the above estimations, we obtain the following inequalities relating the Bregman divergence between two prox-steps:

Proposition B.2. *Let h be a Bregman-Finsler function. Letting $x_1^+ = P_x(v_1)$ and $x_2^+ = P_x(v_2)$, we have:*

$$D(p, x_2^+) \leq D(p, x) + \langle v_2, x_1^+ - p \rangle + [\langle v_2, x_2^+ - x_1^+ \rangle - D(x_2^+, x)] \quad (\text{B.7a})$$

$$\leq D(p, x) + \langle v_2, x_1^+ - p \rangle + \langle v_2 - v_1, x_2^+ - x_1^+ \rangle - D(x_2^+, x_1^+) - D(x_1^+, x). \quad (\text{B.7b})$$

Proof. For the first inequality, by applying Proposition B.1 for $x_2^+ = P_x(v_2)$, we get:

$$\begin{aligned} D(p, x_2^+) &\leq D(p, x) - D(x_2^+, x) + \langle v_2, x_2^+ - p \rangle \\ &= D(p, x) + \langle v_2, x_1^+ - p \rangle + [\langle v_2, x_2^+ - x_1^+ \rangle - D(x_2^+, x)] \end{aligned} \quad (\text{B.8})$$

For the second inequality, we need to bound $\langle v_2, x_2^+ - x_1^+ \rangle - D_h(x_2^+, x)$. In particular, applying again Proposition B.1 for $p = x_2^+$, we get:

$$D(x_2^+, x_1^+) \leq D(x_2^+, x) + \langle v_1, x_1^+ - x_2^+ \rangle - D(x_1^+, x) \quad (\text{B.9})$$

and hence:

$$D(x_2^+, x) \geq D(x_2^+, x_1^+) + D(x_1^+, x) - \langle v_1, x_1^+ - x_2^+ \rangle. \quad (\text{B.10})$$

So, combining the above inequalities we get:

$$\langle v_2, x_2^+ - x_1^+ \rangle - D(x_2^+, x) \leq \langle v_2, x_2^+ - x_1^+ \rangle - D(x_2^+, x_1^+) - D(x_1^+, x) - \langle v_1, x_2^+ - x_1^+ \rangle \quad (\text{B.11})$$

and thus we get the second inequality as well. \square

C MAIN BOUNDS AND ENERGY INEQUALITY

In this appendix, we shall provide the bound of the variation of the operators, i.e.,

$$\|V(X_{t+1/2}) - V(X_t)\|_{X_t, *}^2 \quad (\text{C.1})$$

that lies in the core of our analysis. To begin with, we recall that $(\mathcal{X}, \|\cdot\|_x)$ is a regular Finsler space, i.e., $\|v\|_{x, *} / \|v\|_{x', *} = 1 + \mathcal{O}(\|x - x'\|_x)$. However, in what follows we shall assume the more general condition:

$$\|v\|_{x, *} / \|v\|_{x', *} \leq 1 + \beta [\|x - x'\|_x + \|x - x'\|_{x'}] \text{ for some } \beta > 0 \quad (\text{C.2})$$

Remark 1. It is straightforward for one to observe that a regular Finsler space satisfies (C.2) for $\beta = 1$.

Owning this regularity geometrical property for the problem's domain we shall proceed into showing that

$$\|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2}, *}^2 \quad (\text{C.3})$$

is uniformly bounded. More precisely, we have the following lemma.

Lemma C.1. *Suppose that V satisfies (MB). Then, the sequence $\|V(X_{t+1/2}) - V(X_t)\|_{X_t, *}^2$ is bounded. In particular, the following inequality holds:*

$$\|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2}, *}^2 \leq C^2 \quad (\text{C.4})$$

with $C = 2M + \beta \frac{4M}{K}$.

Proof. It suffices to show that: $\|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2}, *}$ is bounded. More precisely, by the triangle inequality we have:

$$\|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2}, *} \leq \|V(X_{t+1/2})\|_{X_{t+1/2}, *} + \|V(X_t)\|_{X_{t+1/2}, *} \quad (\text{C.5})$$

Let us now bound the (RHS) part of (C.5) term by term. In particular, we have:

- For the first term $\|V(X_{t+1/2})\|_{X_{t+1/2}, *}$ we readily get due to (MB):

$$\|V(X_{t+1/2})\|_{X_{t+1/2}, *} \leq M \quad (\text{C.6})$$

- For the second term $\|V(X_t)\|_{X_{t+1/2}, *}$, we have:

$$\begin{aligned} \|V(X_t)\|_{X_{t+1/2}, *} &\leq \|V(X_t)\|_{X_t, *} + \beta [\|X_t - X_{t+1/2}\|_{X_t} + \|X_t - X_{t+1/2}\|_{X_{t+1/2}}] \\ &\leq M + \beta [\|X_t - X_{t+1/2}\|_{X_t} + \|X_t - X_{t+1/2}\|_{X_{t+1/2}}] \end{aligned} \quad (\text{C.7})$$

Therefore, it suffices to show that the quantity $\|X_t - X_{t+1/2}\|_{X_t} + \|X_t - X_{t+1/2}\|_{X_{t+1/2}}$ is bounded from above. Indeed, we have:

$$\begin{aligned} D(X_t, X_{t+1/2}) + D(X_{t+1/2}, X_t) &= \langle \nabla h(X_t) - \nabla h(X_{t+1/2}), X_t - X_{t+1/2} \rangle \\ &\leq \gamma_t \langle V(X_t), X_t - X_{t+1/2} \rangle \\ &\leq M \gamma_t \|X_t - X_{t+1/2}\|_{X_t} \end{aligned}$$

where the last inequality is obtained due to (MB). Moreover, due to (14) we get:

$$\begin{aligned} D(X_t, X_{t+1/2}) + D(X_{t+1/2}, X_t) &\leq \gamma_t M \sqrt{\frac{2}{K} D(X_{t+1/2}, X_t)} \\ &\leq M \sqrt{\frac{2}{K} [D(X_t, X_{t+1/2}) + D(X_{t+1/2}, X_t)]} \end{aligned}$$

which yields

$$D(X_t, X_{t+1/2}) + D(X_t, X_{t+1/2}) \leq \frac{2M^2}{K} \quad (\text{C.8})$$

Hence, due to the local strong convexity (14) of h , we get:

$$\frac{K}{2} \left[\|X_t - X_{t+1/2}\|_{X_t}^2 + \|X_t - X_{t+1/2}\|_{X_{t+1/2}}^2 \right] \leq \frac{2M^2}{K} \quad (\text{C.9})$$

which in turn implies that:

$$\|X_t - X_{t+1/2}\|_{X_t} \leq \frac{2M}{K} \text{ and } \|X_t - X_{t+1/2}\|_{X_{t+1/2}} \leq \frac{2M}{K} \quad (\text{C.10})$$

and so,

$$\|X_t - X_{t+1/2}\|_{X_t} + \|X_t - X_{t+1/2}\|_{X_{t+1/2}} \leq \frac{4M}{K} \quad (\text{C.11})$$

Moreover, by combining (C.7) and (C.11) we get:

$$\|V(X_t)\|_{X_{t+1/2},*} \leq M + \beta \frac{4M}{K} \quad (\text{C.12})$$

Summarizing, (C.5) combined with (C.7) and (C.12) yields:

$$\|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*} \leq 2M + \beta \frac{4M}{K} \quad (\text{C.13})$$

and hence the result follows. \square

We now proceed to prove the energy inequality stated in Lemma C.2.

