

# Nonlinearly preconditioned gradient flows

Konstantinos Oikonomidis

Alexander Bodard

Jan Quan

Panagiotis Patrinos

**Abstract**—We study a continuous-time dynamical system which arises as the limit of a broad class of nonlinearly preconditioned gradient methods. Under mild assumptions, we establish existence of global solutions and derive Lyapunov-based convergence guarantees. For convex costs, we prove a sublinear decay in a geometry induced by some reference function, and under a generalized gradient-dominance condition we obtain exponential convergence. We further uncover a duality connection with mirror descent, and use it to establish that the flow of interest solves an infinite-horizon optimal-control problem of which the value function is the Bregman divergence generated by the cost. These results clarify the structure and optimization behavior of nonlinearly preconditioned gradient flows and connect them to known continuous-time models in non-Euclidean optimization.

## I. INTRODUCTION

Gradient flows have attracted substantial interest in both continuous-time optimization and control. A major reason is that they offer valuable intuition and enable Lyapunov-based analyses of optimization algorithms. A classical example is the standard gradient flow

$$\dot{x}(t) = -\nabla f(x(t)), \quad (1)$$

which can be viewed as the continuous-time limit of the gradient descent iteration

$$x^{k+1} = x^k - \gamma \nabla f(x^k), \quad (2)$$

obtained by letting the stepsize  $\gamma \rightarrow 0$ .

Recently, this perspective has been extended well beyond the standard gradient flow, with increasingly sophisticated dynamical systems being proposed to capture the behavior of modern optimization algorithms. Prominent examples include second-order dissipative systems that model methods such as Polyak’s heavy-ball algorithm and Nesterov’s accelerated gradient algorithm [1], [2]. Additional developments involve *mirror flows*, which arise as continuous-time analogues of mirror descent and more general Bregman-type methods [3], [4], and *proximal flows*, which correspond to the proximal point algorithm and related schemes [5].

We also draw attention to *normalized flows* [6], which serve as continuous-time counterparts of normalized gradient methods [7]. This family of algorithms is closely connected to *gradient clipping techniques* [8] and has recently attracted

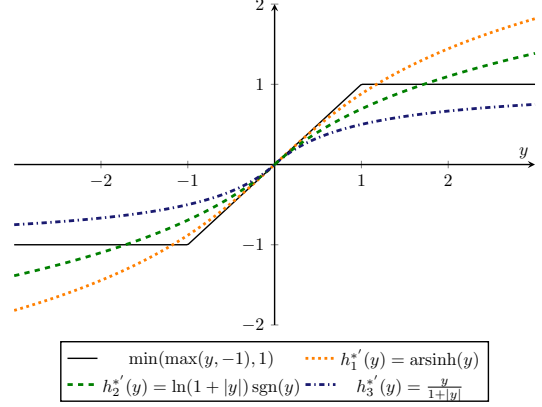


Fig. 1. Preconditioners corresponding to different reference functions.

considerable interest due to its effectiveness in neural network training and its robustness in stochastic optimization settings with heavy-tailed noise.

In [9], it was shown that these methods fit within a broader family of algorithms known as *nonlinearly preconditioned gradient methods* [9], [10], and also referred to as anisotropic gradient descent [11] or dual space preconditioning [12] methods. This class is characterized by updates of the form

$$x^{k+1} = x^k - \gamma \nabla \phi^*(\nabla f(x^k)), \quad (3)$$

where  $\phi : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$  is a convex, proper and lower semi-continuous *reference function*, with  $\phi^*$  denoting its convex conjugate, and  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is the continuously differentiable cost function. The particular choice  $\phi = \frac{1}{2} \|\cdot\|^2$ , yields  $\phi^* = \frac{1}{2} \|\cdot\|^2$ , such that the update rule (3) reduces to the classical gradient descent iteration (2). It is straightforward to verify that gradient normalization also fits within this framework by selecting a reference function  $\phi$  with bounded domain. For instance, consider  $\phi(x) = -\varepsilon(\ln(1 - \|x\|) + \|x\|)$ , with  $\text{dom } \phi = \{x \in \mathbb{R}^n \mid \|x\| < 1\}$ . For this choice, one obtains  $\nabla \phi^*(y) = \frac{y}{\|y\| + \varepsilon}$ . Hence, substituting into (3) yields an update that, as  $\varepsilon \rightarrow 0$ , approximates the classical normalized gradient method [7, Eq.(NSGD)].

A variety of other reference functions  $\phi$  have been proposed. In particular, isotropic choices of the form  $\phi = h \circ \|\cdot\|$ , for suitable scalar functions  $h$ , yield algorithms that mimic gradient clipping, as illustrated in Figure 1. A notable feature of (3) is that its convergence is naturally analyzed under the framework of *anisotropic smoothness* [11], a milder condition than standard Lipschitz smoothness that is satisfied in important applications such as matrix factorization and phase retrieval [13].

