

# Optimizing Optimizers

## Regret-optimal gradient descent algorithms

**Conference Paper****Author(s):**

Casgrain, Philippe; [Kratsios, Anastasis](#) 

**Publication date:**

2021

**Permanent link:**

<https://doi.org/10.3929/ethz-b-000490896>

**Rights / license:**

[In Copyright - Non-Commercial Use Permitted](#)

**Originally published in:**

Proceedings of Machine Learning Research 134

# Optimizing Optimizers: Regret-optimal gradient descent algorithms.

Philippe Casgrain

PHILIPPE.CASGRAIN@MATH.ETHZ.CH

Anastasis Kratsios

ANASTASIS.KRATSIOS@MATH.ETHZ.CH

ETH Zürich

**Editors:** Mikhail Belkin and Samory Kpotufe

## Abstract

This paper treats the task of designing optimization algorithms as an optimal control problem. Using regret as a metric for an algorithm’s performance, we study the existence, uniqueness and consistency of regret-optimal algorithms. By providing first-order optimality conditions for the control problem, we show that regret-optimal algorithms must satisfy a specific structure in their dynamics which we show is equivalent to performing *dual-preconditioned gradient descent* on the value function generated by its regret. Using these optimal dynamics, we provide bounds on their rates of convergence to solutions of convex optimization problems. Though closed-form optimal dynamics cannot be obtained in general, we present fast numerical methods for approximating them, generating optimization algorithms which directly optimize their long-term regret. These are benchmarked against commonly used optimization algorithms to demonstrate their effectiveness.

**Keywords:** meta-optimization, non-convex optimization, optimal control, variational optimization, algorithm generation, hyperparameter optimization, convex optimization, regret

## 1. Introduction

Let  $\mathcal{X} = \mathbb{R}^d$  and consider the unconstrained minimization problem

$$\text{Min}_{x \in \mathcal{X}} f(x), \tag{1}$$

for an objective function  $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{\infty\}$  which satisfies the following regularity assumptions.

**Assumption 1** *We assume that the optimization problem (1) is non-degenerate, in the sense that there exists a minimizer  $x^* \in \text{argmin}_{x \in \mathcal{X}} f(x) \subseteq \mathcal{X}$  for which  $|f(x^*)| < \infty$ .*

This paper considers the problem of optimally selecting amongst a class of optimization algorithms, by minimizing a fixed performance metric of their path. In this sense, we are concerned with a *meta-optimization problem* in which we are optimizing over algorithms, which in turn optimize  $f$ .

We identify algorithms with the paths they take within the optimization domain  $\mathcal{X}$ . More precisely, we define an algorithm  $\mathbf{x}$  as the sequence of points  $\mathbf{x} = \{x_t\}_{t \in \mathbb{N}} \in \mathcal{X}^{\mathbb{N}}$ , where  $\mathcal{X}^{\mathbb{N}}$  is the space of sequences on  $\mathcal{X}$ . Following this notation, we introduce the set of algorithms initialized at  $x \in \mathcal{X}$  which terminate by iteration  $T \in \mathbb{N}$  as

$$\mathcal{A}_x^T := \left\{ \mathbf{x} \in \mathcal{X}^{\mathbb{N}} : x_0 = x \text{ and } \forall u \geq T, \Delta x_u = 0 \right\},$$

where we use the notation  $\Delta x_t = x_{t+1} - x_t$  to represent the increments of  $\mathbf{x}$ . We also define its asymptotic counterpart  $\mathcal{A}_x^\infty := \{\mathbf{x} \in \mathcal{X}^{\mathbb{N}} : x_0 = x\}$ , the set of all sequences with fixed initial point. When necessary, we denote the union of these sets over all initial points as  $\mathcal{A}^T := \bigcup_{x \in \mathcal{X}} \mathcal{A}_x^T$ . Over

the course of the paper, we use the convention that bold symbols  $\mathbf{y} \in \mathcal{X}^{\mathbb{N}}$  represent algorithms and their un-bolded counterparts  $y_t \in \mathcal{X}$  represent their value at a fixed iteration  $t \in \mathbb{N}$ .

Our focus will be on those algorithms which successfully approximate  $x^*$  in the limit, i.e. for which  $x_t \rightarrow x^*$ . As previously stated, we seek the algorithms which achieve this while minimizing a measure of performance along the path they take. We define this measure in such a way as to represent both the speed of convergence of  $x$  with respect to the optimization problem (1) and the *stability* of the path that the algorithm traces over a fixed horizon. Hence, we introduce the *regret* of an algorithm  $x$  as the  $\mathcal{R}_T(x) : \mathcal{X}^{\mathbb{N}} \rightarrow [0, \infty]$ , given by

$$\mathcal{R}_T(x) = \sum_{t=1}^T f(x_t) - f_* + \phi(\Delta x_{t-1}), \quad (2)$$

as our measure of an algorithm's performance. In the above definition, we assume that  $\phi : \mathcal{X} \rightarrow [0, \infty)$  satisfies the following.

**Assumption 2** Assume that  $\phi(0) = 0$ ,  $\phi$  is lower semi-continuous and satisfies the growth condition that  $c\|\cdot\|^p \leq \phi(\cdot)$  for some  $c > 0$  and  $p \geq 1$ , where  $\|\cdot\|$  is the Euclidean norm.

We interpret  $\mathcal{R}_T$  as measuring performance based on two distinct criteria. The first component of  $\mathcal{R}_T$  measures the cumulative distance to optimality through the sum of the terms  $f(x_t) - f_*$ , while the second measures total *path energy* of  $x$  through the sum of the terms  $\phi(\Delta x_{t-1})$ , which we can interpret as a generalization of the notion of its  $p$ -variation<sup>1</sup>. The definition (2) is related to the widely used notion of *adversarial regret* which is the central metric of algorithmic performance in the field of online learning<sup>2</sup>. The definition (2) is also related to the notion of *regularized regret*, which is widely used in the literature on ‘adaptive’ optimization algorithms (e.g. see Xiao (2010); Duchi et al. (2011a,b)), where the main difference lies in that we regularize over the increments  $\Delta x_t$ , rather than the positions  $x_t$ . We note, however, that the definition (2) differs from these related notions in that it is not adversarial since  $f$  remains fixed.

Rather than simply measuring the performance of  $x$  with respect to its last iteration,  $f(x_T) - f_*$ ,  $\mathcal{R}_T$  measures the average performance of  $x$  over  $T$  iterations. This average-iterate formulation, when compared to a last-iterate measure, yields a more stable performance metric as we scale the number of allowed iterations,  $T$ . For example, in the limit, we have that  $\tilde{\mathcal{R}}(x) = \lim_T f(x_T) - f_* \equiv 0$  is the zero functional for all convergent  $x$ , whereas  $\mathcal{R}_\infty(x) = \lim_T \mathcal{R}_T(x)$  is a non-trivial functional which ranks convergent algorithms according to their asymptotic rate of convergence and path stability.

We are interested in algorithms which are optimal with respect to  $\mathcal{R}_T$ . Hence, for each  $T \in \mathbb{N} \cup \{\infty\}$  we define the optimal control problem

$$\mathcal{P}_x^T := \text{Min}_{x \in \mathcal{A}_x^T} \mathcal{R}_T(x), \quad (3)$$

where we use the notation  $x \in \mathcal{P}_x^T$  to represent an element from the set of minimizers of (3). For  $T < \infty$ , elements of  $\mathcal{P}_x^T$  represent algorithms with fixed starting at a point  $x \in \mathcal{X}$ , terminating at iteration  $T$ , which minimize the performance metric  $\mathcal{R}_T$ . Extending previous notation, we also introduce  $\mathcal{P}^T := \bigcup_{x \in \mathcal{X}} \mathcal{P}_x^T$  as the set of solutions from all initial values in  $\mathcal{X}$ . For any  $T \in \mathbb{N} \cup \{\infty\}$ , we say that an algorithm  $x \in \mathcal{P}^T$  is *regret-optimal*.

1. We recover the  $p$ -variation for  $p \geq 1$  whenever  $\phi(x) = \|x\|_p^p$ , where  $\|\cdot\|_p$  is the  $p$ -norm on  $\mathcal{X}$ .  
 2. We refer the reader to the text Hazan (2016) for a comprehensive introduction to the topic.

## Structure and Summary of Main Contributions

The paper is devoted to the design of optimization algorithms which optimize regret. We provide an in-depth theoretical study of *regret-optimal* algorithms including existence and representation theorems, as well as convergence rate guarantees. As a consequence of the theory, we also put forward a meta-optimization algorithm which dynamically adapts to data, which can be applied to automatically tune the coefficients of parametric optimization algorithms. To the authors’ knowledge, this is the first paper to provide a rigorous and in-depth theoretical study of the automated design of algorithms which optimize their long-term regret.

In Section 2, we characterize the existence of regret-optimal algorithms in the general non-smooth and non-convex setting, as well as their consistency across regret horizons  $T \in \mathbb{N} \cup \{\infty\}$ . In Section 3 we study regret-optimality in the setting of differentiable objectives  $f$ , where we derive necessary conditions on their dynamics. We furthermore show that regret-optimal algorithms admit a representation as performing *dual-preconditioned gradient descent* (Maddison et al., 2019) on their value function. Section 4 studies regret-optimality in the context of convex and differentiable objectives. This section culminates in providing a hierarchy of convergence rate bounds for regret-optimal algorithms under varying relative-smoothness and relative-convexity assumptions on the objective function  $f$ , which are presented in Table 1. As a consequence of the theoretical analysis of the previous sections, Section 5 presents an online meta-algorithm for the purpose of learning regret-optimal algorithms from gradient and function data, where we apply this algorithm to the special case of the automated tuning of algorithm hyper-parameters on a host of toy optimization problems.

### 1.1. Related Work

The ideas of this paper are most closely related to the various variational interpretations of optimization. In particular, we highlight Wibisono et al. (2016); Casgrain (2019), which study algorithms which are critical points of an energy functional in a continuous-time setting and their connection to gradient descent algorithms with momentum. We argue that main differences between these and the present work is that we consider the former’s approach *ad hoc*; the variational framework in the former is chosen *a posteriori* to generate momentum-like dynamics, rather than chosen *a priori* to represent a concrete metric of algorithmic performance. Moreover, there are the related works of Betancourt et al. (2018); Shi et al. (2019); Wilson et al. (2019); França et al. (2020) which bring the continuous analysis over to the discrete-time setting through symplectic integration methods. In contrast, our analysis deals with the discrete-time optimization problem from the very beginning without the need for supplementary discretization machinery.

This paper is also related to the body of work on control-theoretic and dynamical systems models of optimization. Of note are Lessard et al. (2016); Hu et al. (2017); Muehlebach and Jordan (2020) which present control-theoretic interpretations of the evolution dynamics of optimization algorithms. These serve to analyze their rate of convergence to optima as well as establish various other stability properties. Though these approaches are control-theoretic, they differ from our approach since they are not concerned with *optimal control*, as they do not seek controls which are optimal with respect to a fixed performance functional. Rather, they take a control as given, and study the convergence of the resulting dynamical system.

The “meta-optimization of optimizers” philosophy used in this paper has also been studied from non control-based points of view. In particular, Taylor et al. (2017); Taylor and Drori (2021); Drori

and Taylor (2021) study the problem of designing algorithms with linear dynamics which maximize their worst-case performance on convex optimization problems. These papers show that this particular meta-optimization problem can be reduced to a semi-definite program which admits approximate solutions in certain simple cases. Taking a more applied perspective, Mitsos et al. (2018) consider automatic optimization algorithm generation by training a parametric algorithm over curated examples. Wichrowska et al. (2017) take a similar approach where algorithms are parametrized by neural networks whose weights are learned by training on a fixed corpus of problems. Hyper-parameter tuning methods such as in Lorraine and Duvenaud (2018), which search for optima in the set of algorithm hyper-parameters, can also be interpreted as trying to solve a finite-dimensional version of the meta-optimization problem. We also mention Li et al. (2018), who propose a continuous-time algorithm tailored to continuous-depth feedforward networks inspired by optimal control techniques.

## 1.2. Notation, Definitions and Conventions

For a Banach space  $\mathcal{Y}$ , the dual space  $\mathcal{Y}^*$  represents the space of continuous linear functionals on  $\mathcal{Y}$ . We say that  $\{y_i\}_{i \in \mathbb{N}} \subset \mathcal{Y}$  converges weakly to  $y_\infty \in \mathcal{Y}$ , which we denote as  $y_i \rightharpoonup y_\infty$ , if  $\ell(y_i) \rightarrow \ell(y_\infty)$  for all  $\ell \in \mathcal{Y}^*$ . For a convex function  $g : \mathcal{X} \rightarrow \mathbb{R}$  we define its convex dual  $g^* : \mathcal{X}^* \rightarrow \mathbb{R}$  as  $g^*(p) = \sup_{x \in \mathcal{X}} \{p(x) - g(x)\}$ . For a convex and differentiable function  $g$  and points  $x, y \in \mathcal{X}$  we define the *Bregman divergence* as  $D_g(x, y) = g(x) - g(y) - \langle \nabla g(y), x - y \rangle$  which is non-negative due to the convexity of  $g$ . For functions  $g, h$  and  $\mu > 0$ , we say that a function  $g$  is  $\mu$ -relatively-convex with respect to  $h$  if  $g - \mu h$  is convex. Conversely, for  $L > 0$ , we say that  $g$  is  $L$ -relatively-smooth with respect to  $h$  if  $Lh - g$  is convex. We say that a function  $g : \mathcal{X} \rightarrow \mathbb{R}$  is a positive-definite quadratic function if there exists a symmetric bi-linear form  $L : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  such that  $g(x) = L(x, x)$  and  $L(x, y) > 0$  for all  $0 \neq x, y \in \mathcal{X}$ .

Consider the functional  $F : U \rightarrow \mathbb{R}$  on a non-empty open subset  $U$  of a Banach space  $\mathcal{X}$ . The Gâteaux derivative of  $F$  at  $x \in U$  in the direction  $dx \in \mathcal{X}$  is the limit  $F'(x)(dx) = \lim_{\epsilon \rightarrow 0} \frac{F(x+\epsilon dx) - F(x)}{\epsilon}$ , which can be interpreted as a directional derivative. We refer the reader to (Ekeland and Temam, 1999, §5.2) for more in-depth coverage of the topic and its role in smooth convex and variational analysis.

We say a function  $g : \mathcal{X} \rightarrow \mathbb{R}$  is locally Lipschitz continuous if for every  $x \in \mathcal{X}$ , there exists a compact set  $K$  with non-empty interior and  $L_K > 0$  such that  $x \in K \subset \mathcal{X}$  and  $|g(y) - g(z)| \leq L_K \|y - z\|$  for all  $y, z \in K$ . For any locally Lipschitz function  $g$ , we define the Clarke directional derivative at  $x \in \mathcal{X}$  in a direction  $v \in \mathcal{X}$  as  $g^\circ(x; v) = \limsup_{y \rightarrow x, t \downarrow 0} \frac{f(y+tv) - f(y)}{t}$  and the generalized gradient as the set  $\partial g(x) = \{\zeta : g^\circ(x; v) \geq \langle \zeta, v \rangle\}$ . If  $g$  is convex then this definition coincides with its subgradient, if  $g$  is differentiable then  $\partial g(x)$  is a singleton containing the classical gradient, and if  $x$  is a minimum of  $g$  then  $0 \in \partial g(x)$ . We point the interested reader to Ferrera (2013), which covers these and other concepts of non-smooth analysis in full detail.

All proofs for theorems, lemmas and corollaries found throughout the paper are relegated to the paper's appendix. As a rule of thumb, all numbered assumptions found within the text are assumed to hold for the remainder of the paper, any other additional assumptions will be explicitly stated in the theorems, lemmas and corollaries that require them.

## 2. Existence and Time-Consistency

We begin by demonstrating the existence of regret optimal algorithms for a finite time horizon  $T \in \mathbb{N}$ , which we show are guaranteed to exist under the mild conditions put forth in Assumptions 1 and 2.

**Theorem 1** *For all  $x \in \mathcal{X}$  and  $T \in \mathbb{N}$ , the set of minima,  $\mathcal{P}_x^T$ , is non-empty.*

Although the control problem (3) enjoys increased analytical tractability when  $T < \infty$ , we are also interested in the case when  $T = \infty$  since the latter admits solutions which are invariant to the iteration number,  $t$ . In order to precisely characterize the relationship between the solutions in the finite and infinite-horizon regimes in Theorem 4, we must first introduce additional notions of regularity on the set of algorithms. For this reason, for each  $\alpha \geq 0$ , we introduce the set of  $\alpha$ -stable algorithms, which we define as

$$\mathcal{A}_x^{\infty:\alpha} := \left\{ \mathbf{x} \in \mathcal{X}^{\mathbb{N}} : x_0 = x \text{ and } \sum_{u \in \mathbb{N}} u^\alpha \|\Delta x_u\|^p < \infty \right\},$$

where  $p > 0$  is the value found in Assumption 2. This set can be loosely interpreted as the set of  $\mathbf{x} \in \mathcal{A}_x^\infty$  for which the increments  $\|\Delta x_t\|^p$  asymptotically decay to zero at a rate  $O(t^{-(1+\alpha)})$ . We also note that the definition above clearly implies that  $\mathcal{A}_x^{\infty:\alpha_1} \subset \mathcal{A}_x^{\infty:\alpha_0} \subset \mathcal{A}_x^\infty$  for any  $0 \leq \alpha_0 \leq \alpha_1$ . Following the above definition, we also define the corresponding optimization problem  $\mathcal{P}_x^{\infty:\alpha} := \operatorname{argmin}_{\mathbf{x} \in \mathcal{A}_x^{\infty:\alpha}} \mathcal{R}_\infty(\mathbf{x})$ . In the theorem that follows, we show that the infinite horizon control problem is well-posed and admits solutions.