Lemma C.2. *For all $x \in \mathcal{X}$, the iterates X_t of AdaProx satisfy the recursive bound:*

$$\begin{aligned} D(x, X_{t+1}) \leq & D(x, X_t) - \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - x \rangle + \gamma_t \langle V(X_{t+1/2}) - V(X_t), X_{t+1} - X_{t+1/2} \rangle \\ & - D(X_{t+1}, X_{t+1/2}) - D(X_{t+1/2}, X_t) \end{aligned} \quad (\text{C.14})$$

Proof. The result follows directly by setting $X_1^+ = X_{t+1/2}$, $X_2^+ = X_{t+1}$, $x = X_t$, $v_1 = -\gamma_t V(X_t)$ and $v_2 = -\gamma_t V(X_{t+1/2})$ in Proposition B.2. \square

D RATE INTERPOLATION GUARANTEES

In this appendix, we provide the proof of the the regime-agnostic rate interpolation guarantees of the AdaProx. In order, to provide the necessary the respective rates we shall provide an intermediate result concerning the case of (MS). Formally, we have the following lemma.

Lemma D.1. *Assume V satisfies (MS) and $X_t, X_{t+1/2}$ are the iterates of AdaProx. Then, the following hold:*

1. $\gamma_t \rightarrow \inf_{t \in \mathbb{N}} \gamma_t = \gamma_\infty > 0$
2. *The sequence $\|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2$ is summable. In particular, we have:*

$$\sum_{t=1}^{+\infty} \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 = \frac{1}{\gamma_\infty^2} - 1 \quad (\text{D.1})$$

Proof. Since γ_t is decreasing and bounded from below ($\gamma_t \geq 0$), then we readily obtain that its limit exists and more precisely we have:

$$\lim_{t \rightarrow +\infty} \gamma_t = \inf_{t \in \mathbb{N}} \gamma_t = \gamma_\infty \geq 0 \quad (\text{D.2})$$

Let us now assume that $\gamma_\infty = 0$. Then, by recalling (C.14):

$$\begin{aligned} D(p, X_{t+1}) \leq & D(p, X_t) - \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - p \rangle + \gamma_t \langle V(X_{t+1/2}) - V(X_t), X_{t+1} - X_{t+1/2} \rangle \\ & - D(X_{t+1/2}, X_t) - D(X_{t+1}, X_{t+1/2}) \end{aligned} \quad (\text{D.3})$$

By rearranging the above and telescoping $t = 1, \dots, T$ we get:

$$\begin{aligned} \sum_{t=1}^T \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - p \rangle &\leq D(p, X_1) + \sum_{t=1}^T \gamma_t \langle V(X_{t+1/2}) - V(X_t), X_{t+1} - X_{t+1/2} \rangle \\ &\quad - \sum_{t=1}^T D(X_{t+1/2}, X_t) - \sum_{t=1}^T D(X_{t+1}, X_{t+1/2}) \end{aligned} \quad (\text{D.4})$$

whereas, by applying Fenchel-Young inequality to the above we readily get:

$$\begin{aligned} \sum_{t=1}^T \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - p \rangle &\leq D(p, X_1) + \frac{1}{2K} \sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \\ &\quad + \frac{K}{2} \sum_{t=1}^T \|X_{t+1} - X_{t+1/2}\|_{X_{t+1/2}}^2 - \sum_{t=1}^T D(X_{t+1/2}, X_t) - \sum_{t=1}^T D(X_{t+1}, X_{t+1/2}) \end{aligned} \quad (\text{D.5})$$

and by considering that by (14):

$$\frac{K}{2} \sum_{t=1}^T \|X_{t+1} - X_{t+1/2}\|_{X_{t+1/2}}^2 - \sum_{t=1}^T D(X_{t+1}, X_{t+1/2}) \leq 0 \quad (\text{D.6})$$

we finally obtain:

$$\begin{aligned} \sum_{t=1}^T \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - p \rangle &\leq D(p, X_1) + \frac{1}{2K} \sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \\ &\quad - \sum_{t=1}^T D(X_{t+1/2}, X_t) \end{aligned} \quad (\text{D.7})$$

Therefore, by the definition (MS) we have:

$$\begin{aligned} \sum_{t=1}^T \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - p \rangle &\leq D(p, X_1) + \frac{1}{2K} \sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \\ &\quad - \frac{K}{2L^2} \sum_{t=1}^T \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \end{aligned} \quad (\text{D.8})$$

which becomes:

$$\begin{aligned} \sum_{t=1}^T \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - p \rangle &\leq D(p, X_1) + \sum_{t=1}^T \left[\frac{\gamma_t^2}{2K} - \frac{K}{4L^2} \right] \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \\ &\quad - \frac{K}{4L^2} \sum_{t=1}^T \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \end{aligned} \quad (\text{D.9})$$

Now, by setting $p = x^*$ with x^* being a solution of (VI) and using the fact that $\langle V(X_{t+1/2}), X_{t+1/2} - x^* \rangle \geq 0$ and $D(x^*, X_1) \leq D'$ (by the compatibility of h), we obtain:

$$\frac{K}{4L^2} \sum_{t=1}^T \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \leq D' + \sum_{t=1}^T \left[\frac{\gamma_t^2}{2K} - \frac{K}{4L^2} \right] \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \quad (\text{D.10})$$

Moreover, by observing that the quantity $\left[\frac{\gamma_t^2}{2K} - \frac{K}{4L^2} \right] \leq 0$, whenever $\gamma_t \leq \sqrt{2}K/2L$ and since we assumed that $\gamma_t \rightarrow 0$, there exists some $t_0 \in \mathbb{N}$ such that:

$$\left[\frac{\gamma_t^2}{2K} - \frac{K}{4L^2} \right] \leq 0 \quad \text{for all } t \geq t_0 \quad (\text{D.11})$$

Therefore, (D.10) becomes:

$$\frac{1}{\gamma_{T+1}^2} - 1 = \sum_{t=1}^T \|V(X_{t+1/2}) - V(X_t)\|_{\bar{X}_{t+1/2,*}}^2 \leq D' + \sum_{t=1}^{t_0} \left[\frac{\gamma_t^2}{2K} - \frac{K}{4L^2} \right] \|V(X_{t+1/2}) - V(X_t)\|_{\bar{X}_{t+1/2,*}}^2 \quad (\text{D.12})$$

In addition, since $1/\gamma_{T+1} \rightarrow +\infty$, by the fact that $\gamma_t \rightarrow 0$, this yields that:

$$+\infty \leq D' + \sum_{t=1}^{t_0} \left[\frac{\gamma_t^2}{2K} - \frac{K}{4L^2} \right] \|V(X_{t+1/2}) - V(X_t)\|_{\bar{X}_{t+1/2,*}}^2 \quad (\text{D.13})$$

which is a contradiction. Hence, we get that:

$$\lim_{t \rightarrow +\infty} \gamma_t = \inf_{t \in \mathbb{N}} \gamma_t = \gamma_\infty > 0 \quad (\text{D.14})$$