In this work we address the open problem of analyzing

This work was supported by the Research Foundation Flanders (FWO) PhD grant 11A8T26N and research projects G081222N, G033822N, and G0A0920N; Research Council KUL grant C14/24/103.

KU Leuven, Department of Electrical Engineering ESAT-STADIUS – Kasteelpark Arenberg 10, box 2446, 3001 Leuven, Belgium [konstantinos.oikonomidis@kuleuven.be](mailto:konstantinos.oikonomidis@kuleuven.be),

nonlinearly preconditioned gradient flows, defined by

$$\dot{x}(t) = -\nabla\phi^*(\nabla f(x(t))). \quad (4)$$

These dynamics arise as the continuous-time analogue of the discrete-time iteration (3) and provide a natural framework for understanding nonlinearly preconditioned gradient methods.

In fact, equations of the form (4) belong to a wider family of nonlinear differential equations, known as *doubly nonlinear equations* [14]. They are described by the general inclusion problem

$$\partial\phi(\dot{u}) + \partial f(u) \ni g,$$

where in this case  $u : [0, T) \rightarrow \mathcal{H}$  with  $\mathcal{H}$  a real Hilbert space,  $g \in L_2([0, \infty), \mathcal{H})$  and  $\partial\phi$  denotes the standard subdifferential from convex analysis [15, Def. 16.1]. As described in [14], this inclusion has a straightforward mechanical interpretation:  $u$  represents the displacement of a body,  $g$  the corresponding energy and  $\phi$  the related dissipation potential. Various works have been devoted to the study of such nonlinear dynamics, including [16] which proposes a generalized inclusion, [17] which studies a second-order system and [18] where a stochastic model is analyzed.

#### A. Contribution

Our contributions can be summarized as follows:

- We study the optimization properties of the nonlinearly preconditioned gradient flow (4) for solving smooth unconstrained minimization problems, obtaining convergence rate guarantees for suitable optimality measures. Through a novel Lyapunov-like function we retrieve standard  $1/t$  rates for the suboptimality gap  $f(x(t)) - \inf f$  and prove the exponential convergence of the method under a generalized PL inequality.
- Through a duality connection with the mirror flow [4] we obtain an optimal control viewpoint on the nonlinearly preconditioned gradient flow (4). This further allows us to show that, under suitable assumptions, among all stabilizing controls,  $-\nabla\phi^*(\nabla f(x))$  constitutes a control that minimizes the cost  $\int_0^\infty q(x(t), u(t))$ , with  $q$  defined later on. The value function  $V(x_0)$  that corresponds to this cost is given by the Bregman divergence generated by  $f$ .

#### B. Notation

We denote the extended real line by  $\overline{\mathbb{R}} = \mathbb{R} \cup \{\infty\}$ , the standard Euclidean inner product on  $\mathbb{R}^n$  by  $\langle \cdot, \cdot \rangle$ , and its induced norm by  $\|\cdot\|$ . We denote by  $\mathcal{C}^k(Y)$  the class of functions that are  $k$  times continuously differentiable on an open set  $Y \subseteq \mathbb{R}^n$ . We say that a function  $f : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$  is *level-bounded* if for every  $\alpha \in \mathbb{R}$ , the set  $\{x \in \mathbb{R}^n \mid f(x) \leq \alpha\}$  is bounded. We adopt the notions of essential smoothness, essential strict convexity and Legendre functions from [19, Section 26]: we say that a proper, lsc and convex function  $f : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$  is *essentially smooth* if  $\text{int}(\text{dom } f) \neq \emptyset$  and  $f$  is differentiable on  $\text{int}(\text{dom } f)$  such that  $\|\nabla f(x^\nu)\| \rightarrow \infty$ ,

whenever  $\text{int}(\text{dom } f) \ni x^\nu \rightarrow x \in \text{bdry } \text{dom } f$ , and *essentially strictly convex*, if  $f$  is strictly convex on every convex subset of  $\text{dom } \partial f$ , and *Legendre*, if  $f$  is both essentially smooth and essentially strictly convex. In particular, a smooth convex function on  $\mathbb{R}^n$  is essentially smooth. For a Legendre function  $f$ , we denote the generated Bregman divergence  $D_f(x, \bar{x}) = f(x) - f(\bar{x}) - \langle \nabla f(\bar{x}), x - \bar{x} \rangle$ .

## II. CONVERGENCE GUARANTEES IN CONTINUOUS TIME

In this section we study the convergence properties of the trajectory solving (4), deriving sublinear convergence rates in the convex case and exponential rates under a generalized gradient dominance condition that fits the geometry of the gradient flow. We begin by formulating the optimization problem that we study

$$\min_{x \in \mathbb{R}^n} f(x). \quad (5)$$

Next, we present our assumptions on the problem data and our reference function.