**Theorem 2** *Let  $x \in \mathcal{X}$ , and  $\alpha \geq 0$ , then  $\mathcal{P}_x^{\infty:\alpha}$  is non-empty. In the case where  $\alpha = 0$ , we have that  $\mathcal{P}_x^\infty = \mathcal{P}_x^{\infty:0}$ , and hence,  $\mathcal{P}_x^\infty$  is also non-empty. Lastly, all solutions  $\mathbf{x} \in \mathcal{P}_x^{\infty:\alpha} \cup \mathcal{P}_x^\infty$  exhibit finite regret, so that  $\mathcal{R}_\infty(\mathbf{x}) < \infty$ .*

**Corollary 3** *For any  $\mathbf{x} \in \mathcal{P}_x^{\infty:\alpha}$  or  $\mathbf{x} \in \mathcal{P}_x^\infty$  we have that  $f(x_t) - f_\star + \phi(\Delta x_t) = o(1)$ . If this sequence is monotone, then we also have that  $f(x_t) - f_\star + \phi(\Delta x_t) = o(\frac{1}{t})$ .*

One important consequence of Theorem 2 and Lemma 36 is that regret-optimal algorithms  $\mathbf{x} \in \mathcal{P}_x^{\infty:\alpha}$  or  $\mathbf{x} \in \mathcal{P}_x^\infty$  must exhibit finite regret, and hence form non-trivial solutions. Moreover, Corollary 3 also shows that these algorithms always satisfy  $f(x_t) \rightarrow f_\star$ , with an asymptotic upper bound on their rate of convergence, provided that  $\{f(x_t)\}_{t \in \mathbb{N}}$  is monotone decreasing. Another consequence is that we have the equivalence between the constrained ( $\mathcal{P}_x^{\infty:0}$ ) and unconstrained ( $\mathcal{P}_x^\infty$ ) solution sets in the  $T = \infty$  regime, yielding the regularity property that  $\sum_{t=0}^\infty \|\Delta x_t\|^p < \infty$  for any  $\mathbf{x} \in \mathcal{P}_x^\infty$ . Regret-optimal algorithms also exhibit a *time-consistency* property across their horizon,  $T$ , which we present in the next theorem.

**Theorem 4** *Let  $x \in \mathcal{X}$ . An algorithm belongs to  $\mathcal{P}_x^{\infty:\alpha}$  if and only if there exists a sequence  $\{\mathbf{x}^n\}_{n \in \mathbb{N}}$  such that  $\mathbf{x}^n \in \mathcal{P}_x^n$  and a subsequence  $\{\mathbf{x}^{n_k}\}_{k \in \mathbb{N}}$  satisfying one of the following conditions.*

1. *If  $\alpha = 0$ , then  $\mathbf{x}^{n_k} \rightharpoonup \mathbf{x}^\infty$  in the weak topology of  $\mathcal{A}_x^{\infty:0}$ .*
2. *If  $\alpha > 0$ , then  $\lim_{k \rightarrow \infty} \sum_{t=1}^\infty \|\mathbf{x}_t^{n_k} - \mathbf{x}_t^\infty\|^p = 0$ .*

Hence, Theorem 4 shows that solutions to the infinite horizon control problem can be represented as the limit of solutions in  $\mathcal{P}_x^T$ , providing another avenue for computation in the  $T = \infty$  regime. Moreover, we find that the required stability level  $\alpha$  dictates the mode of convergence, where we recall that by Theorem 2, since  $\mathcal{P}^{\infty:0} = \mathcal{P}^\infty$ , the statement of Theorem 4-2 holds for the unconstrained problem as well.

### 3. Optimal Dynamics

A natural object of study in the context of optimal control are first-order optimality criteria for critical points of an objective. In the following section, we carry out the analysis of critical points for the regret optimization problem posed in Section 1. In order to carry out this analysis, however, we require smoothness of the control problem. As such, the focus of the remainder of the paper will be the optimization of smooth objectives.

**Assumption 3** *Assume that the following assumptions hold for the remainder of the paper.*

1.  $f$  is everywhere differentiable.
2.  $\phi$  is Legendre convex. That is,  $\phi$  is everywhere finite, strictly convex, differentiable and satisfies the super-coercivity condition that  $\lim_{\|x\| \rightarrow \infty} \|\nabla\phi(x)\| = \infty$ .

Recall that the Legendre convexity condition in Assumption 3 ensures that both  $\phi$  and  $\phi^*$  are strictly convex, differentiable and satisfy the property  $\nabla\phi^*(\nabla\phi(x)) = x$  for all  $x \in \mathcal{X}$ . We refer interested readers to (Rockafellar, 1970, Section 26) for more information on Legendre convex functions and their properties.

We begin our analysis of critical points by computing the Gâteaux derivative of  $\mathcal{R}_T(\mathbf{x})$  over  $\mathcal{A}_x^T$ . Letting the derivative vanish, we find that the dynamics of critical points must satisfy a very specific structure which we present below.

**Theorem 5** *For  $T \in \mathbb{N} \cup \{\infty\}$  and  $\mathbf{x} \in \mathcal{A}_x^T$ , consider the linear functional  $\mathcal{R}' \in (\mathcal{A}_x^T)^*$  defined by*

$$\mathcal{R}'_T(\mathbf{x})(\delta\mathbf{x}) = \sum_{t=1}^T \langle \nabla\phi(\Delta x_{t-1}) - \nabla\phi(\Delta x_t) + \nabla f(x_t), \delta x_{t-1} \rangle. \quad (4)$$

*If  $T \in \mathbb{N}$ ,  $\mathcal{R}'_T$  is the Gâteaux derivative of  $\mathcal{R}_T(\mathbf{x})$  over  $\mathcal{A}_x^T$ .*

**Theorem 6** *For any  $T \in \mathbb{N} \cup \{\infty\}$ , define  $\hat{\mathcal{P}}_x^T \subseteq \mathcal{A}_x^T$  as the set of algorithms  $\mathbf{x} \in \mathcal{A}_x^T$  which satisfy the difference equation*

$$\nabla\phi(\Delta x_t) - \nabla\phi(\Delta x_{t-1}) = \nabla f(x_t) \quad \forall t \leq T. \quad (5)$$

*For  $T < \infty$  we have the two properties that*

- i.  $\hat{\mathcal{P}}_x^T$  is precisely the set of critical points of  $\mathcal{R}_T$ , and hence  $\hat{\mathcal{P}}_x^T \supseteq \mathcal{P}_x^T$ .
- ii. Let  $\mathbf{x} \in \hat{\mathcal{P}}_x^T$ , and for  $h \in \mathbb{N}$ , define the truncation  $\mathbf{x}_{\rightarrow h} = \{x_{u+h}\}_{u \geq t}$ . If  $T \in \mathbb{N}$ , then for any  $0 \leq t < T$  we have the recursive property that  $\mathbf{x}_{\rightarrow t} \in \hat{\mathcal{P}}_{x_t}^{T-t}$ .

Theorem 6 therefore provides a characterization of critical points of  $\mathcal{R}_T$  in terms of the difference equation (5). Readers familiar with optimal control theory can also interpret (5) as the weak Pontryagin maximum principle<sup>3</sup> for the control problem (3), where  $\nabla\phi(\Delta x_t)$  fills the role of what is known as the *co-state process* in optimal control and *momentum* in (discrete) classical mechanics.

Writing the explicit solution to the dynamics (5), we obtain that  $\mathbf{x} \in \hat{\mathcal{P}}^T$  satisfies

$$\nabla\phi(\Delta x_t) = - \sum_{u=t+1}^T \nabla f(x_u) ,$$

which can be loosely interpreted as implying that the dynamics of  $x \in \mathcal{P}^T$  are decelerating when  $T \in \mathbb{N}$ , since the number of items within the sum shrinks at each iteration. It also happens that the optimality dynamics (5) admit another important interpretation in relation to the *value function*. We present the results relevant to this representation below.

**Theorem 7** For  $T \in \mathbb{N} \cup \{\infty\}$ , define the value function  $x \mapsto J^T(x) := \min_{\mathbf{y} \in \mathcal{A}_x^T} \mathcal{R}_T(\mathbf{y})$  over  $x \in \mathcal{X}$ . We assume one of the following.

1. If  $T < \infty$ , assume that for each  $0 < t \leq T$ ,  $J^t$  is locally Lipschitz-continuous.
2. If  $T = \infty$ , assume that  $J^\infty$  is locally Lipschitz-continuous.

Denoting  $\partial J^T(x)$  as the Clarke generalized gradient of  $J^T$ , for any  $x \in \mathcal{P}^T$  and  $t < T$  we have that

$$\begin{aligned} -\nabla\phi(x_{t+1} - x_t) &\in \partial J^{T-t}(x_t) , \quad \text{and} \\ \partial J^{T-t}(x_t) &\subseteq \partial J^{T-(t+1)}(x_{t+1}) + \nabla f(x_{t+1}) . \end{aligned} \tag{6}$$

Hence, if  $T = \infty$ , it is easy to see that under the assumptions of Theorem 7, equation (6) implies that any  $x \in \mathcal{P}^\infty$  will satisfy the optimality dynamics of Theorem 6 (eq. (5)). Therefore, we have a result analogous to Theorem 6-1 that  $\hat{\mathcal{P}}^\infty \supseteq \mathcal{P}^\infty$ , and hence the dynamics of equation 5 are a necessary condition for optimality in the  $T = \infty$  regime.

In order to better understand Theorem (7), we remark that since  $\phi$  is Legendre convex, under the assumptions of Theorem (7) we have the representation

$$x_{t+1} = x_t - \nabla\tilde{\phi}^*(\nu_t) \quad \text{where } \nu_t \in \partial J^{T-t}(x_t) , \tag{7}$$

for the iterates of  $x \in \hat{\mathcal{P}}_x^T$ , where we define  $\tilde{\phi}(x) = \phi(-x)$ . We can therefore interpret  $x \in \hat{\mathcal{P}}^T$  as performing a variant of gradient descent on the generalized gradient of  $J^{T-t}(x)$ . More specifically, this variant of gradient descent happens to generalize *dual-preconditioned gradient descent* (Maddison et al., 2019, Algorithm 1.1). This interpretation will be particularly important in obtaining convergence bounds in Section 4, where their connection becomes more clear.

In the case where  $T = \infty$ , Theorem 7 also implies that any  $x \in \mathcal{P}_x^\infty$  admits a map  $\nu : \mathcal{X} \rightarrow \partial J^\infty(\mathcal{X}) \subseteq \mathcal{X}$  such that

$$x_{t+1} = x_t - \nabla\tilde{\phi}^*(\nu(x_t)) \quad \text{and } \nu(x_t) = \nu(x_{t+1}) + \nabla f(x_{t+1}) \tag{8}$$

for all  $t \in \mathbb{N}$ . Hence the dynamics of such an  $x \in \mathcal{P}^\infty$  can be uniquely represented by a vector field  $\nu$  which is independent of the iteration number  $t$ .

3. We point interested readers to Blot and Hayek (2014) for detailed information on the strong and weak Pontryagin Maximum principle and their role in optimal control.



## 4. Convex Optimization

Over the course of this section, we study the regret optimization problem in the case where  $f$  is convex. In particular, we will focus on the convergence of asymptotically regret-optimal algorithms  $\mathbf{x} \in \mathcal{P}_x^\infty$  to solutions of the optimization problem on  $f$ . We begin by establishing some essential convexity properties of the control problem that arise as a result.

**Lemma 8** *Assume that  $f$  is convex. Then for all  $T \in \mathbb{N} \cup \{\infty\}$ ,  $\mathcal{R}_T(\mathbf{x})$  is a strictly convex functional of  $\mathcal{A}_x^T$  and hence,  $\mathcal{P}_x^T$  is a non-empty singleton and  $\mathcal{P}_x^T = \hat{\mathcal{P}}_x^T$ .*

Lemma 8 therefore implies that the optimality dynamics of Theorem 6 or 7 are both necessary and sufficient conditions of optimality in the context of a convex control problem. Hence, any  $\mathbf{x} \in \mathcal{A}_x^T$  satisfying these dynamics is guaranteed to be the unique solution to the regret minimization control problem.

The assumption that  $f$  is convex also has numerous consequences in terms of the convergence rates of regret-optimal algorithms. We study these from the perspective of the *value function*  $J^T(x) := \min_{\mathbf{x} \in \mathcal{A}_x^T} \mathcal{R}_T(\mathbf{x})$ . We note here that for each  $x \in \mathcal{X}$ , Lemma 8 states that there is a unique  $\mathbf{x} \in \mathcal{P}_x^T$  such that  $J^T(x) = \mathcal{R}_T(\mathbf{x})$ . As is hinted to by Lemma 7 and the discussion that follows, we will see that this function has an important connection with the optimality dynamics of Theorem 6. Before delving directly into this analysis, however, we summarize some geometric and topological facts on the value function in the convex setting.

**Lemma 9** *Assume that  $f$  is convex. Then for all  $T \in \mathbb{N} \cup \{\infty\}$ ,  $J^T : \mathcal{X} \rightarrow \mathbb{R}$  is a convex and differentiable function. Moreover, we have that  $J^T \rightarrow J^\infty$  and  $\nabla J^T \rightarrow \nabla J^\infty$  uniformly on compact sets.*

**Lemma 10** *Let  $T \in \mathbb{N} \cup \{\infty\}$  and assume that  $f$  is convex. If we define  $\tilde{\phi}(x) := \phi(-x)$  then for all  $t < T$  the iterates of  $\mathbf{x} \in \mathcal{P}^T$  satisfy*

$$x_{t+1} = x_t - \nabla \tilde{\phi}^*(\nabla J^{T-t}(x_t)) \quad (9)$$

as well as the recursion  $\nabla J^{T-t}(x_t) = \nabla J^{T-t-1}(x_{t+1}) + \nabla f(x_{t+1})$ .

Hence, Lemma 10 shows a much clearer relationship to *dual-preconditioned gradient descent* (DPGD) of (Maddison et al., 2019, Algorithm 1.1). Indeed, the update rule of equation (9) corresponds to a single step of DPGD applied for descent on the objective  $J^{T-t}$  with preconditioner  $\tilde{\phi}$ . We note that the main difference lies in that we are performing descent on the value function rather than the objective  $f$ . When  $T = \infty$  it is easy to see that the descent is performed on  $J^\infty$  at each iteration  $t \in \mathbb{N}$ .

**Lemma 11** *Assume that  $f$  is strictly convex and let  $\tilde{\phi}(x) := \phi(-x)$ . Then for all  $T \in \mathbb{N} \cup \{\infty\}$ ,*

- i. *For  $T \neq 0$   $J^T$  is Legendre convex and  $\nabla(J^T)^* = (\nabla J^T)^{-1}$ .*
- ii. *Each  $J^T$  has the unique minimum  $x^*$  where  $\min_{x \in \mathcal{X}} J^T(x) = 0$ .*
- iii.  *$(J^T)^*$  is 1-relatively-convex with respect to  $\tilde{\phi}^*$ . If  $\phi$  is a symmetric positive-definite quadratic function then  $J^T$  is also 1-relatively-smooth with respect to  $\phi$ .*

Lemmas 9 and 11 demonstrate that the collection of value functions enjoy many regularity properties in terms of boundedness, differentiability and curvature. In particular Lemma 9 shows that  $J^T$  inherits both the convexity and differentiability of  $f$  and  $\phi$ . Moreover, the dynamics of Lemma 10 along with the observation of Lemma 11 that  $f$  and  $J^\infty$  share minima will be important for the analysis in further sections, where we study the descent of  $x \in \mathcal{P}^\infty$  on  $J^\infty$  and  $f$ .

Lemma 11-iii also has the important implication that the dual of  $J^T$  satisfies a relative convexity condition without any additional smoothness assumptions on  $f$ , which will prove crucial in the convergence analysis. In fact, in the case of a quadratic  $\phi$ , we show that these bounds are tightened if  $f$  is also relatively smooth or convex, which we demonstrate in the following lemma.

**Lemma 12** *Assume that  $f$  is strictly convex, and  $\phi$  is a symmetric positive-definite quadratic function. Define the function  $\Psi : \mathbb{R}_{\geq 0} \rightarrow [0, 1)$  as  $\Psi(x) := \frac{1}{2} \left( \sqrt{x^2 + 4x} - x \right)$ .*

- i. *If  $f$  is  $\lambda$ -relatively smooth w.r.t. to  $\phi$ , then  $J^\infty$  is  $\Psi(\lambda)$ -relatively smooth w.r.t.  $\phi$ .*
- ii. *If  $f$  is  $\mu$ -relatively convex w.r.t.  $\phi$ , then  $J^\infty$  is  $\Psi(\mu)$ -relatively convex w.r.t.  $\phi$ .*

#### 4.1. Convergence on Convex Objectives

Here, we provide bounds on the rate of convergence of regret optimal algorithms in the presence of a convex objective  $f$ . In particular, we focus our analysis on the behaviour of *asymptotically regret-optimal algorithms*  $x \in \mathcal{P}_x^\infty$ . The principal motivation behind the choice of  $T = \infty$  is the time-homogeneous nature of the algorithms implied by Lemma 10. We summarize the convergence rates derived over the course of this section in Table 1.

$f$	$\phi$	Potential	Rate	Reference
non-convex	super-linear	$f(x_t) - f_\star + \phi(\Delta x_t)$	$o(1), o(1/t)$	Corr. 3
strictly conv. / $\lambda$ -rel.-smooth	Legendre conv.	$\phi^*(\sum_{u=t}^\infty \nabla f(x_u))$	$O(1/t)$	Thm. 13
strictly conv. / $\lambda$ -rel.-smooth	p.s.d. quadratic	$f(x_t) - f_\star + \phi(\Delta x_t)$	$O(1/t^2)$	Thm. 15
$\mu$ -relatively-convex	p.s.d. quadratic	$f(x_t) - f_\star + \phi(\Delta x_t)$	$O(e^{-\epsilon t})$	Thm. 15

Table 1: Summary of the convergence rates on convex functions for asymptotically regret-optimal algorithms  $x \in \mathcal{P}_x^\infty$ . We provide a reference to the exact statement with precise constants in the right-most column.

In order to derive tighter rates of convergence for the class of asymptotically regret-optimal algorithms we leverage the connection to *dual-preconditioned gradient descent*, described in the discussion following Lemma 10. In what follows, we present a theorem establishing the rate of convergence of a regret-optimal algorithm with respect to the value function in the case of a convex loss.

**Theorem 13** *Assume that  $f$  is strictly convex,  $\tilde{\phi}$  is Legendre convex, and let  $x \in \mathcal{P}_x^\infty$  with the associated value function  $J^\infty$ . Then we have the bound*

$$\tilde{\phi}^*(\nabla J^\infty(x_t)) \leq \frac{J^\infty(x_0)}{t}. \quad (10)$$

**Theorem 14** *Assume that  $f$  is strictly convex,  $\phi$  is Legendre convex, and let  $\mathbf{x} \in \mathcal{P}_x^\infty$  with the associated value function  $J^\infty$ . Let  $\lambda$  be the  $\phi$ -relative-smoothness constant for  $J^\infty$ , then we have that*

$$J^\infty(x_t) \leq \frac{\lambda \phi(x_0 - x^*)}{t}. \quad (11)$$

*Suppose, in addition that  $J^\infty$  is  $\mu$ -relatively-convex with respect to  $\phi$ . Then we have that*

$$J^\infty(x_t) \leq \lambda \left(1 - \frac{2\mu}{1+\mu}\right)^t \phi(x_0 - x^*). \quad (12)$$

We note here that in the case where  $f$  is just convex, we have by Lemma 11-iii that  $\lambda = 1$ . In the case where  $f$  is either relatively convex or smooth with respect to  $\phi$ , we may obtain  $\lambda$  and  $\mu$  from Lemma 22 which further tightens the rate of convergence. Although the above theorems concern the rate of convergence on the value function,  $J^\infty$ , these also imply rates of convergence on the objective function  $f$  itself, as is shown in the following theorem.