In order to prove our second claim, we first recall the definition of γ_t :

$$\gamma_t = \frac{1}{\sqrt{1 + \sum_{j=1}^{t-1} \|V(X_{j+1/2}) - V(X_j)\|_{\bar{X}_{j+1/2,*}}^2}} \quad (\text{D.15})$$

whereas by developing and rearranging we have:

$$\sum_{j=1}^{t-1} \|V(X_{j+1/2}) - V(X_j)\|_{\bar{X}_{j+1/2,*}}^2 = \frac{1}{\gamma_t^2} - 1 \quad (\text{D.16})$$

Hence, by taking limits on both sides we get:

$$\sum_{t=1}^{+\infty} \|V(X_{t+1/2}) - V(X_t)\|_{\bar{X}_{t+1/2,*}}^2 = \lim_{t \rightarrow +\infty} \sum_{j=1}^{t-1} \|V(X_{j+1/2}) - V(X_j)\|_{\bar{X}_{j+1/2,*}}^2 = \frac{1}{\gamma_\infty^2} - 1 \quad (\text{D.17})$$

where $0 \leq \frac{1}{\gamma_\infty} - 1 < +\infty$, since $0 < \gamma_\infty \leq 1$ and therefore the result follows. \square

Proof of Theorem 2. By recalling (C.14) we have:

$$D(p, X_{t+1}) \leq D(p, X_t) - \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - p \rangle + \gamma_t \langle V(X_{t+1/2}) - V(X_t), X_{t+1} - X_{t+1/2} \rangle - D(X_{t+1/2}, X_t) - D(X_{t+1}, X_{t+1/2}) \quad (\text{D.18})$$

We start our analysis rearranging (C.14). In particular, by telescoping $t = 1, \dots, T$ we get:

$$\begin{aligned} \sum_{t=1}^T \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - p \rangle &\leq D(p, X_1) + \sum_{t=1}^T \gamma_t \langle V(X_{t+1/2}) - V(X_t), X_{t+1} - X_{t+1/2} \rangle \\ &\quad - \sum_{t=1}^T D(X_{t+1/2}, X_t) - \sum_{t=1}^T D(X_{t+1}, X_{t+1/2}) \end{aligned} \quad (\text{D.19})$$

On the other hand, since V is monotone, we readily get:

$$\gamma_t \langle V(p), X_{t+1/2} - p \rangle \leq \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - p \rangle \quad (\text{D.20})$$

Thus, combining (D.20) and (D.19), dividing by $\sum_{t=1}^T \gamma_t$ and setting $\bar{X}_T = \left[\sum_{t=1}^T \gamma_t \right]^{-1} \sum_{t=1}^T \gamma_t X_{t+1/2}$ we get:

$$\begin{aligned} \langle V(p), \bar{X}_T - p \rangle &\leq \left[\sum_{t=1}^T \gamma_t \right]^{-1} \left(D(p, X_1) + \sum_{t=1}^T \gamma_t \langle V(X_{t+1/2}) - V(X_t), X_{t+1} - X_{t+1/2} \rangle - \sum_{t=1}^T D(X_{t+1/2}, X_t) \right. \\ &\quad \left. - \sum_{t=1}^T D(X_{t+1}, X_{t+1/2}) \right) \end{aligned} \quad (\text{D.21})$$

whereas, by applying Fenchel-Young inequality to the above we readily get:

$$\begin{aligned} \langle V(p), \bar{X}_T - p \rangle \leq & \left[\sum_{t=1}^T \gamma_t \right]^{-1} \left(D(p, X_1) + \frac{1}{2K} \sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 + \frac{K}{2} \sum_{t=1}^T \|X_{t+1} - X_{t+1/2}\|_{X_{t+1/2}}^2 \right. \\ & \left. - \sum_{t=1}^T D(X_{t+1/2}, X_t) - \sum_{t=1}^T D(X_{t+1}, X_{t+1/2}) \right) \quad (\text{D.22}) \end{aligned}$$

Thus, if \mathcal{C} is a compact neighbourhood of the solution set \mathcal{X}^* , considering that by (14):

$$\frac{K}{2} \sum_{t=1}^T \|X_{t+1} - X_{t+1/2}\|_{X_{t+1/2}}^2 - \sum_{t=1}^T D(X_{t+1}, X_{t+1/2}) \leq 0 \quad (\text{D.23})$$

and taking suprema on both sides, yields:

$$\begin{aligned} \text{Gap}_{\mathcal{C}}(\bar{X}_T) \leq & \left[\sum_{t=1}^T \gamma_t \right]^{-1} \left(\sup_{p \in \mathcal{C}} D(p, X_1) + \frac{1}{2K} \sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \right. \\ & \left. - \sum_{t=1}^T D(X_{t+1/2}, X_t) \right) \quad (\text{D.24}) \end{aligned}$$

Case 1: Convergence under (MB). Therefore, in order to determine the convergence speed of \bar{X}_T under (MB), we shall examine the asymptotic behaviour of each term of the nominator on the (RHS) of (D.31). In particular, we have the following:

- For the first term: we readily get by the compactness of \mathcal{C} ,

$$\sup_{p \in \mathcal{C}} D(p, X_1) \leq D' \text{ for some constant } D' > 0. \quad (\text{D.25})$$

by the compatibility of the regularizer h .

- For the second term: $\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2$, we have:

$$\begin{aligned} \sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 &= \sum_{t=1}^T (\gamma_t^2 - \gamma_{t+1}^2) \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \\ &\quad + \sum_{t=1}^T \gamma_{t+1}^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \quad (\text{D.26}) \end{aligned}$$

Since, γ_t is non-increasing and therefore $(\gamma_t^2 - \gamma_{t+1}^2 \geq 0)$, and $\gamma_t \leq 1$ the above becomes:

$$\begin{aligned} \sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 &\leq C^2 + \sum_{t=1}^T \gamma_{t+1} \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \\ &\leq C^2 + \sum_{t=1}^T \frac{\|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2}{1 + \sum_{j=1}^t \|V(X_{t+1/2}) - V(X_j)\|_{X_{t+1/2},*}^2} \\ &\leq C^2 + 1 + \log\left(1 + \sum_{t=1}^T \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2\right) \end{aligned}$$

with the last inequality being obtained by [Lemma F.1](#) which combined with (MB) yields:

$$\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \leq C^2 + 1 + \log(1 + C^2 T) \quad (\text{D.27})$$

Finally, for $\sum_{t=1}^T \gamma_t$, we have the following lower-bound

$$\sum_{t=1}^T \gamma_t = \sum_{t=1}^T \frac{1}{\sqrt{1 + \sum_{j=1}^{t-1} \|V(X_{t+1/2}) - V(X_j)\|_{X_{t+1/2},*}^2}} \geq \sum_{t=1}^T \frac{1}{\sqrt{1 + tC^2}} \quad (\text{D.28})$$

which yields:

$$\sum_{t=1}^T \gamma_t = \Omega(\sqrt{T}) \quad \text{and} \quad \sum_{t=1}^T \gamma_t \rightarrow +\infty \quad (\text{D.29})$$

Now, by combining (D.25), (D.27) and (D.29) we readily get that under (MB) we get that:

$$\text{Gap}_C(\bar{X}_T) = \mathcal{O}(1/\sqrt{T}). \quad (\text{D.30})$$

Case 2: Convergence under (MS). We now suppose that V satisfies (MS) condition. By applying Lemma D.1 along with :

$$\begin{aligned} \sum_{t=1}^T \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - p \rangle &\leq D(p, X_1) + \frac{1}{2K} \sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \\ &\quad - \sum_{t=1}^T D(X_{t+1/2}, X_t) \end{aligned} \quad (\text{D.31})$$

by examining the asymptotic behaviour term by term, we get:

- For the first term $D(x^*, X_1)$, since $x^* \in \text{dom } V = \text{dom } h$ and $X_1 \in \text{dom } \partial h$, we have:

$$D(x^*, X_1) < +\infty \quad (\text{D.32})$$

- For the second term $\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2$ we have:

$$\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \leq \sum_{t=1}^T \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \quad (\text{D.33})$$

and by applying Lemma D.1 we have:

$$\sum_{t=1}^T \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \leq \frac{1}{\gamma_\infty^2} - 1 \quad (\text{D.34})$$

with $\gamma_\infty = \inf_t \gamma_t > 0$.