**Assumption II.1.** *We have the following:*

- 1)  $f \in \mathcal{C}^2(\mathbb{R}^n)$  is level-bounded.
- 2)  $\phi : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$  is essentially smooth and strongly convex with constant  $\mu_\phi > 0$ . Moreover,  $\phi \geq 0$  is even and  $\phi(0) = 0$ .

We remark that due to the duality of strong convexity and Lipschitz smoothness [15, Thm. 18],  $\nabla\phi^* : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is  $1/\mu_\phi$ -Lipschitz continuous. Moreover,  $\nabla\phi^*$  is a bijection between  $\mathbb{R}^n$  and  $\text{int } \text{dom } \phi$  [19, Thm. 26.5] and  $\nabla\phi^*(y) = 0$  if and only if  $y = 0$ .

Clearly, since  $f \in \mathcal{C}^2(\mathbb{R}^n)$ ,  $f$  is locally Lipschitz differentiable and thus the mapping  $\nabla\phi^* \circ \nabla f$  is locally Lipschitz continuous. Moreover, since  $f$  is level-bounded and continuous,  $\arg \min f \neq \emptyset$  from [20, Thm. 1.9] and we denote  $f_\star = \inf f$ .

**Remark II.2.** The assumptions on the level-boundedness of  $f$  and the global Lipschitz differentiability of  $\phi^*$  can be relaxed. Nevertheless, we choose the current simple set of assumptions in order to simplify the exposition and emphasize the obtained convergence rates that we present in the following.

We now move on to our first result that describes the global existence and uniqueness of the solution of (4).

**Proposition II.3.** *Let Assumption II.1 hold true. Then, for any initial point  $x(0) = x_0 \in \mathbb{R}^n$ , there exists a unique global solution  $x : [0, \infty) \rightarrow \mathbb{R}^n$  of the system (4). Moreover,  $\dot{x} \in L^2([0, \infty), \mathbb{R}^n)$  and if  $\phi \in \mathcal{C}^2(\text{int } \text{dom } \phi)$ , then  $\lim_{t \rightarrow \infty} \dot{x}(t) = 0$ .*

*Proof.* Since  $F := \nabla\phi^* \circ \nabla f$  is locally Lipschitz, it is Lipschitz continuous on a neighborhood of  $x_0$ . From the Cauchy–Lipschitz theorem, this means that there exists some  $\delta > 0$  such that (4) admits a unique solution on  $[0, \delta]$ . Now note that for any  $t \in [0, \delta]$  we have

$$\frac{d}{dt}(f(x(t))) = \langle \nabla f(x(t)), -\nabla\phi^*(\nabla f(x(t))) \rangle \leq 0, \quad (6)$$

where we substituted (4), and used  $\nabla\phi^*(0) = 0$  and the monotonicity of  $\nabla\phi^*$ . Therefore,  $f(x(t)) \leq f(x_0)$  for all  $t \in [0, \delta]$  and since  $f$  is continuous and level-bounded,  $x$  is contained in a compact set. Hence, using standard arguments, we obtain the existence and uniqueness of  $x$  in  $[0, \infty)$ .

Now integrating (6) from 0 to  $t$  we obtain

$$\int_0^t \langle \nabla f(x(s)), \nabla\phi^*(\nabla f(x(s))) \rangle ds = f(x_0) - f(x(t)).$$

Since this equation is true for any  $t \geq 0$ , we have that  $\int_0^\infty \langle \nabla f(x(s)), \nabla\phi^*(\nabla f(x(s))) \rangle ds \leq f(x_0) - f_*$ . Since  $\phi$  is  $\mu_\phi$ -strongly convex,  $\nabla\phi^*$  is  $\mu_\phi$ -cocoercive from [15, Thm. 18.15] and thus  $\langle \nabla\phi^*(y), y \rangle \geq \mu_\phi \|\nabla\phi^*(y)\|^2$  for all  $y \in \mathbb{R}^n$ . Using this inequality and  $\dot{x} = -\nabla\phi^*(\nabla f(x(t)))$  we then have that

$$\int_0^\infty \|\dot{x}(s)\|^2 ds \leq \frac{1}{\mu_\phi} (f(x_0) - f_*) < +\infty,$$

implying  $\dot{x} \in L^2([0, \infty), \mathbb{R}^n)$ . Assuming moreover that  $\phi \in \mathcal{C}^2(\text{int dom } \phi)$  we have from [21, p. 42] that  $\phi^* \in \mathcal{C}^2(\text{int dom } \phi^*)$  and since  $\text{dom } \phi^* = \mathbb{R}^n$ ,  $\phi^* \in \mathcal{C}^2(\mathbb{R}^n)$ . Differentiating thus (4), we obtain

$$\begin{aligned} \ddot{x}(t) &= -\nabla^2\phi^*(\nabla f(x(t)))\nabla^2 f(x(t))\dot{x}(t) \\ &= \nabla^2\phi^*(\nabla f(x(t)))\nabla^2 f(x(t))\nabla\phi^*(\nabla f(x(t))) \end{aligned}$$