**Theorem 15** *Suppose that the necessary conditions for equation (11) hold. Then*

$$f(x_t) - f_\star + \phi(\Delta x_t) \leq \frac{2\lambda \phi(x_0 - x^*)}{t^2}, \quad (13)$$

*for all but finitely many  $t$ . If the necessary conditions for equation (12) hold, then*

$$f(x_t) - f_\star + \phi(\Delta x_t) \leq \lambda \phi(x_0 - x^*) \left(1 - \frac{2\gamma}{1+\gamma}\right)^{t+1} \quad (14)$$

*for all but finitely many  $t$ .*

The bounds we provide over this section improve upon the general non-convex bound in Corollary 3. Moreover, we show that with additional assumptions on the relative smoothness and relative convexity of the objectives, we can further tighten these rates. In contrast to Corollary 3, we provide exact constants on the rate of convergence. We also point out that the bounds in Table 1 happen to coincide exactly with known lower bounds for the rate of convergence of gradient-based optimization algorithms as shown in Nesterov (2003). In particular, the  $O(t^{-2})$  rate of equation (13) implies that asymptotically-regret-optimal algorithms achieve rates of convergence on relatively-smooth objectives in the same class as the Nesterov accelerated gradient algorithm of Nesterov (1983) and its variants.

**Remark 16** *Although the analysis over the course of this section is applied for deriving rates of convergence for algorithms  $\mathbf{x} \in \mathcal{P}^\infty$ , very similar results can be derived for  $\mathbf{x} \in \mathcal{P}^T$  with  $T < \infty$  using the same techniques.*

## 5. Learning Regret-Optimal Algorithms

For an asymptotically regret-optimal algorithm  $\mathbf{x} \in \mathcal{P}_x^\infty$ , we turn to the problem of computing its dynamics so that it can be applied to tangible optimization problems. We approach the problem by taking the perspective of the optimal dynamics of Theorem 6. Although the optimal dynamics (5) exhibit closed form solutions in the case of descent on a quadratic objective, as is presented in

Appendix A, there are no such solutions available for general  $f$ . Hence, we turn to numerical methods for the approximation of the optimal dynamics presented in equation (5).

Let us assume a vector field  $\hat{v}^\theta : \mathcal{X} \rightarrow \mathcal{X}$  parametrized by  $\theta \in \Theta$ , for some vector space  $\Theta$ . Our approach will be to estimate the  $\theta$  such  $\hat{v}^\theta$  approximately satisfies the necessary optimality criteria of Theorem 7. Hence, we define the loss function

$$\mathcal{L}(\theta; x) = \left\| \hat{v}^\theta(x) - \hat{v}^\theta(\hat{y}^\theta(x)) - \nabla f(\hat{y}^\theta(x)) \right\|^2 \quad (15)$$

$$\text{where } \hat{y}^\theta(x) = x - \nabla \tilde{\phi}^*(\hat{v}^\theta(x)), \quad (16)$$

which is obtained by taking the squared norm of the difference between both sides of the second line of equation (8). Just as Q-learning serves as a method for approximating a function which satisfies the Bellman equation, one can interpret minimizing (15) as identifying a vector field which approximately satisfies the Pontryagin Maximum Principle. The loss function (15) also admits a more direct interpretation in connection with the Gâteaux derivative of the regret, which is presented in the following lemma.

**Lemma 17** *Let  $\mathbf{x}^\theta \in \mathcal{A}^\infty$  be the algorithm induced by the vector field  $\hat{v}^\theta$  according to equation (16) and let  $\mathcal{R}'_T$  be the Gâteaux derivative of  $\mathcal{R}_T$ . Then we have that*

$$\sum_{t=1}^T \mathcal{L}(\theta; x_{t-1}^\theta) = \left\| \mathcal{R}'_T(\mathbf{x}^\theta) \right\|_{2,*}^2, \quad (17)$$

where  $\|\cdot\|_{2,*}$  is the dual norm with respect to the norm on  $\mathcal{A}^\infty$  defined by  $\|\mathbf{x}\|_2^2 = \sum_{t \in \mathbb{N}} \|x_t\|^2$ .

Hence, we may loosely interpret minimizing  $\bar{\mathcal{L}}_\infty(\theta) = \sum_{t \in \mathbb{N}} \mathcal{L}(\theta; x_{t-1}^\theta)$  as minimizing the norm of the Gâteaux derivative of  $\mathcal{R}_\infty$ , and drawing  $\mathbf{x}^\theta$  closer to a critical point of  $\mathcal{R}_\infty$ .

In order to minimize  $\bar{\mathcal{L}}_\infty$ , we rely on online gradient descent or an equivalent algorithm to minimize the online loss function  $\mathcal{L}(\theta_t, x_{t-1})$  at each iteration. We note that each evaluation of  $\mathcal{L}(\theta_t, x_t)$  requires an evaluation of the gradient of  $f$ . In order to reduce the number of gradient evaluations, we consider the following approximation to  $\mathcal{L}$

$$\hat{\mathcal{L}}_t(\theta) = \left\| \hat{v}^\theta(x_t) - \hat{v}^\theta(\hat{y}^\theta(x_t)) - \nabla f(\hat{y}^{\theta_{t-1}}(x_t)) \right\|^2, \quad (18)$$

in which we “freeze”  $\theta$  in the expression  $\nabla f(y^\theta)$  within the above equation. We summarize these ideas in the Algorithm 1, which aims at minimizing  $\bar{\mathcal{L}}_\infty(\theta)$  in an online manner.

In Algorithm 1, we assume that the optimization routine within the inner-loop can be computed quickly and with a small memory footprint. To achieve this, one can ensure that  $\dim(\Theta) \ll \dim(\mathcal{X})$  so that gradients of  $\hat{v}^\theta$  can be computed cheaply. Moreover, at each iteration, the optimization routine does not need to fully optimize  $\mathcal{L}_t$  and can be replaced by a single iteration of gradient descent. As long as the inner optimization loop can be ensured to be fast, Algorithm 1 can serve as a viable option to adaptive optimization algorithms and can be applied to a wide range of optimization problems, regardless of the size of  $\dim(\mathcal{X})$ .

**Algorithm 1:** Online Regret Meta-Optimization

---

**input** :  $x_0 \in \mathcal{X}, T \in \mathbb{N}, \nabla\phi^*, \theta_0, \hat{v}^\theta$

**for**  $t \leftarrow 0$  **to**  $T - 1$  **do**

$\hat{y}_t \leftarrow x_t - \nabla\phi^*(\hat{v}^{\theta_t}(x_t))$  // Compute ``test`` point for gradient evaluation.

$\theta_{t+1} \leftarrow \text{Optimize}(\hat{\mathcal{L}}_t(\theta))$  // Update  $\theta$  estimate.

$x_{t+1} \leftarrow x_t - \nabla\phi^*(\hat{v}^{\theta_{t+1}}(x_t))$  // Step forward with new estimate.

**end**

**return**  $x_T$

---

**5.1. Auto-Tuning Gradient Descent with Momentum**

We present a simple numerical example for Algorithm 1, which we test on some basic optimization objectives. We choose to implement Algorithm 1 using a simple two-parameter model of the form

$$\hat{v}_t^\theta = \alpha \nabla f(y_{t-1}) + \beta \hat{v}_{t-1}^\theta,$$

where  $\theta = \{\alpha, \beta\}$  and we assume that  $\alpha, \beta > 0$  and where  $\nabla f(y_{t-1})$  is the last gradient that has been evaluated. We can interpret this model as a parametrized version of gradient descent with momentum, where  $\alpha, \beta$  control the weights on the gradient and momentum, respectively. We use a single step of gradient descent as the optimization routine for the algorithm. We note that with this formulation, the entire inner-loop optimization routine uses a single gradient computed once per outer-loop iteration so that the gradient complexity scales only with  $T$  and not with the number of inner-loop steps.

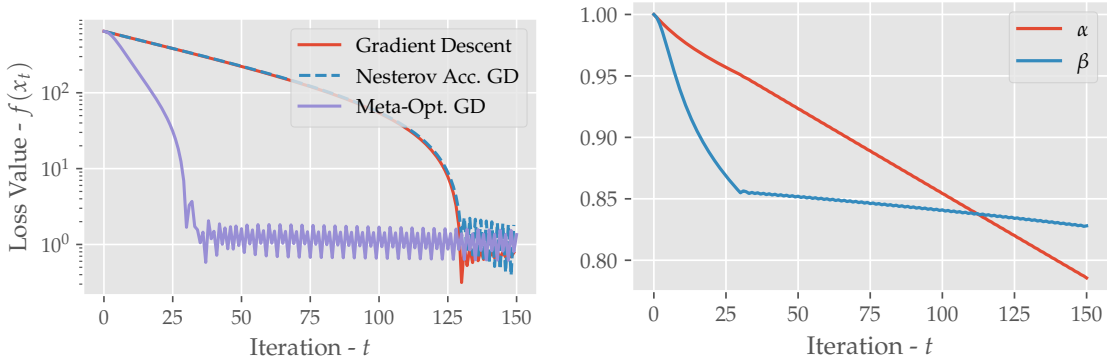


Figure 1: **Rescaled Rosenbrock Objective.** We include (left to right) a plot of the loss function value over each iteration and the evolution of the hyperparameters,  $\theta = \{\alpha, \beta\}$ , of Algorithm 1.

We compare the performance of the resulting online algorithm with gradient descent and Nesterov accelerated gradient descent, each with fixed hyperparameters. We use Algorithm 1 with  $\phi(x) = \frac{\gamma^{-1}}{2} \|x\|^2$ , where  $\gamma > 0$  is the learning rate. We set the learning rates to be the same value for all algorithms that are compared. In order to optimize in the inner-loop of Algorithm 1, we run

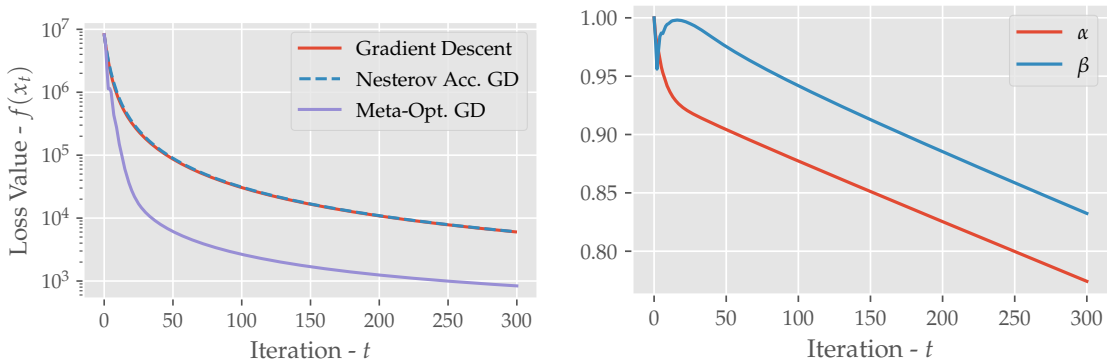


Figure 2:  **$2^{12}$ -Dimensional Quadratic Objective.** We include (left to right) a plot of the loss value over each iteration and the evolution of the algorithm hyperparameters,  $\theta$ .

10 steps of gradient descent with a learning rate of  $10^{-4}$ . We apply each algorithm on two types of examples, first on a rescaled Rosenbrock function<sup>4</sup> on  $\mathbb{R}^2$ , with results displayed in Figure 1 and second on a randomly generated symmetric positive-definite quadratic objective on  $\mathbb{R}^{2^{12}}$  with results displayed in Figure 2. These examples serve as a proof of concept and show that Algorithm 1 works comparatively well to two other well-known optimization algorithms on toy problems.

## 6. Discussion

Over the course of this paper, we characterize the existence and properties of regret optimal algorithms in a wide range of common optimization settings. One shortcoming of our approach, however, is that we do not restrict the measurability of algorithms in the minimization of regret. In light of this fact, it is interesting that we recover in Table 1 bounds that look quite similar to optimal convergence bounds for gradient-based optimization algorithms, in particular the  $O(t^{-2})$  bound that is known to hold for “accelerated” algorithms. A more in depth analysis of these rates and comparison to known lower bounds would also be very interesting.

This paper presents new perspectives on optimization which deserve to be further explored. An interesting potential avenue of research would be the extension of this framework towards stochastic optimization. Another direction would be to see how commonly used optimization algorithms fall within this framework, and to determine whether they satisfy regret optimality in an exact or approximate sense.

## References

Jean-Pierre Aubin and H el ene Frankowska. *Set-valued analysis*. Modern Birkh user Classics. Birkh user Boston, Inc., Boston, MA, 2009. ISBN 978-0-8176-4847-3. doi: 10.1007/978-0-8176-4848-0. URL <https://doi.org/10.1007/978-0-8176-4848-0>. Reprint of the 1990 edition [MR1048347].

4. Let  $f_r(x, y)$  be the Rosenbrock function on  $\mathbb{R}^2$  (Rosenbrock, 1960). We define the rescaled Rosenbrock function as  $(x, y) \mapsto 0.1 \sqrt{f_r(0.5x, 4.5y)}$  which has the property that it is relatively smooth w.r.t.  $\phi$ .

- Michael Betancourt, Michael I Jordan, and Ashia C Wilson. On symplectic optimization. *arXiv preprint arXiv:1802.03653*, 2018.
- Joël Blot and Naïla Hayek. *Infinite-horizon optimal control in the discrete-time framework*. Springer, 2014.
- Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- Andrea Braides. *Local minimization, variational evolution and  $\Gamma$ -convergence*, volume 2094 of *Lecture Notes in Mathematics*. Springer, Cham, 2014. ISBN 978-3-319-01981-9; 978-3-319-01982-6. doi: 10.1007/978-3-319-01982-6. URL <https://doi.org/10.1007/978-3-319-01982-6>.
- Philippe Casgrain. A latent variational framework for stochastic optimization. In *Advances in Neural Information Processing Systems*, pages 5646–5656, 2019.
- John B. Conway. *A course in functional analysis*, volume 96 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, second edition, 1990. ISBN 0-387-97245-5.
- Gianni Dal Maso. *An introduction to  $\Gamma$ -convergence*, volume 8 of *Progress in Nonlinear Differential Equations and their Applications*. Birkhäuser Boston, Inc., Boston, MA, 1993. ISBN 0-8176-3679-X. doi: 10.1007/978-1-4612-0327-8. URL <https://doi.org/10.1007/978-1-4612-0327-8>.
- Ennio De Giorgi, Giuseppe Congedo, and Italo Tamanini. Regularity problems for a new functional in the calculus of variations. *Atti Accad. Naz. Lincei Rend. Cl. Sci. Fis. Mat. Nat. (8)*, 82(4): 673–678 (1990), 1988. ISSN 0392-7881.
- Joseph Diestel. *Sequences and series in Banach spaces*, volume 92 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, 1984. ISBN 0-387-90859-5. doi: 10.1007/978-1-4612-5200-9. URL <https://doi.org/10.1007/978-1-4612-5200-9>.
- Yoel Drori and Adrien Taylor. On the oracle complexity of smooth strongly convex minimization, 2021.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011a.
- John C Duchi, Alekh Agarwal, and Martin J Wainwright. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic control*, 57(3): 592–606, 2011b.
- Ivar Ekeland and Roger Temam. *Convex analysis and variational problems*, volume 28. Siam, 1999.
- Juan Ferrera. *An introduction to nonsmooth analysis*. Academic Press, 2013.
- Guilherme França, Michael I Jordan, and René Vidal. On dissipative symplectic integration with applications to gradient-based optimization. *arXiv preprint arXiv:2004.06840*, 2020.

- Elad Hazan. Introduction to online convex optimization. *Foundations and Trends in Optimization*, 2(3-4):157–325, 2016.
- Bin Hu, Peter Seiler, and Anders Rantzer. A unified analysis of stochastic optimization methods using jump system theory and quadratic constraints. *Proceedings of Machine Learning Research* vol, 65:1–33, 2017.
- S Kakade, Shai Shalev-Shwartz, and Ambuj Tewari. Applications of strong convexity–strong smoothness duality to learning with matrices. *CoRR*, abs/0910.0610, 2009.
- K. Kuratowski. *Topology. Vol. I*. New edition, revised and augmented. Translated from the French by J. Jaworowski. Academic Press, New York-London; Państwowe Wydawnictwo Naukowe, Warsaw, 1966.
- K. Kuratowski. *Topology. Vol. II*. New edition, revised and augmented. Translated from the French by A. Kirkor. Academic Press, New York-London; Państwowe Wydawnictwo Naukowe Polish Scientific Publishers, Warsaw, 1968.
- Laurent Lessard, Benjamin Recht, and Andrew Packard. Analysis and design of optimization algorithms via integral quadratic constraints. *SIAM Journal on Optimization*, 26(1):57–95, 2016.
- Qianxiao Li, Long Chen, Cheng Tai, and Weinan E. Maximum principle based algorithms for deep learning. *Journal of Machine Learning Research*, 18(165):1–29, 2018. URL <http://jmlr.org/papers/v18/17-653.html>.
- Jonathan Lorraine and David Duvenaud. Stochastic hyperparameter optimization through hypernetworks. *arXiv preprint arXiv:1802.09419*, 2018.
- Chris J Maddison, Daniel Paulin, Yee Whye Teh, and Arnaud Doucet. Dual space preconditioning for gradient descent. *arXiv preprint arXiv:1902.02257*, 2019.
- Alexander Mitsos, Jaromił Najman, and Ioannis G Kevrekidis. Optimal deterministic algorithm generation. *Journal of Global Optimization*, 71(4):891–913, 2018.
- Michael Muehlebach and Michael I Jordan. Optimization with momentum: Dynamical, control-theoretic, and symplectic perspectives. *arXiv preprint arXiv:2002.12493*, 2020.
- James R. Munkres. *Topology*, 2000. Second edition of [MR0464128].
- Yurii Nesterov. A method for unconstrained convex minimization problem with the rate of convergence  $O(1/k^2)$ . In *Doklady AN USSR*, volume 269, pages 543–547, 1983.
- Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2003.
- R. Tyrell Rockafellar. *Convex Analysis*. Princeton University Press, 1970. ISBN 9780691015866. URL <http://www.jstor.org/stable/j.ctt14bs1ff>.
- HoHo Rosenbrock. An automatic method for finding the greatest or least value of a function. *The Computer Journal*, 3(3):175–184, 1960.