Finally, by applying Lemma D.1 once more by considering $\gamma_\infty = \inf_{t \in \mathbb{N}} \gamma_t > 0$ we have:

$$\sum_{t=1}^T \gamma_t \geq \gamma_\infty \sum_{t=1}^T 1 = \gamma_\infty T \quad (\text{D.35})$$

which yields:

$$\sum_{t=1}^T \gamma_t = \Omega(T) \quad (\text{D.36})$$

and the result follows. \square

E LAST ITERATE'S CONVERGENCE ANALYSIS

In this appendix, we establish the convergence of the sequence generated by (AdaProx), i.e., its so-called last iterate. In particular, we show that the actual iterates (before averaging) of AdaProx converge towards the solution set \mathcal{X}^* . This result comprises of two parts: first we extract convergent subsequences of $X_t, X_{t+1/2}$ to the said set; then we apply the "trapping" argument described in Section 6.

Lemma E.1. *Suppose that V satisfies (MB) (respectively (MS)) and $X_t, X_{t+1/2}$ are the iterates of AdaProx. Then, the following hold:*

1. $\|X_{t+1/2} - X_t\| \rightarrow 0$ while $t \rightarrow +\infty$
2. $\max\{D(X_{t+1/2}, X_t), D(X_t, X_{t+1/2})\} \leq \frac{2M^2}{K} \gamma_t^2$

Proof. For the proof of the first claim, we shall treat the cases of (MB) and (MS) individually.

Case 1: Under (MB) condition. Since γ_t is decreasing and bounded from below, then we readily obtain that. its limit exists and more precisely:

$$\lim_{t \rightarrow +\infty} \gamma_t = \gamma_\infty \geq 0 \quad (\text{E.1})$$

We shall distinguish two individual cases:

- $\gamma_\infty > 0$: By recalling the definition of the adaptive step-size:

$$\gamma_t = \frac{1}{\sqrt{1 + \sum_{j=1}^{t-1} \|V(X_{j+1/2}) - V(X_j)\|_{X_{j+1/2,*}}^2}} \quad (\text{E.2})$$

whereas by rearranging and developing we have:

$$\sum_{j=1}^{t-1} \|V(X_{j+1/2}) - V(X_j)\|_{X_{j+1/2,*}}^2 = \frac{1}{\gamma_t^2} - 1 \quad (\text{E.3})$$

Therefore, by taking limits on both sides:

$$\sum_{t=1}^{+\infty} \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2,*}}^2 = \lim_{t \rightarrow +\infty} \frac{1}{\gamma_t^2} - 1 = \frac{1}{\gamma_\infty^2} - 1 \geq 0 \quad (\text{E.4})$$

Hence, by recalling (C.14) we have:

$$\begin{aligned} \sum_{t=1}^T D(X_{t+1/2}, X_t) &\leq D(x^*, X_1) + \sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2,*}}^2 \\ &\leq D(x^*, X_1) + \sum_{t=1}^T \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2,*}}^2 \end{aligned}$$

which in turn by (E.4) yields $\sum_{t=1}^{+\infty} D(X_{t+1/2}, X_t) < +\infty$ and hence $D(X_{t+1/2}, X_t) \rightarrow 0$. Moreover, by applying (14):

$$\frac{K}{2} \|X_{t+1/2} - X_t\|_{X_t}^2 \leq D(X_{t+1/2}, X_t) \quad (\text{E.5})$$

Now, by recalling $\mu \|\cdot\| \leq \|\cdot\|_x$, we get:

$$\|X_{t+1/2} - X_t\|^2 \leq \frac{1}{\mu^2} \|X_{t+1/2} - X_t\|_{X_t}^2 \quad (\text{E.6})$$

and the result follows.

- $\gamma_\infty = 0$: By the prox-step, we get:

$$\begin{aligned} \langle \nabla h(X_t) - \nabla h(X_{t+1/2}), X_t - X_{t+1/2} \rangle &\leq \gamma_t \langle V(X_t), X_t - X_{t+1/2} \rangle \\ &\leq \gamma_t \|V(X_t)\|_{X_t,*} \|X_t - X_{t+1/2}\|_{X_t} \end{aligned} \quad (\text{E.7})$$

On the other hand, we have:

$$\langle \nabla h(X_t) - \nabla h(X_{t+1/2}), X_t - X_{t+1/2} \rangle = D(X_t, X_{t+1/2}) + D(X_{t+1/2}, X_t) \quad (\text{E.8})$$

Thus, we get by (14):

$$\begin{aligned} D(X_t, X_{t+1/2}) + D(X_{t+1/2}, X_t) &\leq \gamma_t \|V(X_t)\|_{X_t,*} \|X_t - X_{t+1/2}\|_{X_t} \\ &\leq \gamma_t M \sqrt{\frac{2}{K} [D(X_t, X_{t+1/2}) + D(X_{t+1/2}, X_t)]} \end{aligned}$$

where the last inequality is obtained due to (MB); which in turn yields:

$$D(X_t, X_{t+1/2}) + D(X_{t+1/2}, X_t) \leq \frac{2M^2}{K} \gamma_t^2 \quad (\text{E.9})$$

So, a fortiori we have:

$$D(X_t, X_{t+1/2}) \leq \frac{2M^2}{K} \gamma_t^2 \quad (\text{E.10})$$

Moreover, by (14):

$$\frac{K}{2} \|X_{t+1/2} - X_t\|_{X_{t+1/2}}^2 \leq D(X_t, X_{t+1/2}) \leq \frac{2M^2}{K} \gamma_t^2 \quad (\text{E.11})$$

Now, by recalling $\mu\|\cdot\| \leq \|\cdot\|_x$, we get:

$$\|X_{t+1/2} - X_t\|^2 \leq \frac{1}{\mu^2} \|X_{t+1/2} - X_t\|_{X_t}^2 \quad (\text{E.12})$$

and the result follows since we assumed that $\gamma_t \rightarrow 0$.

Case 2: Under (MS) condition. Following similar reasoning as above, we have:

$$\begin{aligned} \sum_{t=1}^T D(X_{t+1/2}, X_t) &\leq D(x^*, X_1) + \sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \\ &\leq D(x^*, X_1) + \sum_{t=1}^T \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \end{aligned}$$

which by taking limits on both sides and by applying Lemma D.1 we get that:

$$\sum_{t=1}^{+\infty} D(X_{t+1/2}, X_t) < +\infty \quad (\text{E.13})$$

Therefore, $D(X_{t+1/2}, X_t) \rightarrow 0$, whereas by applying (14) we obtain:

$$\frac{K}{2} \|X_{t+1/2} - X_t\|_{X_t}^2 \leq D(X_{t+1/2}, X_t) \quad (\text{E.14})$$

Now, by recalling $\mu\|\cdot\| \leq \|\cdot\|_x$, we get:

$$\|X_{t+1/2} - X_t\|^2 \leq \frac{1}{\mu^2} \|X_{t+1/2} - X_t\|_{X_t}^2 \quad (\text{E.15})$$

and the result follows.