Now note that all the functions in the r.h.s. are continuous and since  $x$  remains in a bounded set, we have that  $\sup_{t \in [0, \infty)} \|\ddot{x}(t)\| < +\infty$ . From Barbalat's lemma we now obtain the claimed result.  $\square$

Note that in the classical setting of the gradient flow described in (4) where  $\phi = \frac{1}{2}\|\cdot\|^2$ , the r.h.s. of (6) becomes  $-\|\nabla f(x(t))\|^2$ , thus leading to convergence guarantees for the standard stationarity measure. In our more general setting, we obtain, in light of the equality case of the Fenchel–Young inequality [20, Prop. 11.3],

$$\frac{d}{dt}(f(x(t))) = -[\phi^*(\nabla f(x(t))) + \phi(\nabla\phi^*(\nabla f(x(t))))], \quad (7)$$

thus extending the convergence guarantees in the nonconvex setting described in [9, Thm. 3.2].

We now move on to the more interesting setting where  $f$  is convex. Note that obtaining convergence rates for the suboptimality gap  $f(x) - f_*$  is not straightforward and requires taking into consideration the properties of the reference function  $\phi$  as well, as also noted in [12, p. 1006].

**Theorem II.4.** *Let Assumption II.1 hold true and  $x : [0, \infty) \rightarrow \mathbb{R}^n$  be the unique solution of (4) with initial condition  $x(0) = x_0$ . Let, moreover,  $f$  be convex. Then, we have the following:*

- (i)  $\phi^*(\nabla f(x(t)))$  is a decreasing function of  $t$ .
- (ii) The function  $V(t) := t\phi^*(\nabla f(x(t))) + f(x(t))$  is monotonically decreasing in  $[0, +\infty)$ . This implies that

$$\phi^*(\nabla f(x(t))) \leq \frac{f(x_0) - f_*}{t}, \quad \text{for } t > 0. \quad (8)$$

Assume moreover that  $\phi = h \circ \|\cdot\|$ , for some  $h : \mathbb{R} \rightarrow \overline{\mathbb{R}}$  satisfying Assumption II.1 and let  $x^* \in \arg \min f$ . Then,

- (iii)  $\|\nabla f(x(t))\|$  is a decreasing function of  $t$  and the same holds for  $\|x(t) - x^*\|$ .
- (iv) If, furthermore,  $h^{*'}(y)/y$  is a decreasing function on  $\mathbb{R}_+$ , the following holds for the suboptimality gap:

$$f(x(t)) - f_* \leq \frac{\|\nabla f(x_0)\| \|x_0 - x^*\|^2}{h^{*'}(\|\nabla f(x_0)\|)t}, \quad \text{for } t > 0. \quad (9)$$

*Proof.* “II.4(i)”: We have that

$$\begin{aligned} \frac{d}{dt}(\phi^*(\nabla f(x(t)))) &= \langle \nabla\phi^*(\nabla f(x(t))), \nabla^2 f(x(t))\dot{x}(t) \rangle \\ &\stackrel{(4)}{=} -\langle \dot{x}(t), \nabla^2 f(x(t))\dot{x}(t) \rangle \leq 0, \end{aligned} \quad (10)$$

since  $f$  is convex and thus  $\nabla^2 f \succeq 0$ .

“II.4(ii)”: Note that

$$\begin{aligned} \frac{d}{dt}V(t) &= \phi^*(\nabla f(x(t))) + t \frac{d}{dt}(\phi^*(\nabla f(x(t)))) \\ &\quad + \langle \nabla f(x(t)), \dot{x}(t) \rangle \end{aligned}$$

Due to (10),

$$\begin{aligned} \frac{d}{dt}V(t) &\leq \phi^*(\nabla f(x(t))) - \langle \nabla f(x(t)), \nabla\phi^*(\nabla f(x(t))) \rangle \\ &= -\phi(\nabla\phi^*(\nabla f(x(t)))) \end{aligned}$$

where we have used (4) and the equality case of the Fenchel–Young inequality [20, Prop. 11.3]. Therefore,  $V(x(t)) \leq V(x_0)$ , which means that

$$\phi^*(\nabla f(x(t))) \leq \frac{f(x_0) - f_*}{t}, \quad (11)$$

for all  $t > 0$ .