- Bin Shi, Simon S Du, Weijie Su, and Michael I Jordan. Acceleration via symplectic discretization of high-resolution differential equations. In *Advances in Neural Information Processing Systems*, pages 5744–5752, 2019.
- Adrien Taylor and Yoel Drori. An optimal gradient method for smooth strongly convex minimization, 2021.
- Adrien Taylor, Julien Hendrickx, and François Glineur. Smooth strongly convex interpolation and exact worst-case performance of first-order methods. *Mathematical Programming*, 161(1):307–345, 01 2017. ISSN 1436-4646.
- L. Tonelli. *Opere scelte*. Vol. II: Calcolo delle variazioni. Edizioni Cremonese, 1961.
- Andre Wibisono, Ashia C Wilson, and Michael I Jordan. A variational perspective on accelerated methods in optimization. *Proceedings of the National Academy of Sciences*, 113(47):E7351–E7358, 2016.
- Olga Wichrowska, Niru Maheswaranathan, Matthew W Hoffman, Sergio Gómez Colmenarejo, Misha Denil, Nando Freitas, and Jascha Sohl-Dickstein. Learned optimizers that scale and generalize. In *International Conference on Machine Learning*, pages 3751–3760, 2017.
- Ashia C Wilson, Lester Mackey, and Andre Wibisono. Accelerating rescaled gradient descent: Fast optimization of smooth functions. In *Advances in Neural Information Processing Systems*, pages 13555–13565, 2019.
- Lin Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *Journal of Machine Learning Research*, 11:2543–2596, 2010.

## Appendix A. Regret-Optimal Dynamics on Quadratic Objectives

**Lemma 18 (Descent on Quadratic Functions)** *Assume that there exist  $A, C \in S_{++}^d$  and  $b \in \mathbb{R}^d$  such that*

$$\nabla f(x) = A(x - b) \quad \text{and} \quad \phi(z) = \frac{1}{2}z^\top C z .$$

*If there exists a matrix  $\tilde{\Phi} \in S_{++}^d$  such that*

$$\tilde{\Phi}^{-1} = C^{-1} + (A + \tilde{\Phi})^{-1} ,$$

*then we have that*

$$\Delta x_t = -C^{-1}\tilde{\Phi}_1(x_t - b) ,$$

*for all  $t \in \mathbb{N}$ .*

**Proof** First note that under the assumptions of the theorem, we have that there exists a constant  $c$  such that  $f(x) = c + \frac{1}{2}(x - b)^\top A(x - b)$ , demonstrating that  $f$  convex function. Moreover, if  $x \in \mathcal{P}^\infty$ , we also have by Lemma 21 that

$$(J^\infty)^*(p) = \phi^*(p) + (f + J^\infty)^*(p) , \tag{19}$$

where, by a simple computation, we have that  $\phi^*(p) = \frac{1}{2}p^\top C^{-1}p$ . Moreover, we note that without loss of generality, we may consider the case where  $b = 0$  since we may simply use the domain transformation  $x \mapsto x + b$ .

Now, assume that there exists  $\tilde{\Phi} \in S_{++}^d$  such that  $\nabla J^\infty(x) = \tilde{\Phi}x$ . Differentiating equation (19) and noting that the linearity of the gradients of  $J^\infty$  and  $f$  imply that  $\nabla(f + J^\infty)^*(p) = (A + \tilde{\Phi})^{-1}p$ , we get that

$$\tilde{\Phi}^{-1}p = C^{-1}p + (A + \tilde{\Phi})^{-1}p ,$$

which must hold for all  $p \in \mathcal{X}$ . Hence, if there exists  $\tilde{\Phi} \in S_{++}^d$  such that

$$\tilde{\Phi}^{-1} = C^{-1} + (A + \tilde{\Phi})^{-1} ,$$

noting that  $\Delta x_t = -C^{-1}\nabla J^\infty(x)$ , we find that the statement of the theorem must hold. ■

## Appendix B. Auxiliary Results

### B.1. Dynamic Programming Principles

**Lemma 19 (Dynamic Programming Principle,  $T < \infty$ )** *Suppose that Assumptions 1 and 2 hold. If for any  $T \in \mathbb{N}$  we define the value function  $J^T(x) = \min_{\mathbf{x} \in \mathcal{A}_x^T} \mathcal{R}_T(\mathbf{x})$ , then for each  $t, T \in \mathbb{N}$ , the iterates of  $\mathbf{x}$  satisfy*

$$J^T(x_t) = \min_{y \in \mathcal{X}} \{ \phi(y - x_t) + f(y) - f_* + J^{T-1}(y) \},$$

where  $x_{t+1} \in \operatorname{argmin}_{y \in \mathcal{X}} \{ \phi(y - x_t) + f(y) - f_* + J^{T-1}(y) \}$ . Moreover, for each  $0 < h < T$ , defining the shifted sequence  $\mathbf{x}_{\rightarrow h}$  such that  $\mathbf{x}_{\rightarrow h} = \{y_{t+h}\}_{t=0}^\infty$ , we have that  $\mathbf{x}_{\rightarrow h} \in \mathcal{P}_{x_h}^{T-h}$ .

**Proof** We assume without loss of generality that  $f_* = 0$ . We first note that by Theorem 1, we have that  $\mathcal{P}_x^T \neq \emptyset$  for all  $T \in \mathbb{N}$  and  $x \in \mathcal{X}$ . Moreover, by the definition of  $\mathcal{P}_x^T$ , we also have that  $J^T(x) = \mathcal{R}_T(\mathbf{x}) < \infty$  for any  $\mathbf{x} \in \mathcal{P}_x^T$ .

Recursively expanding  $\mathcal{R}_T$ , we find that for any  $\mathbf{y} \in \mathcal{A}^T$  such that  $\mathcal{R}_T(\mathbf{y}) < \infty$  and  $0 \leq h \leq T$ , we have that

$$\mathcal{R}_T(\mathbf{y}) = \mathcal{R}_h(\mathbf{y}) + \mathcal{R}_{T-h}(\mathbf{y}_{\rightarrow h}),$$

where  $\mathbf{y}_{\rightarrow h}$  is the shifted sequence  $\mathbf{y}_{\rightarrow h} = \{y_{t+h}\}_{t=0}^h$ . The above recursion also includes the special case that

$$\mathcal{R}_T(\mathbf{y}) = \phi(y_1 - y_0) + f(y_1) + \mathcal{R}_{T-1}(\mathbf{y}_{\rightarrow 1}). \quad (20)$$

Now we show that  $\mathbf{x} \in \mathcal{P}_x^T \Rightarrow \mathbf{x}_{\rightarrow 1} \in \mathcal{P}_{x_1}^{T-1}$ . Assume the converse that  $\mathbf{x}_{\rightarrow 1} \notin \mathcal{P}_{x_1}^{T-1}$  (i.e. that  $\mathcal{R}_{T-1}(\mathbf{x}_{\rightarrow 1}) > J^{T-1}(x_1)$ ). Applying (20) we have that

$$\begin{aligned} J^T(x) &= \mathcal{R}_T(\mathbf{x}) \\ &= \phi(x_1 - x_0) + f(x_1) + \mathcal{R}_{T-1}(\mathbf{x}_{\rightarrow 1}) \\ &> \phi(x_1 - x_0) + f(x_1) + J^{T-1}(x_1). \end{aligned}$$

But, defining a control  $\mathbf{y} \in \mathcal{A}_x^T$  such that  $y_0 = x_1$  and  $\mathbf{y}_{\rightarrow 1} \in \mathcal{P}_{x_1}^{T-1}$ , we have that by the definition of  $\mathbf{x} \in \mathcal{P}_x^{T-1}$ ,

$$\begin{aligned} J^T(x) &\leq \mathcal{R}_T(\mathbf{y}) \\ &= \phi(x_1 - x_0) + f(x_1) + \mathcal{R}_{T-1}(\mathbf{y}_{\rightarrow 1}) \\ &= \phi(x_1 - x_0) + f(x_1) + J^{T-1}(x_1), \end{aligned}$$

which is a contradiction. Hence, we find that  $\mathbf{x}_{\rightarrow 1} \in \mathcal{P}_{x_1}^{T-1}$  and

$$J^T(x_0) = \phi(x_1 - x_0) + f(x_1) + J^{T-1}(x_1), \quad (21)$$

for all  $\mathbf{x} \in \mathcal{P}_x^T$  and  $x \in \mathcal{X}$ . By induction on  $h$ , we also note that  $\mathbf{x}_{\rightarrow 1} \in \mathcal{P}_{x_1}^{T-1}$  also implies that  $\mathbf{x}_{\rightarrow h} \in \mathcal{P}_{x_h}^{T-h}$  for  $0 < h < T$ .

Now suppose that for  $\mathbf{x} \in \mathcal{P}_x^T$ , we have that  $x_1 \notin B_x = \operatorname{argmin}_{y \in \mathcal{X}} \{ \phi(y - x_0) + f(y) + J^{T-1}(y) \}$ . Then by (21), there exists  $y \in \mathcal{X}$  such that

$$\phi(y - x_0) + f(y) + J^{T-1}(y) < \phi(x_1 - x_0) + f(x_1) + J^{T-1}(x_1) = J^T(x_0).$$

Conversely, defining  $\mathbf{y} \in \mathcal{A}_x^T$  such that  $\mathbf{y}_{\rightarrow 1} \in \mathcal{P}_y^{T-1}$ , we have that by the definition of  $J^\infty$  that

$$\begin{aligned} J^T(x) &\leq \mathcal{R}(\mathbf{y}) \\ &= \phi(y - x) + f(y) + \mathcal{R}_{T-1}(\mathbf{y}_{\rightarrow 1}) \\ &= \phi(y - x) + f(y) + J^{T-1}(y), \end{aligned}$$

which is again a contradiction, and hence  $x_1 \in \operatorname{argmin}_{y \in \mathcal{X}} \{\phi(y - x_0) + f(y) + J^{T-1}(y)\}$ . Combining this result with (21), and noting that  $t$  is arbitrary we therefore conclude the proof.  $\blacksquare$

**Lemma 20 (Dynamic Programming Principle,  $T = \infty$ )** *Suppose that  $\mathcal{P}_x^\infty$  is non-empty  $\forall x \in \mathcal{X}$  and that Assumptions 1 and 2 hold. If we define the value function  $J^\infty(x) = \min_{\mathbf{x} \in \mathcal{A}_x^\infty} \mathcal{R}_\infty(\mathbf{x})$ , then for each  $t$ , the iterates of  $\mathbf{x}$  satisfy*

$$J^\infty(x_t) = \min_{y \in \mathcal{X}} \{\phi(y - x_t) + f(y) - f_\star + J^\infty(y)\},$$

where  $x_{t+1} \in \operatorname{argmin}_{y \in \mathcal{X}} \{\phi(y - x_t) + f(y) - f_\star + J^\infty(y)\}$ . Moreover, for each  $T \in \mathbb{N}$ , defining the shifted sequence  $\mathbf{x}_{\rightarrow h}$  such that  $\mathbf{x}_{\rightarrow h} = \{y_{t+h}\}_{t=0}^\infty$ , we have that  $\mathbf{x}_{\rightarrow h} \in \mathcal{P}_{x_h}^\infty$ .

**Proof** We assume without loss of generality that  $f_\star = 0$ . We first note that by Theorem 2, we have that  $0 \leq J^\infty(x) < \infty$  for all  $x \in \mathcal{X}$ . Moreover, by the definition of  $\mathcal{P}_x^\infty$ , we also have that  $J^\infty(x) = \mathcal{R}_\infty(\mathbf{x})$  for any  $\mathbf{x} \in \mathcal{P}_x^\infty$ . The remainder of the proof resembles closely the proof of Lemma 19.

Recursively expanding  $\mathcal{R}_\infty$ , we find that for any  $\mathbf{y} \in \mathcal{A}^\infty$  such that  $\mathcal{R}_\infty(\mathbf{y}) < \infty$  and  $h > 0$ , we have that

$$\mathcal{R}_\infty(\mathbf{y}) = \mathcal{R}_h(\mathbf{y}) + \mathcal{R}_\infty(\mathbf{y}_{\rightarrow h}),$$

where  $\mathbf{y}_{\rightarrow h}$  is the shifted sequence  $\mathbf{y}_{\rightarrow h} = \{y_{t+h}\}_{t=0}^h$ . The above recursion also includes the special case that

$$\mathcal{R}_\infty(\mathbf{y}) = \phi(y_1 - y_0) + f(y_1) + \mathcal{R}_\infty(\mathbf{y}_{\rightarrow 1}). \quad (22)$$

Now we show that  $\mathbf{x} \in \mathcal{P}_x^\infty \Rightarrow \mathbf{x}_{\rightarrow 1} \in \mathcal{P}_{x_1}^\infty$ . Assume the converse that  $\mathbf{x}_{\rightarrow 1} \notin \mathcal{P}_{x_1}^\infty$  (i.e. that  $\mathcal{R}_\infty(\mathbf{x}_{\rightarrow 1}) > J^\infty(x_1)$ ). Applying (22) we have that

$$\begin{aligned} J^\infty(x) &= \mathcal{R}_\infty(\mathbf{x}) \\ &= \phi(x_1 - x_0) + f(x_1) + \mathcal{R}_\infty(\mathbf{x}_{\rightarrow 1}) \\ &> \phi(x_1 - x_0) + f(x_1) + J^\infty(x_1). \end{aligned}$$

But, defining a control  $\mathbf{y} \in \mathcal{A}_x^\infty$  such that  $y_0 = x_1$  and  $\mathbf{y}_{\rightarrow 1} \in \mathcal{P}_{x_1}^\infty$ , we have that by the definition of  $\mathbf{x} \in \mathcal{P}_x^\infty$ ,

$$\begin{aligned} J^\infty(x) &\leq \mathcal{R}_\infty(\mathbf{y}) \\ &= \phi(x_1 - x_0) + f(x_1) + \mathcal{R}_\infty(\mathbf{y}_{\rightarrow 1}) \\ &= \phi(x_1 - x_0) + f(x_1) + J^\infty(x_1), \end{aligned}$$

which is a contradiction. Hence, we find that  $\mathbf{x}_{\rightarrow 1} \in \mathcal{P}_{x_1}^\infty$  and

$$J^\infty(x_0) = \phi(x_1 - x_0) + f(x_1) + J^\infty(x_1), \quad (23)$$

for all  $\mathbf{x} \in \mathcal{P}_x^\infty$  and  $x \in \mathcal{X}$ . By induction on  $h$ , we also note that  $\mathbf{x}_{\rightarrow 1} \in \mathcal{P}_{x_1}^\infty$  also implies that  $\mathbf{x}_{\rightarrow h} \in \mathcal{P}_{x_h}^\infty$  for any  $h \in \mathbb{N}$ .

Now suppose that for  $\mathbf{x} \in \mathcal{P}_x^\infty$ , we have that  $x_1 \notin B_x = \operatorname{argmin}_{y \in \mathcal{X}} \{\phi(y - x_0) + f(y) + J^\infty(y)\}$ . Then by (23), there exists  $y \in \mathcal{X}$  such that

$$\phi(y - x_0) + f(y) + J^\infty(y) < \phi(x_1 - x_0) + f(x_1) + J^\infty(x_1) = J^\infty(x_0).$$

Conversely, defining  $\mathbf{y} \in \mathcal{A}_x^\infty$  such that  $\mathbf{y}_{\rightarrow 1} \in \mathcal{P}_y^\infty$ , we have that by the definition of  $J^\infty$  that

$$\begin{aligned} J^\infty(x) &\leq \mathcal{R}(\mathbf{y}) \\ &= \phi(y - x) + f(y) + \mathcal{R}_\infty(\mathbf{y}_{\rightarrow 1}) \\ &= \phi(y - x) + f(y) + J^\infty(y), \end{aligned}$$

which is again a contradiction, and hence  $x_1 \in \operatorname{argmin}_{y \in \mathcal{X}} \{\phi(y - x_0) + f(y) + J^\infty(y)\}$ . Combining this result with (23), and noting that  $t$  is arbitrary we therefore conclude the proof.  $\blacksquare$

## B.2. Convex Analysis Results

This section compiles some auxiliary results from convex analysis which are needed in the proofs of the main theorems below. In all proofs within this subsection, we assume without loss of generality that  $f_\star = 0$ , since we may simply consider the objective  $\tilde{f}(x) = f(x) - f_\star$  which satisfies this property.

**Lemma 21** *Let Assumptions 1, 2 and 3 hold and assume that  $f$  is strictly convex. Assume that for any  $x, y \in \mathcal{X}$  and  $T \in \mathbb{N} \cup \{\infty\}$ , define the dual points  $q = \nabla J^T(x)$  and  $p = \nabla J^T(y)$  as well as the function  $\tilde{\phi}(x) := \phi(-x)$ . We therefore have that*

$$D_{(J^T)^*}(q, p) = D_{\tilde{\phi}^*}(q, p) + D_{(J^{T-1}+f)^*}(q, p), \quad (24)$$

as well as the recursion

$$(J^T)^*(q) = \tilde{\phi}^*(q) + (J^{T-1} + f)^*(q). \quad (25)$$

**Proof** Recall that Lemma 9 implies that  $J^T$  is convex and differentiable for all  $T \in \mathbb{N}$  under the assumptions of the theorem. Recalling the recursive properties of  $\nabla J^T(x)$  and  $J^T(x)$  (Lemma 10), for  $x, y \in \mathcal{X}$  and  $\mathbf{x} \in \mathcal{P}_x^T$  and  $\mathbf{y} \in \mathcal{P}_y^T$ , we compute

$$\begin{aligned} D_{J^T}(x, y) &\stackrel{(a)}{=} J^T(x_0) - J^T(y_0) - \langle \nabla J^T(y), x_0 - y_0 \rangle \\ &\stackrel{(b)}{=} \{J^{T-1}(x_1) + f(x_1) - J^{T-1}(y_1) - f(y_1) - \langle \nabla J^{T-1}(y_1) + \nabla f(y_1), x_1 - y_1 \rangle\} \\ &\quad + \{\phi(\Delta x_0) - \phi(\Delta y_0) - \langle \nabla J^T(y), (x_1 - x_0) - (y_1 - y_0) \rangle\} \\ &\stackrel{(c)}{=} D_{J^{T-1}+f}(x_1, y_1) + \{\phi(\Delta x_0) - \phi(\Delta y_0) + \langle \nabla J^T(y_0), (x_1 - x_0) - (y_1 - y_0) \rangle\} \\ &\stackrel{(d)}{=} D_{J^{T-1}+f}(x_1, y_1) + \{\phi(\Delta x_0) - \phi(\Delta y_0) - \langle \nabla \phi(\Delta y_0), (x_1 - x_0) - (y_1 - y_0) \rangle\} \\ &\stackrel{(e)}{=} D_{J^{T-1}+f}(x_1, y_1) + D_\phi(\Delta x_0, \Delta y_0), \end{aligned} \quad (26)$$

where (a) follows from the definition of the Bregman Divergence, where (b) follows from the recursive expansion  $J^T(x_0) = f(x_1) + \phi(\Delta x_0) + J^{T-1}(x_1)$  and  $\nabla J^T(x_0) = \nabla(J^{T-1} + f)(x_1)$ , where (c) follows from the definition of the Bregman divergence applied to the left-most curly braces, where (d) follows from the identity  $\nabla J^T(x_0) = -\nabla\phi(\Delta x_0) = \nabla\tilde{\phi}(-\Delta x_0)$  obtained from Lemma 10, and where (e) again follows from the definition of the Bregman divergence.