On the other hand, for the second claim, we have by the prox-step:

$$\begin{aligned} D(X_t, X_{t+1/2}) + D(X_{t+1/2}, X_t) &\leq \gamma_t \langle V(X_t), X_{t+1/2} - X_t \rangle \\ &\leq \gamma_t M \|X_{t+1/2} - X_t\|_{X_t} \end{aligned}$$

Therefore, by following the same reasoning with the first claim, we get:

$$D(X_t, X_{t+1/2}) + D(X_{t+1/2}, X_t) \leq \frac{2M^2}{K} \gamma_t^2 \quad (\text{E.16})$$

and hence since $D(\cdot, \cdot) \geq 0$, we have:

$$D(X_{t+1/2}, X_t) \leq \frac{2M^2}{K} \gamma_t^2 \quad \text{and} \quad D(X_t, X_{t+1/2}) \leq \frac{2M^2}{K} \gamma_t^2 \quad (\text{E.17})$$

and so the result follows \square

Remark 2. We shall point out that (1) in Lemma E.1 establishes the convergence with respect to the global ambient reference norm of \mathbb{R}^d .

Proposition E.1. *Suppose that V satisfies (MB) (respectively (MS)). Then, the iterates $X_t, X_{t+1/2}$ of AdaProx possess convergent subsequences towards the equilibrium set \mathcal{X}^* .*

Proof. By Lemma E.1, it suffices to show that $X_{t+1/2}$ possesses such a subsequence. Assume to the contrary that it does not. That implies that:

$$\liminf_t \text{dist}(X_{t+1/2}, \mathcal{X}^*) = \delta > 0 \quad (\text{E.18})$$

which in turn yields,

$$\liminf_t \langle V(X_{t+1/2}), X_{t+1/2} - x^* \rangle = c > 0 \quad (\text{E.19})$$

Now, by setting $p = x^*$ for some $x^* \in \mathcal{X}^*$ in (C.14), we get:

$$\begin{aligned} D(x^*, X_{t+1}) &\leq D(x^*, X_t) - \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - x^* \rangle + \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2}}^2 \\ &\leq D(x^*, X_t) - c\gamma_t + \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2}}^2 \end{aligned}$$

whereas by telescoping $t = 1, \dots, T$ we obtain:

$$D(x^*, X_T) \leq D(x^*, X_1) - \sum_{t=1}^T \gamma_t \left[c - \frac{\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2}{\sum_{t=1}^T \gamma_t} \right] \quad (\text{E.20})$$

Having this established this general setting, we shall examine the asymptotic behaviour term by term for each regularity case individually, which in both cases shall lead to a contradiction.

Case 1: Under (MB) condition.

- For the first term: $\sum_{t=1}^T \gamma_t$, due to (AdaProx) we have by (D.29) that:

$$\sum_{t=1}^T \gamma_t \rightarrow +\infty \text{ and } \sum_{t=1}^T \gamma_t = \Omega(\sqrt{T}) \quad (\text{E.21})$$

- For the second term $\frac{\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2}{\sum_{t=1}^T \gamma_t}$, we first examine the denominator. In particular, due to (AdaProx) we get:

$$\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 = \sum_{t=1}^T \frac{\|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2}{1 + \sum_{j=1}^{t-1} \|V(X_{j+1/2}) - V(X_j)\|_{X_{j+1/2},*}^2} \quad (\text{E.22})$$

which by recalling (D.27) we obtain:

$$\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 = \mathcal{O}(\log T) \quad (\text{E.23})$$

So, by combining (E.21) and (E.23) we readily obtain:

$$\frac{\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2}{\sum_{t=1}^T \gamma_t} \rightarrow 0 \text{ while } T \rightarrow +\infty \quad (\text{E.24})$$

Therefore, by letting $T \rightarrow +\infty$, the inequality (E.20) yields $D(x^*, X_T) \rightarrow -\infty$, contradiction.

Case 2: Under (MS) condition. Examining the asymptotic behaviour of (E.20) term by term under the light of (MS) condition we get the following:

- For $\sum_{t=1}^T \gamma_t$, (MS) guarantees by (D.36):

$$\sum_{t=1}^T \gamma_t = \Omega(T) \text{ and } \sum_{t=1}^T \gamma_t \rightarrow +\infty \quad (\text{E.25})$$

- For $\frac{\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2}{\sum_{t=1}^T \gamma_t}$, (D.1) guarantees:

$$\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 = \mathcal{O}(1) \quad (\text{E.26})$$

which combined with (D.36) gives us:

$$\frac{\sum_{t=1}^T \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2}{\sum_{t=1}^T \gamma_t} \rightarrow 0 \quad (\text{E.27})$$

Therefore, y letting $T \rightarrow +\infty$, the inequality (E.20) yields that $D(x^*, X_T) \rightarrow -\infty$, a contradiction. \square

Having all this at hand, we are finally in the position to prove the main result of this section; namely the convergence of the actual iterates of the method. For that we will need an intermediate lemma that shall allow us to pass from a convergent subsequence to global convergence (see also [10], [44]).

Lemma E.2. *Let $\chi \in (0, 1]$, $(\alpha_t)_{t \in \mathbb{N}}$, $(\beta_t)_{t \in \mathbb{N}}$ non-negative sequences and $(\varepsilon_t)_{t \in \mathbb{N}} \in l^1(\mathbb{N})$ such that $t = 1, 2, \dots$:*

$$\alpha_{t+1} \leq \chi \alpha_t - \beta_t + \varepsilon_t \quad (\text{E.28})$$

Then, α_t converges.

Proof. First, one shows that $\alpha_{t \in \mathbb{N}}$ is a bounded sequence. Indeed, one can derive directly that:

$$\alpha_{t+1} \leq \chi^{t+1} \alpha_0 + \sum_{k=0}^t \chi^{t-k} \varepsilon_k \quad (\text{E.29})$$

Hence, $(\alpha_t)_{t \in \mathbb{N}}$ lies in $[0, \alpha_0 + \varepsilon]$, with $\varepsilon = \sum_{t=0}^{+\infty} \varepsilon_t$. Now, one is able to extract a convergent subsequence $(\alpha_{k_t})_{t \in \mathbb{N}}$, let say $\lim_{t \rightarrow +\infty} \alpha_{k_t} = \alpha \in [0, \alpha_0 + \varepsilon]$ and fix $\delta > 0$. Then, one can find some t_0 such that $\alpha_{k_{t_0}} - \alpha < \frac{\delta}{2}$ and $\sum_{m > t_{k_0}} \varepsilon_m < \frac{\delta}{2}$. That said, we have:

$$0 \leq \alpha_t \leq \alpha_{k_{t_0}} + \sum_{m > t_{k_0}} \varepsilon_m < \frac{\delta}{2} + \alpha + \frac{\delta}{2} = \alpha + \delta \quad (\text{E.30})$$

Hence, $\limsup_t \alpha_t \leq \liminf_t \alpha_t + \delta$. Since, δ is chosen arbitrarily the result follows. \square

Proof of Theorem 3. Once more, we shall treat each regularity class individually.