“II.4(iii)”: In light of [9, Lem. 1.3], we have that  $h^*$  is an increasing function on  $\mathbb{R}_+$ , while  $\phi^* = h^* \circ \|\cdot\|$ , and

$$\nabla\phi^*(y) = \frac{h^{*'}(\|y\|)}{\|y\|}y,$$

for all  $y \in \mathbb{R}^n \setminus \{0\}$  and 0 otherwise. Now, from II.4(i) we have that for all  $t_2 \geq t_1 \geq 0$ ,  $h^*(\|\nabla f(x(t_2))\|) \leq h^*(\|\nabla f(x(t_1))\|)$  and since  $h^*$  is increasing we obtain the claimed result. For the second claim,

$$\begin{aligned} \frac{d}{dt}(\tfrac{1}{2}\|x(t) - x^*\|^2) &= \langle x(t) - x^*, \dot{x}(t) \rangle \\ &= \langle x(t) - x^*, -\nabla\phi^*(\nabla f(x(t))) \rangle \\ &= -\frac{h^{*'}(\|\nabla f(x(t))\|)}{\|\nabla f(x(t))\|} \langle x(t) - x^*, \nabla f(x(t)) \rangle \\ &\leq \frac{h^{*'}(\|\nabla f(x(t))\|)}{\|\nabla f(x(t))\|} (f_* - f(x(t))) \leq 0, \end{aligned} \quad (12)$$

where the first inequality follows by the convex gradient inequality for  $f$  and the last one by the fact that  $f_* \leq f(\bar{x})$  for all  $\bar{x} \in \mathbb{R}^n$ .

“II.4(iv)”: Let  $W(t) = t \frac{h^{*'}(\|\nabla f(x_0)\|)}{\|\nabla f(x_0)\|} (f(x(t)) - f_*) + \frac{1}{2}\|x(t) - x^*\|^2$ . Using (12) and (4), we have that

$$\begin{aligned} \frac{d}{dt}W(t) &\leq -t \frac{h^{*'}(\|\nabla f(x_0)\|)}{\|\nabla f(x_0)\|} \langle \nabla f(x(t)), \nabla\phi^*(\nabla f(x(t))) \rangle \\ &\quad - \left( \frac{h^{*'}(\|\nabla f(x(t))\|)}{\|\nabla f(x(t))\|} - \frac{h^{*'}(\|\nabla f(x_0)\|)}{\|\nabla f(x_0)\|} \right) (f(x(t)) - f_*). \end{aligned}$$

We have that  $\|\nabla f(x(t))\| \leq \|\nabla f(x_0)\|$  from II.4(iii) and since  $h^{*'}(y)/y$  is decreasing,  $\frac{h^{*'}(\|\nabla f(x(t))\|)}{\|\nabla f(x(t))\|} \geq \frac{h^{*'}(\|\nabla f(x_0)\|)}{\|\nabla f(x_0)\|}$ . Therefore, since  $f(x(t)) \geq f_*$ ,  $W(t) \leq W(0)$  and the claimed result follows.  $\square$

*Remark II.5.* The conditions imposed on  $\phi$  in order to obtain convergence rates for  $f(x) - f_*$  are in fact mild: the assumption  $\phi = h \circ \|\cdot\|$  covers a wide variety of methods in the related literature and is often used in order to obtain convergence guarantees. The assumption that  $h^{*'}(y)/y$  is a decreasing function on  $\mathbb{R}_+$  actually covers a plethora of interesting reference functions as stressed in [9]. When  $h \in \mathcal{C}^2(\text{int dom } h)$ , it is equivalent to  $yh^{*''}(y) \leq h^{*'}(y)$  for all  $y > 0$ .

The standard gradient flow is known to converge exponentially in function values when  $f$  satisfies a growth condition known as the PL inequality:  $\|\nabla f(x)\|^2 \geq 2\mu(f(x) - f_*)$  for all  $x \in \mathbb{R}^n$  and some  $\mu > 0$ . In our general setting, where  $f$  might be nonconvex, the equivalent condition utilized to prove linear rates of convergence in the function values for the sequence generated by (3) is known as anisotropic gradient dominance [11, Def. 5.6]. It takes the following form: there exists some  $\mu > 0$  such that for all  $x \in \mathbb{R}^n$ ,

$$\phi(\nabla\phi^*(\nabla f(x))) \geq \mu(f(x) - f_*). \quad (13)$$

Clearly, choosing  $\phi = \frac{1}{2}\|\cdot\|^2$  we obtain the aforementioned PL inequality. Under this condition, we can show exponential convergence, as described in the following proposition.

**Proposition II.6.** *Let Assumption II.1 hold true and  $x : [0, \infty) \rightarrow \mathbb{R}^n$  be the unique solution of (4) with initial point  $x_0$ . Let, moreover, (13) hold. Then,*

$$f(x(t)) - f_* \leq \exp(-\mu t)(f(x_0) - f_*).$$

*Proof.* Combining (7) and (13) we have that

$$\frac{d}{dt}(f(x(t))) \leq -\mu(f(x(t)) - f_*).$$

The claimed result now follows from Grönwall's inequality.  $\square$

*Remark II.7.* Note that the normalized gradient flow studied in [6],  $\dot{x}(t) = -\frac{\nabla f(x(t))}{\|\nabla f(x(t))\|}$  can be brought into the generalized form of (4), where equality is replaced by inclusion,

$$\dot{x}(t) \in -\partial\phi^*(\nabla f(x(t)))$$

with  $\phi^* = \|\cdot\|$  [22, Ex. 3.34]. In this case, by standard convex conjugacy [20, Ex. 11.4] and [22, Ex. 2.31],  $\phi$  is the indicator of the unit norm ball.