Now, recall the property of the Bregman divergence that

$$D_g(x, y) = D_{g^*}(\nabla g(y), \nabla g(x)) \quad (27)$$

for convex and differentiable  $g$  and  $g^*$ . Since  $f$  is strictly convex and differentiable and  $J^{T-1}$  is convex and differentiable,  $J^{T-1} + f$  is strictly convex and differentiable. Hence by (Rockafellar, 1970, Theorem 26.3)  $(J^{T-1} + f)^*$  is also differentiable and strictly convex, so we have we have

$$\begin{aligned} D_{J^{T-1}+f}(x_1, y_1) &= D_{(J^{T-1}+f)^*}(\nabla(J^{T-1} + f)(x_1), \nabla(J^{T-1} + f)(y_1)) \\ &= D_{(J^{T-1}+f)^*}(\nabla J^T(x_0), \nabla J^T(y_0)) , \end{aligned}$$

where in the second line, we use the recursive property that  $\nabla J^T(x_0) = \nabla J^{T-1}(x_1) + \nabla f(x_1)$  from Lemma 10. Applying (27) and that  $\nabla J^T(x_0) = \nabla\tilde{\phi}(-\Delta x_0)$ , we get

$$D_\phi(\Delta x_0, \Delta y_0) = D_{\tilde{\phi}^*}(\nabla J^T(y_0), \nabla J^T(x_0)) .$$

Combining these results with equations (26) and (27), we obtain

$$D_{(J^T)^*}(\nabla J^T(y), \nabla J^T(x)) = D_{\tilde{\phi}^*}(\nabla J^T(y_0), \nabla J^T(x_0)) + D_{(J^{T-1}+f)^*}(\nabla J^T(x_0), \nabla J^T(y_0)) ,$$

and hence, letting  $p = \nabla J^T(x)$  and  $q = \nabla J^T(y)$ , we have

$$D_{(J^T)^*}(q, p) = D_{\tilde{\phi}^*}(q, p) + D_{(J^{T-1}+f)^*}(q, p) .$$

Letting  $p = 0 = \nabla J^T(x^*)$  in the above, we obtain the second result.  $\blacksquare$

**Lemma 22 (Bregman Relative Duality)** *Consider a convex, differentiable function  $g : \mathcal{X} \rightarrow \mathbb{R}$ . If  $\phi$  is a symmetric positive-definite quadratic function on  $\mathcal{X}$ , then*

- i. If  $g$  is  $\lambda$ -relatively-smooth with respect to  $\phi$ , then  $g^*$  is  $\frac{1}{\lambda}$ -relatively-convex with respect to  $\phi^*$ .*
- ii. If  $g$  is  $\mu$ -relatively-convex with respect to  $\phi$ , then  $g^*$  is  $\frac{1}{\mu}$ -relatively-smooth with respect to  $\phi^*$ .*

**Proof** We first note that *i* and *ii* are identical statements, where we obtain the other by interchanging  $g$  and  $g^*$ . Hence, we only show the proof of *i*.

We begin by establishing a few results regarding  $\phi$ . First, since  $\phi$  is symmetric positive-definite and quadratic, there exists a norm  $\|\cdot\|_0$  on  $x$  such that  $\phi(x) = \frac{1}{2}\|x\|_0^2$ . This fact is easy to verify by setting  $\|\cdot\|_0 = \sqrt{2\phi(x)}$ , and verifying that this satisfies the necessary conditions of a norm. Next, it is also easy to verify (either by simple computation or see (Boyd et al., 2004, Example 3.27)) that  $\phi^*(p) = \frac{1}{2}\|p\|_{0,*}^2$ , where  $\|\cdot\|_{0,*}$  represents the dual norm of  $\|\cdot\|_0$ .

Since  $\phi$  is quadratic, we also have that  $D_\phi(x, y) = \frac{1}{2}\|x - y\|_0^2$ , and therefore the definition of relative smoothness and strong smoothness with respect to  $\|\cdot\|_0$  of (Kakade et al., 2009, Definition 5) coincide. Hence,  $g$  is  $\lambda$ -strongly-smooth with respect to the norm  $\|\cdot\|_0$ . Applying (Kakade et al., 2009, Theorem 6), we obtain that  $g^*$  must be  $\frac{1}{\lambda}$ -strongly-convex with respect to the norm  $\|\cdot\|_{0,*}$ , and hence  $g^*$  is  $\frac{1}{\lambda}$ -relatively-convex with respect to  $\phi^*$ , yielding the desired result.  $\blacksquare$

### B.3. Descent Lemmas

**Theorem 23 (Primal Descent Lemma)** *Let Assumptions 1, 2 and 3 hold, and let  $f$  be strictly convex. For  $T \in \mathbb{N} \cup \{\infty\}$ , let  $\mathbf{x} \in \mathcal{P}_x^T$  and assume that  $\phi$  is a symmetric positive-definite quadratic function on  $\mathcal{X}$ , and that  $J^T$  is 1-relatively-smooth with respect to  $\phi$ . For  $0 \leq t \leq T - 1$ , we have*

1. For all  $y \in \mathcal{D}$ ,  $J^T(x_{t+1}) - J^T(y) \leq -D_{J^T}(y, x_t) + D_\phi(y, x_t) - D_\phi(y, x_{t+1})$ ,

from which it follows that

2.  $J^T(x_{t+1}) - J^T(x_t) \leq -D_\phi(x_t, x_{t+1})$ , and
3.  $D_{J^T}(x_{t+1}, x^*) \leq -D_{J^T}(x^*, x_t) + D_\phi(x^*, x_t) - D_\phi(x^*, x_{t+1})$ .

**Proof** We first note that by Lemma 11, the assumptions of the theorem guarantee that  $J^T$  is Legendre convex. We first note that since  $\phi$  is quadratic, we have the properties that  $\phi(x) = \phi(-x)$ . Due to the assumed linearity of  $\nabla\phi$ , we find that the update rule for  $\mathbf{x} \in \mathcal{P}_{t,x}^T$  can be expressed as

$$\begin{aligned} \nabla\phi(x_{t+1} - x_t) &= \nabla\phi(x_{t+1}) - \nabla\phi(x_t) \\ &= -\nabla J^T(x_t), \end{aligned}$$

where the second equality holds due to the linearity of  $\nabla\phi$ . It is then easy to verify that this update rule can also be expressed as the proximal update rule

$$x_{t+1} = \arg \min_z \{ \langle \nabla J^T(x_t), z - x_t \rangle + D_\phi(z, x_t) \} \quad (28)$$

which for any  $z \in \mathcal{D}$  satisfies the proximal inequality

$$\langle \nabla J^T(x_t), z - x_t \rangle + D_\phi(z, x_t) \geq \langle \nabla J^T(x_t), x_{t+1} - x_t \rangle + D_\phi(x_{t+1}, x_t) + D_\phi(z, x_{t+1}). \quad (29)$$

Hence, we have that

$$\begin{aligned} J^T(x_{t+1}) &\stackrel{(a)}{\leq} J^T(x_t) + \langle \nabla J^T(x_t), x_{t+1} - x_t \rangle + D_\phi(x_{t+1}, x_t) \\ &\stackrel{(b)}{\leq} J^T(x_t) + \langle \nabla J^T(x_t), x - x_t \rangle + D_\phi(x, x_t) - D_\phi(x, x_{t+1}) \\ &\stackrel{(c)}{\leq} J^T(x) - D_{J^T}(x, x_t) + D_\phi(x, x_t) - D_\phi(x, x_{t+1}), \end{aligned}$$

where (a) follows from 9-i, (b) follows from the proximal inequality and (c) follows from simple algebra, yielding the first result in the statement of the lemma. The second and third follow by applying the special cases  $y = x_t$  and  $y = x^*$  to the first.  $\blacksquare$

**Theorem 24 (Dual Descent Lemma)** *Let Assumptions 1, 2 and 3 hold, and let  $f$  be strictly convex. For  $T \in \mathbb{N} \cup \{\infty\}$ , let  $\mathbf{x} \in \mathcal{P}^T$  and assume that  $(J^T)^*$  is  $\lambda$ -relatively-convex with respect to  $\phi$ .*

1. For all  $y \in \mathcal{D}$  and  $0 \leq t \leq T - 1$ , and letting  $\tilde{\phi}(x) := \phi(-x)$ , we have

$$\begin{aligned} \tilde{\phi}^*(\nabla J^T(x_{t+1})) - \tilde{\phi}^*(\nabla J^T(y)) &\leq -D_{\tilde{\phi}^*}(\nabla J^T(y), \nabla J^T(x_t)) \\ &\quad + D_{J^T}(x_t, y) - D_{J^T}(x_{t+1}, y). \end{aligned}$$

Moreover, the above inequality implies that for all  $0 \leq t \leq T - 1$ ,

2.  $\tilde{\phi}^* (\nabla J^T(x_{t+1})) + D_{J^T} (x_{t+1}, x_t) \leq \tilde{\phi}^* (\nabla J^T(x_t))$ , and
3. For any  $x^* \in \arg \min_{y \in \mathcal{D}} J^T(y)$ ,

$$\begin{aligned} D_{J^T} (x_{t+1}, x^*) + D_{\phi^*} (\nabla J^T(x_{t+1}), \nabla J^T(x^*)) \\ + D_{\phi^*} (\nabla J^T(x^*), \nabla J^T(x_t)) \leq D_{J^T} (x_t, x^*) \end{aligned}$$

**Proof** We first note that by Lemma 11, the assumptions of the theorem guarantee that  $J^T$  is Legendre convex. Now let  $k = \tilde{\phi}^*$ ,  $f = J^T$  and  $L^* = 1$  in the statement of (Maddison et al., 2019, Lemma 4.6). Noting that a single iterate of  $x \in \mathcal{P}_x^T$  coincides exactly with the iterates of (Maddison et al., 2019, Algorithm 1.1) with  $L^* = 1$ , we apply (Maddison et al., 2019, Lemma 4.6) to obtain the desired result. ■

#### B.4. Descent on Convex Objectives

**Lemma 25** *Let  $f$  be convex, differentiable and  $\phi$  be a symmetric positive-definite quadratic function, and consider  $x \in \mathcal{P}^\infty$ . Let us define the Lagrangian*

$$\mathcal{L}_t = f(x_{t+1}) + \phi(\Delta x_t).$$

Then  $\mathcal{L}_t$  is a non-increasing sequence.

**Proof** We first note that by Lemma 11, the assumptions of the theorem guarantee that  $J^\infty$  is Legendre convex. We first note that since  $\phi$  is quadratic, we have the properties that  $\phi(x) = \phi(-x) = \frac{1}{2} \langle x, \nabla \phi(x) \rangle$  for all  $x \in \mathcal{X}$ , and that  $\phi^*$  is quadratic as well. We begin by separating the expression for  $\mathcal{L}_t$  into two parts and showing that each is monotone. Using the short-hand notation  $\nabla J^\infty(x_t) = \nabla J_t^\infty$ , we separate the expression as

$$\begin{aligned} \mathcal{L}_t &= \{f(x_{t+1}) - \phi^*(\nabla J_{t+1}^\infty)\} + \phi^*(\nabla J_{t+1}^\infty) + \phi(\Delta x_t) \\ &= \mathcal{H}_{t+1} + \phi^*(\nabla J_{t+1}^\infty) + \phi(\Delta x_t). \end{aligned} \tag{30}$$

Note that since  $-\nabla \phi(\Delta x_t) = \nabla J_t^\infty$  and that and hence that  $\Delta x_t = \nabla \phi^*(-\nabla J_t^\infty)$ , we have

$$\begin{aligned} \phi(\Delta x_t) &= \langle \nabla \phi(\Delta x_t), \Delta x_t \rangle - \phi^*(\nabla \phi(\Delta x_t)) \\ &= \langle \nabla \phi(\nabla \phi^*(-\nabla J_t^\infty)), \nabla \phi^*(-\nabla J_t^\infty) \rangle - \phi^*(\nabla \phi(\Delta x_t)) \\ &= \langle \nabla J_t^\infty, \nabla \phi^*(\nabla J_t^\infty) \rangle - \phi^*(\nabla J_t^\infty). \end{aligned}$$

Since  $\phi^*$  is quadratic,  $\phi^*(x) = \frac{1}{2} \langle \nabla \phi^*(x), x \rangle$ , and hence we conclude that

$$\phi(\Delta x_t) = \phi^*(\nabla J_t).$$

By Lemma 24-ii, we have that  $\phi^*(\nabla J_{t+1}) + \phi(\Delta x_t) = \phi^*(\nabla J_{t+1}^\infty) + \phi^*(\nabla J_t^\infty)$  is a monotone non-increasing sequence, and hence the latter two terms of equation (30) are monotone non-increasing.



To show that  $\mathcal{H}_t$  is non-increasing, we use that  $\Delta x_t = \nabla \phi^*(-\nabla J_t^\infty)$  and that  $\nabla J_t^\infty = \nabla f(x_{t+1}) + \nabla J_{t+1}^\infty$  to compute

$$\begin{aligned}
\mathcal{H}_{t+1} - \mathcal{H}_t &= \{f(x_{t+1}) - f(x_t)\} - \{\phi^*(\nabla J_{t+1}^\infty) - \phi^*(\nabla J_t^\infty)\} \\
&= -\{D_f(x_t, x_{t+1}) + D_{\phi^*}(\nabla J_{t+1}^\infty, \nabla J_t^\infty)\} - \langle \Delta x_t, \nabla f(x_{t+1}) \rangle + \langle \nabla J_{t+1}^\infty - \nabla J_t^\infty, \phi^*(\nabla J_t^\infty) \rangle \\
&= -\{D_f(x_t, x_{t+1}) + D_{\phi^*}(\nabla J_{t+1}^\infty, \nabla J_t^\infty)\} - \langle \Delta x_t, \nabla f(x_{t+1}) \rangle - \langle \nabla J_{t+1}^\infty - \nabla J_t^\infty, \Delta x_t \rangle \\
&= -\{D_f(x_t, x_{t+1}) + D_{\phi^*}(\nabla J_{t+1}^\infty, \nabla J_t^\infty)\} \\
&\leq 0.
\end{aligned}$$

Hence  $\mathcal{H}_t$  is non-increasing and we have the desired result. ■

## Appendix C. Proofs for Section 2

Most of the results in this section, rely on establishing the interchangeability of the  $\lim$  and  $\operatorname{argmin}$ . In general, these operations do not commute. The theory of  $\Gamma$ -convergence, introduced in [De Giorgi et al. \(1988\)](#), is designed specifically to address these types of problems. We briefly overview this theory before applying its to derive proofs of the results in [Section 2](#).

### C.1. Overview of $\Gamma$ -Convergence Theory

Fix a metric space  $(X, d)$  and consider a functional  $F : X \rightarrow (-\infty, \infty] \cup \{\infty\}$  which we would like to minimize on  $X$ . When it exists, a minimum of  $F$  describes an  $x \in X$  achieving the lowest value of  $F$ , or equivalently, it represents the lowest point on  $F$ 's epigraph, as defined by  $\operatorname{epi}(F) := \{(x, r) \in X \times (-\infty, \infty] : r \geq F(x)\}$ ; this set describes all points on or above  $F$ 's graph.

We can expect  $F$  to admit a minimum precisely if the sub-level set of  $\operatorname{epi}(F)$  are not too large, if limits of points in  $\operatorname{epi}(F)$  stay within  $\operatorname{epi}(F)$ , and if  $\operatorname{epi}(F)$  does not trail off to  $-\infty$  at some point. These three conditions respectively describe the intuition behind the need for  $F$ 's coercivity, lower semi-continuity, and boundedness from below. Formally, a function  $F$  is said to be *coercive* if its sub-level sets  $\operatorname{Lev}_s(F) := \{x \in X : F(x) \leq s\}$ , for every  $s \in \mathbb{R}$ , are compact in  $X$ ,  $F$  is said to be *lower semi-continuous* if  $\operatorname{epi}(F)$  is a closed subset of  $X \times (-\infty, \infty]$ , and it is bounded from below if  $F(x) \geq M$  for some  $M \in \mathbb{R}$ . The *Direct Method* of [Tonelli \(1961\)](#), guarantees that together these three conditions are sufficient for  $F$  to be minimized over  $X$ .

Suppose now that  $F$  can be described as the point-wise limit of a sequence of functionals  $(F_n)_{n \in \mathbb{N}}$  on  $X$ . It is tempting, to solve  $\operatorname{argmin}_{x \in X} F(x)$  by interchanging the limit and  $\operatorname{argmin}$  operations via

$$\lim_{n \uparrow \infty} \operatorname{argmin}_{x \in X} F_n(x) = \operatorname{argmin}_{x \in X} \lim_{n \uparrow \infty} F_n(x); \quad (31)$$

however, (31) is generally false even when the convergence of  $F_n$  to  $F$  is uniform (see [Dal Maso \(1993\)](#)). This is because any of the three aforementioned properties can fail for the limiting functional  $F$  even if they hold for each of the  $F_n$ .

Any inconsistency with the lower semi-continuity of the functionals  $(F_n)_{n \in \mathbb{N}}$  and  $F$  is avoided when the epigraphs  $\operatorname{epi}(F_n)$  converge to  $F$ 's epigraph. There are various modes of convergence of sets, see [Aubin and Frankowska \(2009\)](#); however, the correct notion of set-convergence here is *Kuratowski convergence*, introduced in [Kuratowski \(1966\)](#), which intuitively describes the set of accumulation points of  $\operatorname{epi}(F_n)$  and is defined to be  $\{(x, t) \in X \times (-\infty, \infty] : \limsup_{n \uparrow \infty} d((x, t), \operatorname{epi}(F_n)) = 0\}$ , whenever that set coincides with  $\{(x, t) \in X \times (-\infty, \infty] : \liminf_{n \uparrow \infty} d((x, t), \operatorname{epi}(F_n)) = 0\}$ . If this happens, we say that the sequence of functionals  $(F_n)_{n \in \mathbb{N}}$   $\Gamma$ -converges to  $F$  and we write  $\Gamma - \lim_{n \uparrow \infty} F_n = F$ .

**Remark 26**  $\Gamma$ -convergence can still be formulated when  $X$  is not a metric space, for example, this is the case when  $X$  is an infinite-dimensional Banach space such as  $\mathcal{A}_x^{\infty;0}$  equipped with its weak topology. For details on this general case, we point the reader to ([Dal Maso, 1993, Chapter 4](#)).