Case 1: Under (MB) condition. For the (MB), by denoting $\lim_{t \rightarrow +\infty} \gamma_t = \gamma_\infty$ case we shall consider two cases for the asymptotic behaviour of the step-size γ_t .

- $\gamma_\infty > 0$: By recalling the definition of γ_t :

$$\gamma_t = \frac{1}{\sqrt{1 + \sum_{j=1}^{t-1} \|V(X_{j+1/2}) - V(X_j)\|_{X_{j+1/2}}^2}} \quad (\text{E.31})$$

whereas by rearranging we get:

$$\sum_{j=1}^{t-1} \|V(X_{j+1/2}) - V(X_j)\|_{X_{j+1/2}}^2 = \frac{1}{\gamma_t^2} - 1 \quad (\text{E.32})$$

and hence:

$$\sum_{t=1}^{+\infty} \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2}}^2 = \frac{1}{\gamma_\infty^2} - 1 < +\infty \quad (\text{E.33})$$

Therefore, by recalling (C.14), we have for solution of (VI), $x^* \in \mathcal{X}$

$$D(x^*, X_{t+1}) \leq D(x^*, X_t) - \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - x^* \rangle + \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \quad (\text{E.34})$$

which enables us to directly apply Lemma E.2 for $\alpha_t = D(x^*, X_t)$, $\beta_t = \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - x^* \rangle$ and $\varepsilon_t = \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2$.

- $\gamma_\infty = 0$: Fix an equilibrium $x^* \in \mathcal{X}^*$ and consider the "Bregman zone":

$$D_\varepsilon = \{x \in \mathcal{X} : D(x^*, x) < \varepsilon\} \quad (\text{E.35})$$

By the assumption for the regularizer h , it follows that there exists some $\delta > 0$ such that:

$$B_\delta = \{x \in \mathcal{X} : \|x^* - x\| < \delta\} \quad (\text{E.36})$$

is contained in D_ε . Hence, by regularity assumption for the (3), it follows that:

$$\langle V(x), x - x^* \rangle \geq c > 0 \text{ for some } c \equiv c(\varepsilon) > 0 \text{ and for all } x \notin D_\varepsilon, \quad (\text{E.37})$$

in particular, for all $x \in D_{2\varepsilon} \setminus D_\varepsilon$. Assume now that x^* is a limit point of X_t , i.e., $X_t \in D_{2\varepsilon}$ for infinitely many $t \in \mathbb{N}$. Now, by the prox-step, we get: and hence,

$$\gamma_t \langle V(X_t), X_t - x^* \rangle \leq \langle \nabla h(X_t) - \nabla h(X_{t+1/2}), X_t - x^* \rangle \quad (\text{E.38})$$

whereas by [Lemma B.2](#) and after rearranging we get:

$$\begin{aligned} D(x^*, X_{t+1/2}) &\leq D(x^*, X_t) - \gamma_t \langle V(X_t), X_t - x^* \rangle + D(X_t, X_{t+1/2}) \\ &\leq D(x^*, X_t) - \gamma_t \langle V(X_t), X_t - x^* \rangle + \max\{D(X_t, X_{t+1/2}), D(X_t, X_{t+1/2})\} \end{aligned}$$

Therefore due to [Lemma E.1](#) we obtain:

$$D(x^*, X_{t+1/2}) \leq D(x^*, X_t) - \gamma_t \langle V(X_t), X_t - x^* \rangle + \frac{2M^2}{K} \gamma_t^2 \quad (\text{E.39})$$

We consider two cases:

1. $X_t \in D_{2\varepsilon} \setminus D_\varepsilon$: Then, $\langle V(X_t), X_t - x^* \rangle \geq c > 0$. So,

$$D(x^*, X_{t+1/2}) \leq D(x^*, X_t) - c\gamma_t + \frac{2M^2}{K} \gamma_t^2 \quad (\text{E.40})$$

Now, provided that $\frac{2M^2\gamma_t^2}{K} \leq c\gamma_t$ or equivalently $\gamma_t \leq \frac{cK}{2M^2}$. we get: $D(x^*, X_{t+1/2}) \leq 2\varepsilon$.

2. $X_t \in D_\varepsilon$: Then, in this case we have:

$$D(x^*, X_{t+1/2}) \leq D(x^*, X_t) + \frac{2M^2}{K} \gamma_t^2 \quad (\text{E.41})$$

Again, provided that $\frac{2M^2}{K} \gamma_t^2 \leq \varepsilon$ or equivalently $\gamma_t \leq \frac{\sqrt{2\varepsilon K}}{2M}$ we get $D(x^*, X_{t+1/2}) \leq 2\varepsilon$

Therefore, by summarizing the above we get that if $\gamma_t \leq \min\{\frac{\sqrt{2\varepsilon K}}{2M}, \frac{cK}{2M^2}\}$, we have that $X_{t+1/2} \in D_{2\varepsilon}$ whenever $X_t \in D_{2\varepsilon}$. Going further, due to [Proposition B.2](#) by setting $p = x^*$, $x_1 = X_{t+1/2}$, $x_2^+ = X_{t+1}$, $x = X_t$, $v_1 = -\gamma_t V(X_{t+1/2})$ and $v_2 = -\gamma_t V(X_{t+1/2})$ we get:

$$\begin{aligned} D(x^*, X_{t+1}) &\leq D(x^*, X_t) - \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - x^* \rangle + \gamma_t \langle V(X_{t+1/2}) - V(X_t), X_{t+1} - X_{t+1/2} \rangle \\ &\quad - D(X_{t+1}, X_{t+1/2}) - D(X_{t+1/2}, X_t) \quad (\text{E.42}) \end{aligned}$$

whereas by applying Fenchel's inequality we obtain:

$$\begin{aligned} D(x^*, X_{t+1}) &\leq D(x^*, X_t) - \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - x^* \rangle + \frac{\gamma_t^2}{2K} \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \\ &\quad + \frac{K}{2} \|X_{t+1} - X_{t+1/2}\|_{X_{t+1/2}}^2 - D(X_{t+1}, X_{t+1/2}) - D(X_{t+1/2}, X_t) \quad (\text{E.43}) \end{aligned}$$

Now, since $\frac{K}{2} \|X_{t+1} - X_{t+1/2}\|_{X_{t+1/2}}^2 - D(X_{t+1}, X_{t+1/2}) \leq 0$ by [\(14\)](#) we get:

$$D(x^*, X_{t+1}) \leq D(x^*, X_t) - \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - x^* \rangle + \frac{\gamma_t^2}{2K} \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \quad (\text{E.44})$$

which, in turn, by [\(C.13\)](#) the above yields:

$$D(x^*, X_{t+1}) \leq D(x^*, X_t) - \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - x^* \rangle + \frac{C^2}{2K} \gamma_t^2 \quad (\text{E.45})$$

with $C = 2M + \beta \frac{4M}{K}$. Recall that $X_{t+1/2} \in D_{2\varepsilon}$ by our previous claim. We now consider the following two cases:

1. $X_{t+1/2} \in D_{2\varepsilon} \setminus D_\varepsilon$: In this case: $\langle V(X_{t+1/2}), X_{t+1/2} - x^* \rangle \geq c > 0$, so,

$$D(x^*, X_{t+1}) \leq D(x^*, X_t) - c\gamma_t + \frac{C^2}{2K} \gamma_t^2 \quad (\text{E.46})$$

which holds provided that $\frac{C^2\gamma_t^2}{2K} \leq c\gamma_t$ or equivalently $\gamma_t \leq \frac{2cK}{C^2}$,