### III. MIRROR DESCENT AND MIRROR FLOW

Throughout the remainder of this section we assume the following.

**Assumption III.1.**  $f$  is strictly convex and supercoercive.

Note that since  $f$  is also assumed smooth, it is Legendre,  $\nabla f : \mathbb{R}^n \rightarrow \text{int dom } f^*$  is bijective and the set  $\arg \min f$  is

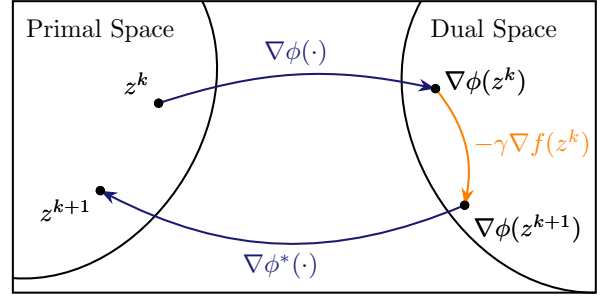


Fig. 2. Visualization of a mirror descent update.

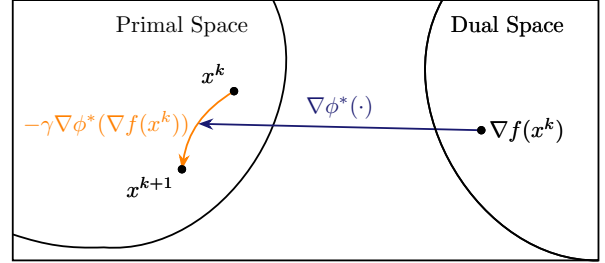


Fig. 3. Visualization of a nonlinearly preconditioned gradient method update.

a singleton, i.e.,  $\arg \min f = \{x^*\}$ . Moreover, due to the duality of supercoercivity and full domain [20, Thm. 11.8(d)],  $f^*$  has full domain.

*Remark III.2.* Strict convexity does not in general imply supercoercivity and in fact does not even imply simple coercivity as evidenced by the example  $f(x) = \exp(x)$ , which is a strictly convex function but  $\lim_{x \rightarrow -\infty} f(x) = 0$ . A standard example of strictly convex functions that are also supercoercive are strongly convex ones, something which can be seen from the strong convexity inequality. On the other hand, there exist strictly convex, supercoercive functions that are not strongly convex, with a standard example being  $f(x) = \|x\|^4$ .

In discrete-time, the algorithm described in (3) has a very interesting duality relation with the celebrated mirror descent method, which was explored in [23]. Consider the standard mirror descent update:

$$z^{k+1} = \nabla\phi^*(\nabla\phi(z^k) - \gamma\nabla f(z^k)),$$

where  $z^0 \in \text{int dom } \phi$ . Now, consider the iteration (3) and set  $z^k := \nabla f(x^k)$ . The latter is equivalent to  $x^k = \nabla f^*(z^k)$ . Substituting we obtain

$$\nabla f^*(z^{k+1}) = \nabla f^*(z^k) - \gamma\nabla\phi^*(z^k)$$

and thus  $z^{k+1} = \nabla f(\nabla f^*(z^k) - \gamma\nabla\phi^*(z^k))$ . It is thus clear that (3) is equivalent to applying the mirror descent algorithm on  $\phi^*$  with mirror potential  $f$ , i.e., swapping the roles of  $f$  and  $\phi^*$ . Nevertheless, the logic behind the two methods differs. In mirror descent, the primal iterate  $z^k$  is mapped to some dual iterate  $\nabla\phi(z^k)$ , a gradient step is performed and then it is mapped back to  $z^{k+1}$  through  $\nabla\phi^*$ . In the dual space preconditioning framework (3),  $\nabla f(x^k)$  is mapped to the primal space through  $\nabla\phi^*$  and then a gradient step is



performed directly on  $x^k$  to produce  $x^{k+1}$ . This difference is visualized in Figures 2 and 3.

A similar relation is present also in the continuous-time analog of the two methods. To better see this, consider the gradient flow described in [4, Eq. (4)]:

$$\dot{z}(t) = -\nabla f(\nabla\phi^*(z(t))),$$

with  $z(0) = z_0$ . Comparing the equation above with (4) it is clear that one can be obtained from the other by interchanging the roles of  $f$  and  $\phi^*$ . The described difference on the logic of the two methods is also evident in the continuous-time setup.  $z(t)$  is described as the dual trajectory, while  $y(t) = \nabla\phi^*(z(t))$  is the primal trajectory, the one that tends to the solutions of (5).