Analogously, any inconsistency with the coercivity of the functionals  $(F_n)_{n \in \mathbb{N}}$  is avoided when the sequence is *equi-coercive*. This means that the sub-level sets are uniformly small in the same places; mathematically, this means that for every  $s \in \mathbb{R}$  there exists a compact subset  $K$  of  $X$  containing each  $\operatorname{Lev}_s(F_n)$ , for  $n \in \mathbb{N}$ .

When the sequence  $(F_n)_{n \in \mathbb{N}}$  is both equi-coercive and its  $\Gamma$ -limit is  $F$ , then the *Fundamental Theorem of  $\Gamma$ -Convergence* (Braides, 2014, Theorem 2.1) is a sequential extension of Tonelli's Direct method. Accordingly, the result guarantees that  $F$  and each  $F_n$  is minimized over  $X$  and that any accumulation point of these minima are a minimizer of  $F$ .

## C.2. Proof of Theorem 1

We establish Theorem 1 using Tonelli's Direct method, described in the previous section. This amounts to showing that the functional  $\mathcal{R}_T$ , for each  $T \in \mathbb{N}$ , is lower semi-continuous, coercive, and bounded-below.

We simplify our task by breaking up the regret functional into the sum  $\mathcal{R}_T = F_T + \Phi_T$ , where the functionals  $\Phi_T$  and  $F_T$  on  $\mathcal{A}_x^T$  are defined by

$$\Phi_T : \mathbf{x} \mapsto \sum_{t=1}^T \phi(\Delta x_t) \text{ and } F_T : \mathbf{x} \mapsto \sum_{t=1}^T f(x_t) - f^*.$$

The convenience arises from the fact that we can establish each of these two functionals' individual properties, since they rely on different assumptions, before combining them back together to infer the relevant properties of  $\mathcal{R}_T$ . Accordingly, our task is divided into establishing a sequence of lemmas, each dedicated to showing a different property of  $F_T$  or  $\Phi_T$ , at which point theorems' proof reduces Tonelli's method.

### C.2.1. AUXILIARY LEMMAS

We now note that  $\mathcal{A}_x^{\infty,0}$  will be equipped with the topology generated by the norm  $\mathbf{x} \mapsto (\sum_{t=1}^{\infty} \|\Delta x_t\|^p)^{\frac{1}{p}}$  and we note that  $\mathcal{A}_x^{\infty,0}$  is a Banach space which is isometrically isomorphic to  $\ell^p(X)$  via the map  $\mathbf{x} \mapsto \Delta \mathbf{x}$ . We also observe that  $\mathcal{A}_x^{\infty,\alpha} \subset \mathcal{A}_x^{\infty,0} \subset \mathcal{A}_x^{\infty}$  for every  $\alpha > 0$ .

**Lemma 27 (Regularity of Algorithmic Penalty Function)** *Under Assumption 2, for each  $T \in \mathbb{N} \cup \{\infty\}$ ,  $\Phi_T$  is coercive on each  $\mathcal{A}_x^T$  (resp.  $\mathcal{A}_x^{\infty,0}$  when  $T = \infty$ ). Moreover, it is weakly lower semi-continuous, not identically  $\infty$ , and bounded-below by 0 on each  $\mathcal{A}_x^T$  with  $T < \infty$  (resp. on  $\mathcal{A}_x^{\infty,\alpha}$  for  $\alpha \geq 0$  when  $T = \infty$ ). In particular, it is lower semi-continuous on  $\mathcal{A}_x^{\infty,\alpha}$  for  $\alpha \geq 0$ . Furthermore, the following hold:*

1.  $\Phi_{\infty}$  is coercive on  $\mathcal{A}_x^{\infty,\alpha}$ , for every  $\alpha > 0$ ,
2.  $\Phi_{\infty}$  is weakly lower semi-continuous, and weakly coercive on  $\mathcal{A}_x^{\infty,0}$ .

### Proof

Since the weak lower semi-continuity of any function from  $\mathcal{A}_x^{\infty,0}$  to  $(-\infty, \infty]$  implies its lower semi-continuity (l.s.c.), then it is enough for us to establish the former on  $\mathcal{A}_x^{\infty,0}$ . This is because, the restriction of weakly lower semi-continuous (weakly l.s.c.) functions to closed sets of a topological space preserves weak lower semi-continuity; thus, it is enough to show that  $\Phi_{\infty}$  is l.s.c. on  $\mathcal{A}_x^{\infty,0}$  to conclude that it must also be l.s.c. on each  $\mathcal{A}_x^{\infty,\alpha}$  for  $\alpha \geq 0$ .

We begin with the following remark. Since  $\mathcal{X}$  is a finite-dimensional normed space, it must admit a Hamel basis  $\{e_n\}_{n=1}^{\dim(\mathcal{X})}$  such that every  $x \in \mathcal{X}$  is uniquely expressed as  $x = \sum_{n=1}^{\dim(\mathcal{X})} \beta_n e_n$ . In particular, for  $1 \leq n \leq \dim(\mathcal{X})$ , the map  $x \mapsto \pi_n(x) = \beta_n$  mapping  $x$  to the coefficient  $\beta_n$  in its Hamel basis expansion is a bounded linear map and it is therefore, continuous. Next, for every

$T \in \mathbb{N}$ , define the map  $p_T : \mathcal{A}_x^{\infty:0} \rightarrow \mathcal{X}$ , taking  $\mathbf{x}$  to its value at the  $T$ -th component,  $x_T$ . By definition this is a bounded linear map. Hence, for every  $T \in \mathbb{N}$  and every  $1 \leq n \leq \dim(\mathcal{X})$ , the composition  $p_T \circ \pi_n : \mathcal{A}_x^{\infty:0} \rightarrow \mathbb{R}$  is a bounded linear functional. Moreover, for every  $T \in \mathbb{N}$ , we have the representation  $p_T = \sum_{n=1}^{\dim(\mathcal{X})} p_T \circ \pi_n e_n$ .

Since  $\mathcal{A}_x^{\infty:0}$  is a Banach space then, by definition, in the weak topology on  $\mathcal{A}_x^{\infty:0}$  all bounded linear functionals are continuous. Thus, for  $T \in \mathbb{N}$  and every  $1 \leq n \leq \dim(\mathcal{X})$ , the map sending  $p_T \circ \pi_n$  is bounded and linear then it is weakly continuous on  $\mathcal{A}_x^{\infty:0}$ . Now, since the sum of weakly lower continuous functions is again weakly lower continuous and since  $p_T = \sum_{n=1}^{\dim(X)} p_T e_n$ , then  $p_T$  is weakly-to-strong continuous<sup>5</sup> from  $\mathcal{A}_x^{\infty:0}$  to  $\mathcal{X}$ .

Next, since the pre-composition of a lower semi-continuous function by a weak-to-strong continuous map is weakly lower semi-continuous then, for every  $T \in \mathbb{N}$ , the map  $\phi \circ p_T$  is weakly lower semi-continuous. By (Dal Maso, 1993, Proposition 1.9) the sum of weakly lower semi-continuous functions is again weakly lower semi-continuous and therefore, for every  $T \in \mathbb{N}$ ,  $\Phi_T : \mathcal{A}_x^{\infty:0} \rightarrow (-\infty, \infty]$  is weakly lower semi-continuous. Since the point-wise supremum of a family of weakly lower semi-continuous functions is weakly lower semi-continuous, see (Dal Maso, 1993, Proposition 1.8), then  $\Phi_\infty = \sup_{T \in \mathbb{N}} \Phi_T$  is weakly lower semi-continuous on  $\mathcal{A}_x^{\infty:0}$ .

Now, consider the sequence  $\mathbf{x}^*$  defined by  $x_t^* = x$ . Since  $\phi(0) = 0$  and  $\Delta x_t^* = 0$  for every  $t \in \mathbb{N}$  then  $\Phi$  is not identically  $\infty$  on  $\mathcal{A}_x^{\infty:0}$ . Moreover, by construction  $\Phi \geq 0$ .

By (Dal Maso, 1993, Definition 1.12), it  $\Phi$  is coercive on  $\mathcal{A}_x^{\infty:\alpha}$  if its sub-level sets are compact; i.e. for every  $s \geq 0$  the sub-level set  $\text{Lev}_s(\Phi) := \{\mathbf{x} \in \mathcal{A}_x^{\infty:\alpha} : \Phi(\mathbf{x}) \leq s\}$ . Assumption 2 implies that for every  $s \geq 0$ , we have the inclusion

$$\text{Lev}_s(\Phi) \subseteq \left\{ \mathbf{x} \in \mathcal{A}_x^{\infty:0} : c \sum_{t=1}^{\infty} n^\alpha \|\Delta x_t\|^p \leq s \right\}. \quad (32)$$

Since  $\mathcal{A}_x^{\infty:0}$  is isometrically isomorphic to  $\ell^p(X)$  via the map  $\mathbf{x} \mapsto \Delta \mathbf{x}$  then Grothendieck's compactness principle (here we use the formulation of (Diestel, 1984, Exercises 1.6)) implies that the right-hand side of (32) is a compact subset of the Banach space  $\mathcal{A}_x^{\infty:0}$ ; which, by construction, is a subset of  $\mathcal{A}_x^{\infty:\alpha}$ . Since we assume that  $\Phi$  is lower semi-continuous, then the sub-level set  $\text{Lev}_s(\Phi)$  is closed and in particular it is a closed subset of the compact set  $\{\mathbf{x} \in \mathcal{A}_x^{\infty:0} : c \sum_{t=1}^{\infty} n^\alpha \|\Delta x_t\|^p \leq s\}$ ; thus  $\text{Lev}_s(\Phi)$  is compact by (Munkres, 2000, Theorem 26.2).

Consider the case where  $\phi$  is convex. Let  $\mathbf{x}, \mathbf{y} \in \mathcal{A}_x^{\infty:0}$  and  $\rho \in [0, 1]$  then the convexity of  $\phi$  and the linearity of  $\Delta$  imply that

$$\begin{aligned} \sum_{t=1}^{\infty} \phi(\Delta(\rho x_t + (1-\rho)y_t)) &= \sum_{t=1}^{\infty} \phi(\rho \Delta x_t + (1-\rho)\Delta y_t) \\ &\leq \sum_{t=1}^{\infty} \rho \phi(\Delta x_t) + (1-\rho)\phi(\Delta y_t); \end{aligned}$$

hence  $\Phi_\infty : \mathbf{x} \mapsto \sum_{t=1}^{\infty} \phi(\Delta x_t)$  is convex on  $\mathcal{A}_x^{\infty:0}$ .

5. Let  $X$  be a topological space endowed with the weak topology and  $Y$  be a normed space. We say that a map  $f : X \rightarrow Y$  is *weak-to-strong* continuous if it satisfies the usual definition of continuity; that is, for any point  $x \in X$  and neighborhood  $\epsilon_y$  of  $f(x) = y \in Y$ , there exists a neighborhood  $\delta_x$  of  $x$  such that  $f(\delta_x) \subseteq \epsilon_y$ .

Next, since there exists a constant  $c > 0$  satisfying  $c\|x\|^p \leq \phi(x)$  for every  $x \in \mathcal{X}$  and since  $\mathcal{X}$  is a linear space then  $\Delta x_t \in X$  for every  $\mathbf{x} \in \mathcal{A}_x^{\infty:0}$ . Thus,  $c\|\Delta x_t\|^p \leq \phi(\Delta x_t)$  for every  $\mathbf{x} \in \mathcal{A}_x^{\infty:0}$  and therefore

$$c\|\mathbf{x}\|_{\mathcal{A}_x^{\infty:0}}^p = c \sum_{t=1}^{\infty} \|\Delta x_t\|^p \leq \sum_{t=1}^{\infty} \phi(\Delta x_t) = \Phi(\mathbf{x});$$

whence  $\Phi_{\infty}$  is weakly coercive, since if  $\|\mathbf{x}\|_{\mathcal{A}_x^{\infty:0}} \rightarrow \infty$  implies that  $\Phi_{\infty}(\mathbf{x}) \rightarrow \infty$ .  $\blacksquare$

**Lemma 28 (Regularity of Unregularized Regret)** *Under Assumption 1, for every  $T \in \mathbb{N} \cup \{\infty\}$ , the functions  $F_T(\mathbf{x}) := \sum_{t=1}^T f(x_t)$  are weakly lower semi-continuous on  $\mathcal{A}_x^{\infty:\alpha}$ , for every  $\alpha \geq 0$ . In particular, they are lower semi-continuous on  $\mathcal{A}_x^{\infty:\alpha}$ , for every  $\alpha \geq 0$ . Moreover, each  $F_T$  is weakly lower semi-continuous on  $\mathcal{A}_x^{\infty:0}$ .*

**Proof** Analogously to the proof of Lemma 27, it is enough for us to show that, for each  $T \in \mathbb{N}$ , the function  $F_T$  is weakly lower semi-continuous on  $\mathcal{A}_x^{\infty:0}$  to conclude that it is both lower semi-continuous and weakly lower semi-continuous on each  $\mathcal{A}_x^{\infty:\alpha}$  for every  $\alpha \geq 0$ .

As in the proof of Lemma 27, we know that each map  $p_n : \mathcal{A}_x^{\infty:0} \rightarrow \mathcal{X}$  weakly continuous. Since  $f$  is lower semi-continuous then, for each  $T \in \mathbb{N}$ , the composition  $f \circ p_T : \mathcal{A}_x^{\infty:0} \rightarrow (-\infty, \infty]$  is weakly lower semi-continuous. Applying (Dal Maso, 1993, Proposition 1.9) we conclude that, for every  $T \in \mathbb{N}$ ,  $F_T$  is also weakly lower semi-continuous. By (Dal Maso, 1993, Proposition 1.8), since the supremum of any family of lower semi-continuous functions is itself lower semi-continuous; thus,

$$F_{\infty}(\mathbf{x}) = \sum_{t=1}^{\infty} f(x_t) - f(x^*) = \sup_{T \in \mathbb{N}} \sum_{t=1}^T f(x_t) - f(x^*) = \sup_{T \in \mathbb{N}} F_T(\mathbf{x}),$$

is weakly lower semi-continuous.

If in addition, if  $f$  is convex then arguing similarly to the proof of Lemma 27 we find that each  $F_T$ , for  $T \in \mathbb{N} \cup \{\infty\}$ , is convex. Since the point-wise supremum of any family of convex functions is itself convex, then  $F_{\infty}$  is also convex if  $f$  is convex.  $\blacksquare$

**Lemma 29 (Regularity of Regret Functionals)** *Under Assumptions 1 and 2, for every  $T \in \mathbb{N}$ , the regret functionals  $\mathcal{R}_T = F_T + \Phi_T$  are not identically  $\infty$ , bounded-below by 0, and weakly coercive on  $\mathcal{A}_x^{\infty:0}$  and weakly lower semi-continuous on  $\mathcal{A}_x^{\infty:0}$ .*

**Proof** Since  $f(x) \geq f^*$  for every  $x \in X$  and since  $\Delta x_t \in X$  for any  $\mathbf{x} \in \mathcal{A}_x^{\infty}$  then  $f(x_t) - f^* \geq 0$  for every  $\mathbf{x} \in \mathcal{A}_x^{\infty}$ . Thus, each  $F_T$  takes values in  $[0, \infty]$ . Moreover, since  $\phi(x^* - x) < \infty$  then the sequence  $\mathbf{z}$  given by  $z_0 = x, z_t = x^*$  for  $t > 0$  satisfies  $(F_T + \Phi_T)(\mathbf{z}) = f(x) + \phi(x^* - x) < \infty$ . Hence, each  $F_T + \Phi_T$  is not identically  $\infty$  on  $\mathcal{A}_x^{\infty:0}$ . Moreover, since, for every  $T \in \mathbb{N}$ ,  $F_T$  and  $\Phi_T$  are weakly lower semi-continuous on  $\mathcal{A}_x^{\infty:0}$  then (Dal Maso, 1993, Proposition 1.9) implies that  $\mathcal{R}_T$  is weakly lower semi-continuous.

It remains only to demonstrate the weak coercivity of the regret functionals. We denote the sub-level set at  $s \in \mathbb{R}$  of any  $G : \mathcal{A}_x^{\infty:0} \rightarrow \mathbb{R} \cup \{\infty\}$ , by  $\text{Lev}_s(G) := \{\mathbf{y} \in \mathcal{A}_x^{\infty:0} : G(\mathbf{y}) \leq s\}$ . Since each  $F_T$  and  $\Phi_T$  are bounded-below by 0, then  $(F_T + \Phi_T)(\mathbf{y}) \geq \Phi_T(\mathbf{y})$  for every  $\mathbf{y} \in \mathcal{A}_x^{\infty:0}$ . Hence,  $\text{Lev}_s(F_T + \Phi_T) \subseteq \text{Lev}_s(\Phi_T)$  for every  $s > 0$  and every  $T \in \mathbb{N}$ . Moreover, since  $F_T$

and  $\Phi_T$  are lower semi-continuous then, for every  $s > 0$  and every  $T \in \mathbb{N}$ , each  $\text{Lev}_s(F + \Phi)$  (resp.  $\text{Lev}_s(F_T + \Phi)$ ) is a weakly closed subset of  $\text{Lev}_s(\Phi_T)$ . Since  $\Phi_T$  is weakly coercive then, by definition, each  $\text{Lev}_s(\Phi)$  is weakly compact. Since the weak-topology is metrizable on every compact subset thereof, and since every closed subset of a compact set is itself compact in a metric space then  $\text{Lev}_s(F_T + \Phi_T)$  is weakly-compact; therefore,  $\mathcal{R}_T = F_T + \Phi_T$  is weakly coercive on  $\mathcal{A}_x^{\infty:0}$ . ■

### C.2.2. PROOF OF THEOREM 1

**Proof** By Lemmas 28 and 27 each  $F_T$  is lower semi-continuous on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. weakly lower semi-continuous on  $\mathcal{A}_x^{\infty:0}$ ). By construction,  $F_T \geq 0$  and by Lemma 27,  $\Phi_T$  is bounded-below and not identically  $\infty$ . By Assumption 1, for  $T \in \mathbb{N}$ ,  $F_T$  is not identically  $\infty$  either since the sequence  $\mathbf{x}^T$  defined by  $\mathbf{x}_0^T := x$ ,  $\mathbf{x}_t^T := x^*$  for  $0 < t \leq T$ , satisfies

$$(F_T + \Phi_T)(\mathbf{x}^T) = f(x) + \phi(x^* - x) < \infty. \quad (33)$$

By Lemma 29, for every  $T \in \mathbb{N}$ ,  $F_T + \Phi_T$ . Therefore, for every  $T \in \mathbb{N}$ ,  $F_T + \Phi_T$  is bounded-below by 0, lower semi-continuous, and coercive on  $\mathcal{A}_x^{\infty:\alpha}$ ; hence, by (Dal Maso, 1993, Theorem 1.15) each  $\mathcal{P}_x^{\infty:0}$  is non-empty. Now since  $\mathcal{R}_T$  is point-wise monotonically increasing in  $T$ , then  $\mathbf{x} = (x_t)_{t \in \mathbb{N}} \in \mathcal{P}^{\infty:0}$  only if the halted sequence

$$\tilde{\mathbf{x}} \triangleq \begin{cases} x_t & : t \leq T \\ 0 & : \text{else} \end{cases},$$

belongs to  $\mathcal{P}_x^T$ . Hence,  $\mathcal{P}_x^T$  is non-empty. Moreover, for each  $T \in \mathbb{N}$ , since  $F_T + \Phi_T$  is not identically  $\infty$  then by definition of  $\mathcal{P}_x^{\infty:0}$  and (33) we have any  $\mathbf{x}^{*:T} \in \mathcal{P}_x^T$  satisfies

$$(F_T + \Phi_T)(\mathbf{x}^{*:T}) \leq (F_T + \Phi_T)(\mathbf{x}^T) < \infty.$$

Hence,  $\mathcal{P}_x^T$  is non-empty and every element therein has finite regret. ■

### C.3. Proof of Theorem 2

The results of Theorems 2 and 4 both follow as consequences of the Fundamental Theorem of  $\Gamma$ -convergence, described in Section C.1. As in the proof of Theorem 1, we begin by establishing the required properties which will allow to apply this result. These amount to showing that the regret functional  $\mathcal{R}_\infty$  is the  $\Gamma$ -limit of an equi-coercive family of functionals expressing the finite-time regret-optimal algorithm selection problem posed on each  $\mathcal{A}_x^T$  embedded within the larger space  $\mathcal{A}_x^{\infty:0}$ .