2. $X_{t+1/2} \in D_\varepsilon$: First recall that:

$$D(X_{t+1/2}, X_{t+1}) + D(X_{t+1}, X_{t+1/2}) \leq \frac{2\gamma_t^2}{K} \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2 \leq \frac{2\gamma_t^2}{K} C^2 \quad (\text{E.47})$$

Therefore, we get that:

$$\|X_{t+1} - X_{t+1/2}\|^2 \leq \frac{4\mu^2 C^2}{K^2} \gamma_t^2 \quad (\text{E.48})$$

Now, let us define the following:

$$D_\varepsilon(\alpha) = \max\{D(x^*, x) : \text{dist}(x, D_\varepsilon(x^*)) < \alpha\} \quad (\text{E.49})$$

Clearly, $D_\varepsilon(\alpha)$ is continuous relative to α and $\lim_{\alpha \rightarrow 0^+} D_\varepsilon(\alpha) = \varepsilon$. Therefore, we have:

$$D_\varepsilon(\alpha) \leq \varepsilon \quad \text{for all } \alpha \leq \alpha^* \text{ with } \alpha^* \text{ sufficiently small.} \quad (\text{E.50})$$

Moreover, due to (E.48), we conclude that $D(x^*, X_{t+1}) \leq 2\varepsilon$, provided that $\gamma_t \leq \frac{\alpha^*}{2\mu C} K$.

We conclude that $X_{t+1} \in U_{2\varepsilon}$ provided that $X_t \in D_{2\varepsilon}$ and $\gamma_t \leq \min\{\frac{2\alpha^* K}{M^2}, \frac{\sqrt{2\varepsilon K}}{2M}, \frac{\alpha^*}{2\mu C} K\}$. Since, $\gamma_t \rightarrow 0$ and $X_t \in D_{2\varepsilon}$ infinitely often (due to Proposition E.1) we conclude that $X_t \in D_{2\varepsilon}$ for all sufficiently large t . With $\varepsilon > 0$ being arbitrary, the result follows.

Case 2: Under (MS) condition. By plugging in $\alpha_t = D(x^*, X_t)$, $\beta_t = \gamma_t \langle V(X_{t+1/2}), X_{t+1/2} - x^* \rangle$ and $\varepsilon_t = \gamma_t^2 \|V(X_{t+1/2}) - V(X_t)\|_{X_{t+1/2},*}^2$ in Lemma E.2 and combine it with Lemma D.1, we get $\inf_{x^* \in \mathcal{X}^*} \|x^*, X_t\|$ converges. Thus, the result follows by applying Proposition E.1 \square

F PROPERTIES OF NUMERICAL SEQUENCES

In this appendix, we provide the necessary inequality of numerical sequences. This inequality is due to Bach & Levy [2] and Levy et al. [27] and will play an indispensable role for establishing the last iterate convergence and universality of our method.

Lemma F.1. For all non-negative numbers $\alpha_1, \dots, \alpha_T$, the following inequality holds:

$$\sum_{i=1}^T \frac{\alpha_i}{1 + \sum_{i=1}^i \alpha_i} \leq 1 + \log(1 + \sum_{i=1}^T \alpha_i) \quad (\text{F.1})$$

Proof. The lemma will be proved by induction. The induction base $T = 1$ holds, since:

$$\frac{\alpha_1}{1 + \alpha_1} \leq 1 \leq 1 + \log(1 + \alpha_1) \quad (\text{F.2})$$

Assume now that the lemma holds for $T - 1$. Then, we are left to show that it also holds for T . Indeed, by the induction hypothesis, we get:

$$\sum_{i=1}^T \frac{\alpha_i}{1 + \sum_{i=1}^i \alpha_i} \leq 1 + \log(1 + \sum_{i=1}^{T-1} \alpha_i) + \frac{\alpha_T}{1 + \sum_{i=1}^T \alpha_i} \quad (\text{F.3})$$

Thus, in order to complete the induction it suffices to show that:

$$1 + \log(1 + \sum_{i=1}^{T-1} \alpha_i) + \frac{\alpha_T}{1 + \sum_{i=1}^T \alpha_i} \leq 1 + \log(1 + \sum_{i=1}^T \alpha_i) \quad (\text{F.4})$$

By denoting $x = \alpha_T / (1 + \sum_{i=1}^{T-1} \alpha_i)$, the above equation is equivalent:

$$\log(x + 1) - \frac{x}{1 + x} \geq 0 \quad (\text{F.5})$$

which can be straightforwardly checked since $H(x) = \log(x + 1) - \frac{x}{1+x} \geq 0$ for all $x \geq 0$. Therefore, the result follows. \square

REFERENCES

- [1] Kimon Antonakopoulos, E. Veronica Belmega, and Panayotis Mertikopoulos. An adaptive mirror-prox algorithm for variational inequalities with singular operators. In *NeurIPS '19: Proceedings of the 33rd International Conference on Neural Information Processing Systems*, 2019.
- [2] Francois Bach and Kfir Yehuda Levy. A universal algorithm for variational inequalities adaptive to smoothness and noise. In *COLT '19: Proceedings of the 32nd Annual Conference on Learning Theory*, 2019.
- [3] David Balduzzi, Sebastien Racaniere, James Martens, Jakob Foerster, Karl Tuyls, and Thore Graepel. The mechanics of n -player differentiable games. In *ICML '18: Proceedings of the 35th International Conference on Machine Learning*, 2018.
- [4] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- [5] Dimitri P. Bertsekas and Robert Gallager. *Data Networks*. Prentice Hall, Englewood Cliffs, NJ, 2 edition, 1992.
- [6] Lev M. Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7(3):200–217, 1967.
- [7] Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends in Machine Learning*, 8(3-4):231–358, 2015.
- [8] Tatjana Chavdarova, Gauthier Gidel, François Fleuret, and Simon Lacoste-Julien. Reducing noise in GAN training with variance reduced extragradient. In *NeurIPS '19: Proceedings of the 33rd International Conference on Neural Information Processing Systems*, 2019.
- [9] Gong Chen and Marc Teboulle. Convergence analysis of a proximal-like minimization algorithm using Bregman functions. *SIAM Journal on Optimization*, 3(3):538–543, August 1993.
- [10] Patrick L. Combettes. Quasi-Fejérian analysis of some optimization algorithms. In Dan Butnariu, Yair Censor, and Simeon Reich (eds.), *Inherently Parallel Algorithms in Feasibility and Optimization and Their Applications*, pp. 115–152. Elsevier, New York, NY, USA, 2001.
- [11] Constantinos Daskalakis, Paul W. Goldberg, and Christos H. Papadimitriou. The complexity of computing a Nash equilibrium. *SIAM Journal on Computing*, 39(1):195–259, 2009.
- [12] Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training GANs with optimism. In *ICLR '18: Proceedings of the 2018 International Conference on Learning Representations*, 2018.
- [13] Gérard Debreu. A social equilibrium existence theorem. *Proceedings of the National Academy of Sciences of the USA*, 38(10):886–893, October 1952.
- [14] Francisco Facchinei and Christian Kanzow. Generalized Nash equilibrium problems. *4OR*, 5(3):173–210, September 2007.
- [15] Francisco Facchinei and Jong-Shi Pang. *Finite-Dimensional Variational Inequalities and Complementarity Problems*. Springer Series in Operations Research. Springer, 2003.
- [16] Lampros Flokas, Emmanouil Vasileios Vlatakis-Gkaragkounis, and Georgios Piliouras. Poincaré recurrence, cycles and spurious equilibria in gradient-descent-ascent for non-convex non-concave zero-sum games. In *NeurIPS '19: Proceedings of the 33rd International Conference on Neural Information Processing Systems*, 2019.
- [17] A.V. Gasnikov, P.E. Dvurechensky, F.S. Stonyakin, and A.A. Titov. An adaptive proximal method for variational inequalities. *Computational Mathematics and Mathematical Physics*, 59:836–841, 2019.
- [18] Gauthier Gidel, Hugo Berard, Gaëtan Vignoud, Pascal Vincent, and Simon Lacoste-Julien. A variational inequality perspective on generative adversarial networks. In *ICLR '19: Proceedings of the 2019 International Conference on Learning Representations*, 2019.
- [19] Gauthier Gidel, Reyhane Askari Hemmat, Mohammad Pezehski, Rémi Le Priol, Gabriel Huang, Simon Lacoste-Julien, and Ioannis Mitliagkas. Negative momentum for improved game dynamics. In *AISTATS '19: Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, 2019.
- [20] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS '14: Proceedings of the 28th International Conference on Neural Information Processing Systems*, 2014.
- [21] Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Fundamentals of Convex Analysis*. Springer, Berlin, 2001.
- [22] Yu-Guan Hsieh, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. On the convergence of single-call stochastic extra-gradient methods. In *NeurIPS '19: Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pp. 6936–6946, 2019.