The aforementioned duality relation allows us to apply the results of [4] directly to the setting of the current paper. To that end, consider the point  $x^*$  and set  $y^* = \nabla f^*(x^*)$ . Now consider the following control system:

$$\dot{x}(t) = u(t), \quad (14)$$

the candidate Lyapunov function  $V(x) := D_f(x, y^*)$  and  $q(x, u) := \phi^*(\nabla f(x)) + \phi(-u) + \langle u, x^* \rangle$ . In the following,  $q$  will play the role of the optimal control cost that we want to minimize. Now note that  $V$  satisfies the properties listed in [4, Lem. 1] by definition, while

$$\begin{aligned} \frac{d}{dt} V(x(t)) &= \langle \nabla V(x(t)), u(t) \rangle \\ &= \langle \nabla f(x(t)) - x^*, u(t) \rangle, \end{aligned}$$

where we have used  $\nabla f(\nabla f^*(x^*)) = x^*$ , and thus  $\frac{d}{dt} V(x(t)) + q(x, u) \geq 0$  from the Fenchel–Young inequality. Therefore, denoting  $J_\infty(x_0, u) := \int_0^\infty q(x(t), u(t)) dt$  the cost we want to minimize over all stabilizing controls  $u$  (see [4, p. 1543]), i.e., functions  $u : [0, \infty) \rightarrow \mathbb{R}^n$  such that  $\nabla f(x(t))$  converges to the unique minimizer of  $f$ , we retrieve the following result [4, Thm. 1]:

**Theorem III.3.** *For any stabilizing control  $u$  and for all  $t \geq 0$ ,*

$$\int_0^t q(x(t), u(t)) dt \geq V(x_0) - V(x(t)).$$

*In particular,  $J_\infty(x_0, u) \geq V(x_0)$ . Moreover, for each  $x_0$ , the closed-loop system  $\dot{x}(t) = -\nabla\phi^*(\nabla f(x(t)))$  gives rise to an optimal stabilizing control  $u(t) = -\nabla\phi^*(\nabla f(x(t)))$ , such that*

$$J_\infty(x_0, u) = V(x_0) = D_f(x_0, y^*).$$

*Furthermore,  $V(x)$  is a global Lyapunov function for the closed-loop system.*

**Remark III.4.** The duality relation that is discussed throughout this section, allows us to obtain some convergence results directly from the analysis of the mirror descent method in [4]. To begin with, note that since  $D_{f^*}(\nabla f^*(x^*), \nabla f(x_0)) = D_f(x_0, x^*) = f(x_0) - f_*$ , by swapping  $f$  and  $\phi^*$  in [4, Thm. 2] (and noting that in that paper  $y(t) = \nabla f(x(t))$ ,  $y_0 = \nabla f(x_0)$ ) we get the result described in II.4(ii). Similarly, the

second item in [4, Thm. 2] describes exponential convergence under relative strong convexity of  $\phi^*$  w.r.t.  $f^*$ , i.e., the following inequality holding for all  $y, \bar{y} \in \mathbb{R}^n$ :

$$D_{\phi^*}(y, \bar{y}) \geq \mu D_{f^*}(y, \bar{y}).$$

Since this holds for all  $y, \bar{y} \in \mathbb{R}^n$ , it holds also for  $\bar{y} = \nabla f(x)$  and  $y = \nabla f(x^*) = 0$ ,  $x \in \mathbb{R}^n$ . Then,  $D_{f^*}(y, \bar{y}) = D_f(x, x^*)$  and  $D_{\phi^*}(y, \bar{y}) = \langle \nabla\phi^*(\nabla f(x)), \nabla f(x) \rangle - \phi^*(\nabla f(x)) = \phi(\nabla\phi^*(\nabla f(x)))$  and we can see that it implies (13). Nevertheless, the direct analysis of the previous section allows us to go into further detail and obtain tighter convergence guarantees under less restrictive assumptions. Note, moreover, that choosing  $u = -\nabla\phi^*(\nabla f(x))$  allows us to enforce constraints on the input by choosing a suitable reference function  $\phi$ , i.e., one with a compatible domain. For example,  $\phi(x) = 1/2\|x\|^2 + \delta_{[0,1]}(\|x\|)$ , leads to  $u$  taking values in  $B(0, 1) = \{x \in \mathbb{R}^n \mid \|x\| \leq 1\}$  and in this case  $\nabla\phi^*$  plays the role of the projection on  $B(0, 1)$ .