We now precisely define the process of embedding the finite-time-horizon control problems as control problems on  $\mathcal{A}_x^{\infty:0}$ . For any  $T \in \mathbb{N} \cup \{\infty\}$ , we introduce the *indicator functions*<sup>6</sup>  $\chi_{\mathcal{A}_x^T} : \mathcal{A}_x^{\infty:0} \rightarrow \mathbb{R} \cup \{\infty\}$  of the subsets  $\mathcal{A}_x^T \subseteq \mathcal{A}_x^{\infty:0}$ , defined by  $\chi_{\mathcal{A}_x^T}(\mathbf{x}) = 0$  only if  $\mathbf{x} \in \mathcal{A}_x^T$  and

6. We note here that we use the *indicator function* terminology belonging to convex-analysis and not that of probability theory.

$\infty$  otherwise. These functions allow us to express the constrained optimization problem (3) as an unconstrained optimization problem on all of  $\mathcal{A}_x^{\infty:0}$  through

$$\min_{\mathbf{x} \in \mathcal{A}_x^{\infty:0}} \mathcal{R}_T(\mathbf{x}) + \chi_{\mathcal{A}_x^T}(\mathbf{x}). \quad (34)$$

Thus, Theorem 4 can be stated concisely as a guarantee that the argmin and lim operations can be interchanged regret-optimization problem's time-horizon becomes unbounded. Hence, we seek to show that

$$\lim_{T \uparrow \infty} \operatorname{argmin}_{\mathbf{x} \in \mathcal{A}_x^{\infty:0}} \mathcal{R}_T(\mathbf{x}) + \chi_{\mathcal{A}_x^T}(\mathbf{x}) \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{A}_x^{\infty:0}} \lim_{T \uparrow \infty} \mathcal{R}_T(\mathbf{x}) + \chi_{\mathcal{A}_x^T}(\mathbf{x}). \quad (35)$$

As with the proof of Theorem 1, our approach is to first establish the relevant properties of the sequence of functionals (34) through various lemmas. These properties include lower semi-continuity and lower-boundedness on  $\mathcal{A}_x^{\infty:0}$ , their  $\Gamma$ -convergence to  $\mathcal{R}_\infty$ , and their equi-coercivity. The proof of the aforementioned results then follows from the Fundamental Theorem of  $\Gamma$ -convergence ((Braides, 2014, Theorem 2.1)) which guarantees that (35) holds.

### C.3.1. AUXILIARY LEMMAS FOR THEOREMS 2 AND 4

**Lemma 30** *Under Assumption 1, the sequence  $\{F_T + \Phi_T\}_{T \in \mathbb{N}}$   $\Gamma$ -converges to  $\mathcal{R}_\infty$  on  $\mathcal{A}_x^{\infty:\alpha}$ . Moreover,  $\{F_T + \Phi_T\}_{T \in \mathbb{N}}$  also  $\Gamma$ -converges to  $\mathcal{R}_\infty$  on  $\mathcal{A}_x^{\infty:0}$  in the weak topology.*

**Proof** First note that since  $f$  is bounded below, by the minimum values  $f^* \in \mathbb{R}$ , then  $f(x) - f^* \geq 0$  for every  $x \in \mathcal{X}$ . By Assumption 2  $\phi \geq 0$  for every  $x \in \mathcal{X}$ . Thus, for each  $T \in \mathbb{N}$  and each  $\mathbf{x} \in \mathcal{A}_x^\infty$ , we that

$$(F_T + \Phi_T)(\mathbf{x}) = \sum_{t=1}^T f(x_t) - f^* + \phi(\Delta x_t) \geq \sum_{t=1}^T 0 + 0 = 0; \quad (36)$$

hence,  $(F_T + \Phi_T)(\mathbf{x})$  is non-negatively valued. Next, observe that for any  $\mathbf{x} \in \mathcal{A}_x^\infty$  the sequence of real-numbers  $\{F_T + \Phi_T(\mathbf{x})\}_{T \in \mathbb{N}}$  is non-negative monotonically increasing as a function of  $T$ . Hence, the Monotone Convergence Theorem implies that, for each  $T \in \mathbb{N}$  and each  $\mathbf{x} \in \mathcal{A}_x^\infty$ ,  $\{(F_T + \Phi_T)(\mathbf{x})\}_{T \in \mathbb{N}}$  converges point-wise to  $(F_\infty + \Phi_\infty)(\mathbf{x}) = \mathcal{R}_\infty(\mathbf{x})$ ; moreover the convergence is monotone. In particular, the sequence of functionals  $\{F_T + \Phi_T\}_{T \in \mathbb{N}}$  is monotonically increasing and converges point-wise to  $\mathcal{R}_\infty$  on all of  $\mathcal{A}_x^\infty$ .

We may therefore apply (Dal Maso, 1993, Proposition 5.4) to conclude that the lower semi-continuous relaxation  $\mathcal{R}_\infty^{lsc}$  of  $\mathcal{R}_\infty$  on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. on  $\mathcal{A}_x^{\infty:0}$  for the weak topology) and therefore is the  $\Gamma$ -limit of  $\{F_T + \Phi_T\}_{T \in \mathbb{N}}$  on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. on  $\mathcal{A}_x^{\infty:0}$  for the weak topology). Lemma 28 guarantees that  $F_\infty$  is weakly lower semi-continuous and Lemma 27 guarantees that  $\Phi_\infty$  is lower semi-continuous on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. weakly lower semi-continuous on  $\mathcal{A}_x^{\infty:0}$ ). Since  $F_\infty, \Phi_\infty > -\infty$  then the sum  $R_\infty = F_\infty + \Phi_\infty$  is well-defined on all of  $\mathcal{X}$  and therefore (Dal Maso, 1993, Proposition 1.9) implies that  $\mathbb{R}_\infty$  is itself lower semi-continuous (resp. weakly lower semi-continuous on  $\mathcal{A}_x^{\infty:0}$ ). Hence,  $\mathcal{R}_\infty^{lsc} = \mathcal{R}_\infty$  and therefore  $\mathcal{R}_\infty$  is the  $\Gamma$ -limit of  $\{F_T + \Phi_T\}_{T \in \mathbb{N}}$  on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. on  $\mathcal{A}_x^{\infty:0}$  in the weak topology).  $\blacksquare$

We continue our analysis by exhibiting some helpful properties of these indicator functions.

**Lemma 31 (Regularity of Indicator Functions  $\chi_{\mathcal{A}_x^T}$ )** For any  $T \in \mathbb{N}$ ,  $\mathcal{A}_x^T$  is convex and weakly closed in  $\mathcal{A}_x^{\infty:0}$ . Thus, the function  $\chi_{\mathcal{A}_x^T} : \mathcal{A}_x^{\infty:0} \rightarrow \mathbb{R} \cup \{\infty\}$  is weakly lower semi-continuous and convex. In particular, the function  $\chi_{\mathcal{A}_x^T} : \mathcal{A}_x^{\infty:0} \rightarrow \mathbb{R} \cup \{\infty\}$  is lower semi-continuous and convex on  $\mathcal{A}_x^{\infty:\alpha}$ .

**Proof** Fix  $k \in \mathbb{R}$  and  $\mathbf{x}, \mathbf{y} \in \mathcal{A}_x^T$ . By linearity of  $\Delta$  we have that  $\Delta(\mathbf{x} + k\mathbf{y})_u = 0$  for every  $u \geq T$ ; thus,  $\mathcal{A}_x^T$  is a linear space and it is therefore convex. Thus,  $\chi_{\mathcal{A}_x^T}$  is convex; for each  $T \in \mathbb{N}$ .

Let  $\{\mathbf{x}^n\}_{n \in \mathbb{N}}$  be a sequence of algorithms in  $\mathcal{A}_x^T$  converging to some algorithm  $\mathbf{x} \in \mathcal{A}_x^{\infty:0}$ . Then,  $\sum_{t \geq T} \|\Delta x_t^n - \Delta x_t\|^p \rightarrow 0$ . However, by the definition of  $\mathcal{A}_x^T$ , for all  $n \in \mathbb{N}$  we have that  $x_t^n = 0$  if  $t \geq T$ . Therefore,

$$\sum_{t \geq T} \|\Delta x_t^n - \Delta x_t\|^p = \sum_{t \geq T} \|\Delta x_t\|^p;$$

hence  $\sum_{t \geq T} \|\Delta x_t^n - \Delta x_t\|^p \mapsto 0$  only if  $\Delta x_t = 0$  for every  $t \geq T$ . Thus,  $\mathcal{A}_x^T$  is a closed convex subset of  $\mathcal{A}_x^{\infty:0}$ ; for every  $T \in \mathbb{N}$ . Now, by (Conway, 1990, 1.5 Corollary) we conclude that  $\mathcal{A}_x^T$  is weakly closed.

Since, for every  $T \in \mathbb{N}$ , the set  $\mathcal{A}_x^T$  is closed and convex, then (Dal Maso, 1993, Example 3.4) implies that each  $\chi_{\mathcal{A}_x^T}$  is lower semi-continuous. Since it is lower semi-continuous and convex (Dal Maso, 1993, Proposition 1.18) implies that  $\chi_{\mathcal{A}_x^T}$  is also weakly lower semi-continuous.  $\blacksquare$

In what follows, we frequently make use of the notion of *continuous convergence* (see (Kuratoski, 1968, Chapter 20, Section 6)). In general, this mode of continuous is strictly stronger than point-wise convergence but strictly weaker than uniform convergence. Continuous convergent functionals are abundant enough to exhibit while simultaneously being regular enough to control.

**Definition 32 (Continuous Convergence)** A sequence  $\{F^n\}_{n \in \mathbb{N}}$  from a topological space  $Z$  to  $\mathbb{R} \cup \{\infty\}$  converge continuously to  $F : Z \rightarrow \mathbb{R} \cup \{\infty\}$  if, for every  $z \in Z$  and every open neighbourhood  $U_{F(z)} \subseteq \mathbb{R} \cup \{\infty\}$  of  $F(z)$ , there exists some  $N_z \in \mathbb{N}$  and some open neighbourhood  $U_z \subseteq Z$  of  $z$  for which

$$F_n(y) \in U_{F(z)},$$

for every  $n \leq N_z$  and every  $y \in U_z$ .

**Lemma 33** The sequence of functionals  $\chi_{\mathcal{A}_x^T}$  converge continuously to the constant-zero functional on  $\mathcal{X}$  mapping any  $x \in \mathcal{X}$  to 0.

**Proof** If  $\mathbf{x} \in \mathcal{A}_x^T$  then  $x_t = 0$  for every  $t \geq T+1 > T$  and therefore,  $\mathbf{x} \in \mathcal{A}_x^{T+1}$ ; thus,  $\mathcal{A}_x^T \subseteq \mathcal{A}_x^{T+1}$ . Moreover, since the algorithm  $\mathbf{x}^{T:T+1}$  defined by  $\mathbf{x}_{T+1}^{T:T+1} = \mathbf{x}$  is in  $\mathcal{A}_x^{T+1} - \mathcal{A}_x^T$  then  $\{\mathcal{A}_x^T\}_{T \in \mathbb{N}}$  is a strictly nested increasing sequence of closed linear subsets of  $\mathcal{A}_x^{\infty:0}$ . Therefore,  $\{\chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$  converges point-wise and in a strictly monotonically decreasing fashion to the functional  $\chi_{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}$ . Hence, by (Dal Maso, 1993, Proposition 5.7) the lower semi-continuous relaxation  $\chi_{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}^{lsc}$  of  $\chi_{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}$  is the  $\Gamma$ -limit of the sequence of functionals  $\{\chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$ .

Similarly,  $\{-\chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$  is strictly monotonically increasing with limit point-wise limit the functional  $-\chi_{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}$ . Hence, by (Dal Maso, 1993, Proposition 5.4) the sequence  $\{-\chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$   $\Gamma$ -converges to  $-\chi_{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}^{lsc}$ . Thus, by definition the sequence  $\{-\chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$  converges continuously to  $\chi_{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}^{lsc}$ .



It remains to compute  $\chi_{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}^{lsc}$ . By (Dal Maso, 1993, Example 3.4)

$$\chi_{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}^{lsc} = \chi_{\overline{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}}, \quad (37)$$

where  $\overline{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}$  denotes the closure of  $\cup_{T \in \mathbb{N}} \mathcal{A}_x^T$  in  $\mathcal{A}_x^{\infty:0}$ . We show that this closure is all of  $\mathcal{A}_x^{\infty:0}$ , in other words, we show that  $\cup_{T \in \mathbb{N}} \mathcal{A}_x^T$  is dense in  $\mathcal{A}_x^{\infty}$ . First, observe that  $\mathbf{y} \in \cup_{T \in \mathbb{N}} \mathcal{A}_x^T$  only if there exists some  $T \in \mathbb{N}$  for which  $\Delta y_t = 0$  for all  $t \geq T$ . Indeed, if  $\mathbf{x} \in \mathcal{A}_x^{\infty:0}$  then, by definition, given any  $\epsilon > 0$  there exists some  $T_\epsilon \in \mathbb{N}$  for satisfying

$$\sum_{t \geq T_\epsilon} \|\Delta x_t\|^p < \epsilon. \quad (38)$$

Observe that, the sequence  $\mathbf{x}^\epsilon$  defined by

$$\mathbf{x}^\epsilon := \begin{cases} x_t & : t \leq T_\epsilon \\ 0 & : t > T_\epsilon \end{cases}$$

belongs to  $\cup_{T \in \mathbb{N}} \mathcal{A}_x^T$ . In particular, the tail estimate (38) implies that

$$\sum_{t \geq T_\epsilon} \|\Delta x_t^\epsilon - \Delta x_t\|^p = \sum_{t \geq T_\epsilon} \|\Delta x_t\|^p < \epsilon.$$

Thus, any  $\mathbf{x} \in \mathcal{A}_x^{\infty:0}$  is an accumulation point of some sequence in  $\cup_{T \in \mathbb{N}} \mathcal{A}_x^{\infty:0}$ , for the strong topology. Since  $\cup_{T \in \mathbb{N}} \mathcal{A}_x^T$  is the union of a nested sequence of convex sets it is itself convex and since the closure of a convex set is itself convex then  $\overline{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}$  is convex. Note that, (Conway, 1990, 1.5 Corollary) implies that the weak closure of  $\overline{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}$  equals to  $\mathcal{A}_x^{\infty:0}$ . Therefore, our remaining computations hold both in the weak and strong topologies on  $\mathcal{A}_x^{\infty:0}$  (and consequentially also on the subsets  $\mathcal{A}_x^{\infty:\alpha}$ ).

Therefore, (37) simplifies to  $\chi_{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}^{lsc} = \chi_{\overline{\cup_{T \in \mathbb{N}} \mathcal{A}_x^T}} = \chi_{\mathcal{A}_x^{\infty:0}} = 0$ .  $\blacksquare$

We are now in place to take the first main step to proving Theorems 2 and 4. Namely, we are in place to conclude that  $\mathcal{R}_\infty$  is the  $\Gamma$ -limit of the sequence of functionals  $\{\mathcal{R}_T + \chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$ .

**Lemma 34** *The  $\mathcal{R}_\infty$  is the  $\Gamma$ -limit of  $\{\mathcal{R}_T + \chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$  on  $\mathcal{A}_x^{\infty:\alpha}$  and  $\mathcal{R}_\infty$  is also the  $\Gamma$ -limit of  $\{\mathcal{R}_T + \chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$  on  $\mathcal{A}_x^{\infty:0}$ .*

**Proof** By Lemma 30,  $\{F_T + \Phi_T\}_{T \in \mathbb{N}}$   $\Gamma$ -converges to  $\mathcal{R}_\infty$  on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. on  $\mathcal{A}_x^{\infty:0}$  for the weak topology). By Lemma 33,  $\{\chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$  converges to the constant 0 functional continuously. Moreover, since 0 is everywhere finite on  $\mathcal{A}_x^{\infty:0}$  then (Dal Maso, 1993, Proposition 6.20) implies that

$$\Gamma - \lim_{T \uparrow \infty} \mathcal{R}_T + \chi_{\mathcal{A}_x^T} = \sup_{T \in \mathbb{N}} \mathcal{R}_T^{lsc} + 0 = \sup_{T \in \mathbb{N}} \mathcal{R}_T^{lsc}; \quad (39)$$

(where we take the lower semi-continuous relaxation with respect to the strong topology on  $\mathcal{A}_x^{\infty:\alpha}$ , otherwise, mutatis mutandis, we take it with respect to the weak topology on  $\mathcal{A}_x^{\infty:0}$ ).