- [23] Yu-Guan Hsieh, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. Explore aggressively, update conservatively: Stochastic extragradient methods with variable stepsize scaling. <https://arxiv.org/abs/2003.10162>, 2020.
- [24] Anatoli Juditsky, Arkadi Semen Nemirovski, and Claire Tauvel. Solving variational inequalities with stochastic mirror-prox algorithm. *Stochastic Systems*, 1(1):17–58, 2011.
- [25] G. M. Korpelevich. The extragradient method for finding saddle points and other problems. *Èkonom. i Mat. Metody*, 12:747–756, 1976.
- [26] Rida Laraki, Jérôme Renault, and Sylvain Sorin. *Mathematical Foundations of Game Theory*. Universitext. Springer, 2019.
- [27] Kfir Yehuda Levy, Alp Yurtsever, and Volkan Cevher. Online adaptive methods, universality and acceleration. In *NeurIPS '18: Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2018.
- [28] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR '18: Proceedings of the 2018 International Conference on Learning Representations*, 2018.
- [29] Yura Malitsky. Projected reflected gradient methods for monotone variational inequalities. *SIAM Journal on Optimization*, 25(1):502–520, 2015.
- [30] Yura Malitsky. Golden ratio algorithms for variational inequalities. *Mathematical Programming*, 2019.
- [31] Panayotis Mertikopoulos and Zhengyuan Zhou. Learning in games with continuous action sets and unknown payoff functions. *Mathematical Programming*, 173(1-2):465–507, January 2019.
- [32] Panayotis Mertikopoulos, Christos H. Papadimitriou, and Georgios Piliouras. Cycles in adversarial regularized learning. In *SODA '18: Proceedings of the 29th annual ACM-SIAM Symposium on Discrete Algorithms*, 2018.
- [33] Panayotis Mertikopoulos, Bruno Lecouat, Houssam Zenati, Chuan-Sheng Foo, Vijay Chandrasekhar, and Georgios Piliouras. Optimistic mirror descent in saddle-point problems: Going the extra (gradient) mile. In *ICLR '19: Proceedings of the 2019 International Conference on Learning Representations*, 2019.
- [34] Aryan Mokhtari, Asuman Ozdaglar, and Sarath Pattathil. Convergence rate of $\mathcal{O}(1/k)$ for optimistic gradient and extra-gradient methods in smooth convex-concave saddle point problems. <https://arxiv.org/pdf/1906.01115.pdf>, 2019.
- [35] Arkadi Semen Nemirovski. Information-based complexity of linear operator equations. *Journal of Complexity*, 8(2):153–175, 1992.
- [36] Arkadi Semen Nemirovski. Prox-method with rate of convergence $O(1/t)$ for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1):229–251, 2004.
- [37] Arkadi Semen Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- [38] Yurii Nesterov. Dual extrapolation and its applications to solving variational inequalities and related problems. *Mathematical Programming*, 109(2):319–344, 2007.
- [39] Yurii Nesterov. Primal-dual subgradient methods for convex problems. *Mathematical Programming*, 120(1):221–259, 2009.
- [40] Noam Nisan, Tim Roughgarden, Éva Tardos, and V. V. Vazirani (eds.). *Algorithmic Game Theory*. Cambridge University Press, 2007.
- [41] Yuyuan Ouyang and Yangyang Xu. Lower complexity bounds of first-order methods for convex-concave bilinear saddle-point problems. *Mathematical Programming*, 2019. URL <https://doi.org/10.1007/s10107-019-01420-0>.
- [42] Georgios Piliouras and Jeff S. Shamma. Optimization despite chaos: Convex relaxations to complex limit sets via Poincaré recurrence. In *SODA '14: Proceedings of the 25th annual ACM-SIAM Symposium on Discrete Algorithms*, 2014.
- [43] Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. Robust adversarial reinforcement learning. In *ICML '17: Proceedings of the 34th International Conference on Machine Learning*, 2017.
- [44] Boris Teodorovich Polyak. *Introduction to Optimization*. Optimization Software, New York, NY, USA, 1987.
- [45] Leonid Denisovich Popov. A modification of the Arrow–Hurwicz method for search of saddle points. *Mathematical Notes of the Academy of Sciences of the USSR*, 28(5):845–848, 1980.
- [46] Alexander Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable sequences. In *NIPS '13: Proceedings of the 27th International Conference on Neural Information Processing Systems*, 2013.

- [47] Ralph Tyrrell Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, NJ, 1970.
- [48] Gesualdo Scutari, Francisco Facchinei, Daniel Pérez Palomar, and Jong-Shi Pang. Convex optimization, game theory, and variational inequality theory in multiuser communication systems. *IEEE Signal Process. Mag.*, 27(3):35–49, May 2010.
- [49] Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends in Machine Learning*, 4(2):107–194, 2011.
- [50] Fedor Stonyakin, Alexander Gasnikov, Pavel Dvurechensky, Mohammad Alkousa, and Alexander Titov. Generalized mirror prox for monotone variational inequalities: Universality and inexact oracle. <https://arxiv.org/abs/1806.05140>, 2018.
- [51] Fedor Stonyakin, Alexander Gasnikov, Alexander Tyurin, Dmitry Pasechnyuk, Artem Agafonov, Pavel Dvurechensky, Darina Dvinskikh, Alexey Kroshnin, and Victorya Piskunova. Inexact model: A framework for optimization and variational inequalities. <https://arxiv.org/abs/1902.00990>, 2019.
- [52] John von Neumann. Zur Theorie der Gesellschaftsspiele. *Mathematische Annalen*, 100:295–320, 1928. Translated by S. Bargmann as “On the Theory of Games of Strategy” in A. Tucker and R. D. Luce, editors, *Contributions to the Theory of Games IV*, volume 40 of *Annals of Mathematics Studies*, pages 13–42, 1957, Princeton University Press, Princeton.