#### IV. CONCLUSION

In this paper, we analyzed nonlinearly preconditioned gradient flows, which arise as the continuous-time analogue of the nonlinearly preconditioned gradient method. By means of a novel Lyapunov-like function, we established a standard sublinear convergence rate in the convex setting, and proved exponential convergence under a generalized PL inequality. Moreover, we described a duality connection with mirror flows, which allows existing results for the mirror descent method to be directly transferred. Interesting future work involves extending our analysis to second-order equations in the spirit of [1] and studying the differential equation under the lens of  $\Phi$ -convexity similar to [11].

#### REFERENCES

- [1] H. Attouch, X. Goudou, and P. Redont, “The heavy ball with friction method, I. the continuous dynamical system: global exploration of the local minima of a real-valued function by asymptotic analysis of a dissipative dynamical system,” *Communications in Contemporary Mathematics*, vol. 2, no. 01, pp. 1–34, 2000.
- [2] W. Su, S. Boyd, and E. J. Candes, “A differential equation for modeling Nesterov’s accelerated gradient method: Theory and insights,” *Journal of Machine Learning Research*, vol. 17, no. 153, pp. 1–43, 2016.
- [3] W. Krichene, A. Bayen, and P. L. Bartlett, “Accelerated mirror descent in continuous and discrete time,” *Advances in neural information processing systems*, vol. 28, 2015.
- [4] B. Tzen, A. Raj, M. Raginsky, and F. Bach, “Variational principles for mirror descent and mirror Langevin dynamics,” *IEEE Control Systems Letters*, vol. 7, pp. 1542–1547, 2023.
- [5] H. Attouch and J. Peypouquet, “Convergence of inertial dynamics and proximal algorithms governed by maximally monotone operators,” *Mathematical Programming*, vol. 174, no. 1, pp. 391–432, 2019.
- [6] J. Cortés, “Finite-time convergent gradient flows with applications to network consensus,” *Automatica*, vol. 42, no. 11, pp. 1993–2000, 2006.
- [7] F. Hübler, I. Fatkhullin, and N. He, “From gradient clipping to normalization for heavy tailed SGD,” in *International Conference on Artificial Intelligence and Statistics*, pp. 2413–2421, PMLR, 2025.
- [8] J. Zhang, T. He, S. Sra, and A. Jadbabaie, “Why gradient clipping accelerates training: A theoretical justification for adaptivity,” in *International Conference on Learning Representations*.
- [9] K. Oikonomidis, J. Quan, E. Laude, and P. Patrinos, “Nonlinearly preconditioned gradient methods under generalized smoothness,” in *Forty-second International Conference on Machine Learning*.
- [10] K. Oikonomidis, J. Quan, and P. Patrinos, “Nonlinearly preconditioned gradient methods: Momentum and stochastic analysis,” *arXiv preprint arXiv:2510.11312*, 2025.

- [11] E. Laude and P. Patrinos, "Anisotropic proximal gradient," *Mathematical Programming*, pp. 1–45, 2025.
- [12] C. J. Maddison, D. Paulin, Y. W. Teh, and A. Doucet, "Dual space preconditioning for gradient descent," *SIAM Journal on Optimization*, vol. 31, no. 1, pp. 991–1016, 2021.
- [13] A. Bodard and P. Patrinos, "Escaping saddle points without Lipschitz smoothness: the power of nonlinear preconditioning," *arXiv preprint arXiv:2509.15817*, 2025.
- [14] U. Stefanelli, "The Brezis–Ekeland principle for doubly nonlinear equations," *SIAM Journal on Control and Optimization*, vol. 47, no. 3, pp. 1615–1642, 2008.
- [15] H. H. Bauschke and P. L. Combettes, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer, 2017.
- [16] A. Mielke, R. Rossi, and G. Savaré, "Nonsmooth analysis of doubly nonlinear evolution equations," *Calculus of Variations and Partial Differential Equations*, vol. 46, no. 1, pp. 253–310, 2013.
- [17] E. Emmrich and M. Thalhammer, "Doubly nonlinear evolution equations of second order: Existence and fully discrete approximation," *Journal of Differential Equations*, vol. 251, no. 1, pp. 82–118, 2011.
- [18] G. A. Bonaschi and M. A. Peletier, "Quadratic and rate-independent limits for a large-deviations functional," *Continuum Mechanics and Thermodynamics*, vol. 28, no. 4, pp. 1191–1219, 2016.
- [19] R. T. Rockafellar, *Convex analysis*, vol. 28. Princeton university press, 1997.
- [20] R. T. Rockafellar and R. J. Wets, *Variational Analysis*. New York: Springer, 1998.
- [21] R. T. Rockafellar, "Higher derivatives of conjugate convex functions," *Int. J. Applied Analysis*, no. 1, pp. 41–43, 1977.
- [22] A. Beck, *First-order methods in optimization*. SIAM, 2017.
- [23] J. Kim, C. Park, A. Ozdaglar, J. Diakonikolas, and E. K. Ryu, "Mirror duality in convex optimization," *arXiv preprint arXiv:2311.17296*, 2023.