By Lemmas 28, 27, and (Dal Maso, 1993, Proposition 1.9) we find that each  $\mathcal{R}_T$  is lower semi-continuous. Thus, the right-hand side of (39) yields equal to  $\sup_{T \in \mathbb{N}} \mathcal{R}_T^{lsc}$ . Hence, the conclusion follows since  $\Gamma - \lim_{T \uparrow \infty} \mathcal{R}_T + \chi_{\mathcal{A}_x^T} = \sup_{T \in \mathbb{N}} \mathcal{R}_T^{lsc} = \mathcal{R}_\infty$ .  $\blacksquare$

**Lemma 35 (Equi-coercivity Lemma)** *The family of functionals  $\{\mathcal{R}_T + \chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$  is equi-coercive on  $\mathcal{A}_x^{\infty:\alpha}$ . In addition,  $\phi$  the family of functionals  $\{\mathcal{R}_T + \chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$  is equi-coercive on  $\mathcal{A}_x^{\infty:0}$  for the weak topology.*

**Proof** Fix  $T \in \mathbb{N}$  and  $\mathbf{x} \in \mathcal{A}_x^{\infty:0}$ . If  $\mathbf{x} \in \mathcal{A}_x^T$  then by definition  $\chi_{\mathcal{A}_x^T}(\mathbf{x}) = 0$  and  $\Delta x_t = 0$  for every  $t \geq T$ . Since Assumption 2 guarantees that  $\phi(0) = 0$ , then  $\phi(\Delta x_t) = 0$  for every  $t \geq T$ . Therefore, we may compute:

$$\begin{aligned}
 (\mathcal{R}_T + \chi_{\mathcal{A}_x^T})(\mathbf{x}) &= \sum_{t=1}^T f(x_t) - f^* + \phi_t(\Delta x_t) + \chi_{\mathcal{A}_x^T}(\mathbf{x}) \\
 &= \sum_{t=1}^T f(x_t) - f^* + \sum_{t=1}^T \phi_t(\Delta x_t) \\
 &= \sum_{t=1}^T f(x_t) - f^* + \sum_{t=1}^T \phi_t(\Delta x_t) + 0 \\
 &= \sum_{t=1}^T f(x_t) - f^* + \sum_{t=1}^{\infty} \phi_t(\Delta x_t) \\
 &\geq \sum_{t=1}^{\infty} \phi_t(\Delta x_t);
 \end{aligned} \tag{40}$$

where the last inequality in (40) holds since  $-\infty < f^* \leq f(x)$  by Assumption 1.

If  $\mathbf{x} \notin \mathcal{A}_x^T$ , then  $\chi_{\mathcal{A}_x^T}(\mathbf{x}) = \infty$ . In which case we necessarily have

$$(\mathcal{R}_T + \chi_{\mathcal{A}_x^T})(\mathbf{x}) = \infty \geq \sum_{t=1}^{\infty} \phi_t(\Delta x_t). \tag{41}$$

Thus, together (40) and (41) imply that for every  $T \in \mathbb{N}$  and every  $\mathbf{x} \in \mathcal{A}_x^{\infty:0}$  the following bound must hold:

$$(\mathcal{R}_T + \chi_{\mathcal{A}_x^T})(\mathbf{x}) \geq \sum_{t=1}^{\infty} \phi_t(\Delta x_t) = \Phi_{\infty}(\mathbf{x}). \tag{42}$$

Therefore, by Lemma 27 and (Dal Maso, 1993, Proposition 7.7), we may conclude that  $\{\mathcal{R}_T + \chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$  forms an equi-coercive family on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. on  $\mathcal{A}_x^{\infty:0}$  for the weak topology if  $\phi$  is also convex) of functionals since  $\Phi_{\infty}$  is itself coercive on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. on  $\mathcal{A}_x^{\infty:0}$  for the weak topology).  $\blacksquare$

**Lemma 36** *Suppose that  $\mathcal{P}_x^{\infty:0} \neq \emptyset$ . Then  $\mathcal{P}_x^{\infty} = \mathcal{P}_x^{\infty:0}$ , and hence  $\mathcal{P}_x^{\infty}$  is also non-empty.*

**Proof** The proof follows from the fact that  $\mathcal{R}_{\infty}$  is bounded over both  $\mathcal{P}_x^{\infty}$  and  $\mathcal{P}_x^{\infty}$ .

Let us separate  $\mathcal{A}_x^{\infty} = \mathcal{P}_x^{\infty:0} \cup C_x^{\infty} \cup D_x^{\infty}$  into the three disjoint parts, where we define  $C_x^{\infty} = \mathcal{A}_x^{\infty:0} \setminus \mathcal{P}_x^{\infty:0}$  and  $D_x^{\infty} = \mathcal{A}_x^{\infty} \setminus \mathcal{A}_x^{\infty:0}$ . We show that the claim of the lemma holds by demonstrating the equivalent claim that for any  $\mathbf{x} \in \mathcal{P}_x^{\infty:0}$  and  $\mathbf{y} \in C_x^{\infty} \cup D_x^{\infty}$ , we have that  $\mathcal{R}_{\infty}(\mathbf{x}) < \mathcal{R}_{\infty}(\mathbf{y})$ .

Note that the above claim is equivalent to the claim of the lemma since it implies that for any  $\mathbf{y} \in \mathcal{A}_x^\infty$ , we have  $\mathcal{R}_\infty(\mathbf{x}) \leq \mathcal{R}_\infty(\mathbf{y})$  if and only if  $\mathbf{x} \in \mathcal{P}_x^{\infty:0}$ .

First, we note that by the definition of  $\mathcal{P}_x^{\infty:0}$ , since  $C_x^\infty \subseteq \mathcal{A}_x^{\infty,0}$ , we have that for any  $\mathbf{x} \in \mathcal{P}_x^\infty$  and  $\mathbf{y} \in C_x^\infty$ , we have that  $\mathcal{R}_\infty(\mathbf{x}) < \mathcal{R}_\infty(\mathbf{y})$ .

Next, we begin by recalling that if  $\mathbf{x} \in \mathcal{P}_x^{\infty:0}$ , then  $\mathcal{R}(\mathbf{x}) < \infty$ . This holds since by defining  $\mathbf{y} \in \mathcal{A}_x^{\infty:0}$  such that  $y_0 = x$  and  $y_u = x^*$  for all  $u > 0$ , we have that  $\mathcal{R}_\infty(\mathbf{y}) = \phi(x - x^*) < \infty$ , which by the definition of  $\mathbf{x} \in \mathcal{P}_x^{\infty:0}$  implies that  $\infty > \mathcal{R}_\infty(\mathbf{y}) \geq \mathcal{R}_\infty(\mathbf{x}) \geq 0$ . By the definition of  $\mathcal{A}_x^{\infty:0}$ , we have that for any  $\mathbf{y} \in D_x^\infty$ , the sum  $\sum_{u=0}^\infty \|\Delta y_u\|^p = \infty$  diverges. Moreover, by Assumption 2, we have that there exists a constant  $c > 0$  such that  $\phi(z) \geq c\|z\|^p$ . Hence,

$$\mathbf{y} \in D_x^\infty \implies \mathcal{R}_\infty(\mathbf{y}) \geq \sum_{u=0}^\infty \phi(\Delta y_u) \geq c \sum_{u=0}^\infty \|\Delta y_u\|^p = \infty.$$

We have shown that  $\mathcal{R}_\infty(\mathbf{x}) < \infty$  for all  $\mathbf{x} \in \mathcal{P}_x^{\infty:0}$  and  $\mathcal{R}_\infty(\mathbf{y}) = \infty$  for all  $\mathbf{y} \in D_x^\infty$ . Combining these facts, we obtain that  $\mathcal{R}_\infty(\mathbf{x}) < \mathcal{R}_\infty(\mathbf{y})$ , concluding the proof.  $\blacksquare$

### C.3.2. PROOF OF THEOREM 2

We are now in place to prove Theorem 2. Since the proof of the  $\mathcal{A}_x^{\infty,0}$  case and the  $\mathcal{A}_x^{\infty,\alpha}$  (for  $\alpha > 0$ ) case are analogous, we simplify our exposition by combining them and highlight their differences when necessary.

**Proof** By Lemma 34 and 35 the family of functionals  $\left\{ \mathcal{R}_T + \chi_{\mathcal{A}_x^T} \right\}_{T \in \mathbb{N}}$  is lower semi-continuous and equi-coercive on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. weakly lower semi-continuous and equi-coercive with respect to its weak topology on  $\mathcal{A}_x^{\infty:0}$ ). Therefore, (Dal Maso, 1993, Theorem 7.8) implies that  $\mathcal{R}_\infty$  is coercive on  $\mathcal{A}_x^{\infty:\alpha}$  (on  $\mathcal{A}_x^{\infty:0}$  with respect to the weak topology).

Lemma 29 guaranteed that, for every  $T \in \mathbb{N}$ , the regret functionals  $\mathcal{R}_T$  are all coercive on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. on  $\mathcal{A}_x^{\infty:0}$  with respect to the weak topology). Since  $\mathcal{R}_\infty = \sup_{T \in \mathbb{N}} \mathcal{R}_T$  then (Dal Maso, 1993, Proposition 1.8) guarantees that  $\mathcal{R}_\infty$  is lower semi-continuous on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. on  $\mathcal{A}_x^{\infty:0}$  with respect to the weak topology).

Next, Lemma 29 guaranteed that, for every  $T \in \mathbb{N}$ , the regret functional  $\mathcal{R}_T$  takes non-negative values. Hence, for every  $\mathbf{x} \in \mathcal{A}_x^\infty$  we compute

$$\mathcal{R}_\infty(\mathbf{x}) = \sup_{T \in \mathbb{N}} \mathcal{R}_T(\mathbf{x}) \geq 0. \quad (43)$$

Therefore,  $\mathcal{R}_\infty$  is both lower semi-continuous and coercive on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. weakly lower semi-continuous and coercive with respect to the weak topology on  $\mathcal{A}_x^{\infty:0}$ ) and it is bounded below by 0. Hence, (Dal Maso, 1993, Theorem 1.15) implies that  $\mathcal{P}_x^{\infty:\alpha} \neq \emptyset$  ( $\mathcal{P}_x^\infty \neq \emptyset$ ).

Now, Assumptions 1 and (2) imply that the sequence  $\mathbf{x}^\infty$  defined by  $x_0^\infty = x$  and  $x_t^\infty = x^*$  for  $t \geq 1$  satisfies  $\mathbf{x} \in \mathcal{A}_x^{\infty:\alpha}$  (resp. in  $\mathcal{A}_x^{\infty:0}$ ) and  $\mathcal{R}_\infty(\mathbf{x}) < \infty$ . Hence, by definition, any  $\mathbf{x}^* \in \mathcal{P}_x^{\infty:\alpha}$  (resp. in  $\mathcal{P}_x^{\infty:0}$ ) must satisfy  $\mathcal{R}_\infty(\mathbf{x}^*) \leq \mathcal{R}_\infty(\mathbf{x}^\infty) < \infty$ . Therefore,  $\mathcal{P}_x^{\infty:\alpha}$  (resp.  $\mathcal{P}_x^{\infty:0}$ ) is non-empty and any algorithm therein has finite regret.

Lastly, we apply Lemma 36 to conclude that  $\mathbf{x} \in \mathcal{P}^{\infty:0} \implies \mathbf{x} \in \mathcal{P}^\infty$  and hence that  $\mathcal{P}^\infty$  is non-empty.  $\blacksquare$

**C.4. Proof of Corollary 3**

**Proof** Since  $x \in \mathcal{P}$ , Theorem (2) implies that

$$\mathcal{R}_\infty(x) = \sum_{t=0}^{\infty} f(x_t) - f^* + \phi(\Delta x_t) = \sum_{t=0}^{\infty} a_t < \infty,$$

where the summand  $a_t = f(x_t) - f^* + \phi(\Delta x_t)$  is non-negative. Hence, we have that  $\lim_t a_t = 0$ .

Now, assume that the sequence is monotone. Since the  $a_t$  are summable, the partial sums  $S_t = \sum_{u=t}^{\infty} a_u \rightarrow 0$  asymptotically vanish. Hence, we have the bound

$$2S_t \geq 2 \sum_{u=t}^{2t} a_u \geq 2ta_{2t} \geq 0.$$

We therefore have that  $\lim_{t \rightarrow \infty} ta_t = 0$ , proving the claim of the theorem. ■

**C.5. Proof of Theorem 4**

The Lemmas established in C.3.1 reduce the proof of Theorem 4 to a simple consequence of the Fundamental Theorem of  $\Gamma$ -convergence. As with the proof of Theorem 2, since the proof of the convex and the non-convex cases are analogous, mutatis mutandis, and are therefore combined.

**Proof** Since  $\mathcal{R}_\infty$  is the  $\Gamma$ -limit of the equi-coercive sequence of functionals  $\{\mathcal{R}_T + \chi_{\mathcal{A}_x^T}\}_{T \in \mathbb{N}}$  on  $\mathcal{A}_x^{\infty:\alpha}$  (resp. on  $\mathcal{A}_x^{\infty:0}$  with respect to the weak topology), then result therefore follows directly from (Braides, 2014, Corollary 2.1). ■

## Appendix D. Proofs for Section 3

### D.1. Proof of Theorem 5

**Proof** We begin by recalling the definition of the Gâteaux derivative,

$$\mathcal{R}'_T(\mathbf{x})(\delta\mathbf{x}) = \lim_{\epsilon \rightarrow 0} \frac{\mathcal{R}_T(\mathbf{x} + \epsilon\delta\mathbf{x}) - \mathcal{R}_T(\mathbf{x})}{\epsilon},$$

where  $\delta\mathbf{x} = \mathbf{y} - \mathbf{x}$  for some  $\mathbf{y} \in \mathcal{A}_x^T$ . Expanding the definition of  $\mathcal{R}_T$ , exchanging the limit with the sum and applying the assumed smoothness of Assumption 3, we obtain that

$$\mathcal{R}'_T(\mathbf{x})(\delta\mathbf{x}) = \sum_{t=1}^T \langle \nabla f(x_t), \delta x_t \rangle + \langle \nabla \phi(\Delta x_{t-1}), \Delta \delta x_{t-1} \rangle.$$

Lastly, noting that  $\delta x_0 = 0$  since  $x_0 = y_0 = x$  and re-arranging the sum, we obtain the expression in the statement of the theorem.  $\blacksquare$

### D.2. Proof of Theorem 6

**Proof** First we show that  $\hat{\mathcal{P}}_x^T$  is precisely the set of  $\mathbf{x} \in \mathcal{A}_x^T$  which make  $\mathcal{R}'_T$  vanish. To see this, note that  $\mathcal{R}'(\mathbf{x})(\delta\mathbf{x})$  given in equation (4) is a bounded linear functional in  $\delta\mathbf{x}$ , and hence vanishes if and only if we have

$$\nabla f(x_{t+1}) - \nabla \phi(x_t) - \nabla \phi(x_{t+1}) = 0$$

for all  $t = 0, 1, \dots, T-1$ . Hence, we have that

$$\mathcal{P}_x^T = \{\mathbf{x} \in \mathcal{A}_x^T : \mathcal{R}'_T(\mathbf{x}) \equiv 0\},$$

and is by definition the set of critical points. Since  $\mathcal{A}_x^T$  is open and since  $\mathcal{P}_x^T \subseteq \mathcal{A}_x^T$  is non-empty, we must have that  $\mathcal{R}'_T(\mathbf{x}) = 0$  for any  $\mathbf{x} \in \mathcal{P}_x^T$ . Hence, we obtain the inclusion  $\hat{\mathcal{P}}_x^T \supseteq \mathcal{P}_x^T$ . Lastly, it is easy to see that the recursion-ii holds by the stationarity of equation (6).  $\blacksquare$

### D.3. Proof of Theorem 7

**Proof** We begin by the case where  $T < \infty$ . Recall that by the dynamic programming principle (Lemma 19) that

$$J^{T-t}(x_t) = \min_{y \in \mathcal{X}} \{\phi(y - x_t) + f(y) + J^{T-(t+1)}(y)\},$$

where  $x_{t+1} \in C_x^{T-t} = \operatorname{argmin}_{y \in \mathcal{X}} \{\phi(y - x_t) + f(y) + J^{T-(t+1)}(y)\}$ . Since  $x_{t+1} \in C_x^{T-t}$ , by the differentiability of  $f$  and  $\phi$ , as well as the assumed local-Lipschitz property of  $J^{T-(t+1)}$ , we find that

$$0 \in \nabla \phi(x_{t+1} - x_t) + \nabla f(x_{t+1}) + \partial J^{T-(t+1)}(x_{t+1}),$$

and hence that

$$-\nabla \phi(x_{t+1} - x_t) - \nabla f(x_{t+1}) \in \partial J^{T-(t+1)}(x_{t+1}). \quad (44)$$

Since  $\mathbf{x} \in \hat{\mathcal{P}}^T$ , we have that the optimal dynamics of equation (5) must hold and hence we get that  $\nabla\phi(x_{t+1} - x_t) + \nabla f(x_{t+1}) = \nabla\phi(x_{t+2} - x_{t+1})$ , yielding the claim of the theorem.

In the case of  $T = \infty$ , we can arrive at equation (44) with the DPP (Lemma 20) and applying the same sequence of steps, which yield that for all  $t \in \mathbb{N}$ ,

$$-\nabla\phi(\Delta x_t) - \nabla f(x_{t+1}) \in \partial J^\infty(x_{t+1}), \quad (45)$$

where the assumed local Lipschitzness of  $J^\infty$  ensures that  $\partial J^\infty$  is always non-empty.

In order to show that the recursion, applying the DPP of Lemma 20 twice implies that for all  $t$ ,

$$J^\infty(x_t) = \phi(\Delta x_t) + f(x_{t+1}) + \phi(\Delta x_{t+1}) + f(x_{t+2}) + J^\infty(x_{t+1}).$$

and that  $(x_{t+1}, x_{t+2}) \in \operatorname{argmin}_{y, z \in \mathcal{X}} \{\phi(y - x_t) + \phi(z - y) + f(y) + f(z) + J^\infty(z)\}$ . By the local Lipschitz smoothness of the above function, we know that its generalized derivative must contain zero at  $(x_{t+1}, x_{t+2})$ . Taking the (generalized) derivative at  $y = x_t$  and letting it vanish we obtain that

$$\nabla\phi(\Delta x_t) = \nabla\phi(\Delta x_{t+1}) + \nabla f(x_{t+1}), \quad (46)$$

yielding the desired recursion.

Hence, combining with equation (45), we find that we must have that

$$-\nabla\phi(\Delta x_{t+1}) \in \partial J^\infty(x_{t+1})$$

for all  $t$ , as well as the recursion (46). Noting that  $x = x_t$  and  $t \in \mathbb{N}$  can be chosen arbitrarily, we have that  $\nabla J^\infty(x_0) = -\nabla\phi(\Delta x_0)$  for all  $x_0 \in \Upsilon$ , where

$$\Upsilon_1 = \{y \in \mathcal{X} : \exists \mathbf{x} \in \mathcal{P}^\infty \text{ such that } y = x_1\} \subseteq \mathcal{X},$$

the set of points that can be reached in a single step of an algorithm  $\mathbf{x} \in \mathcal{P}^\infty$ .

We now show that  $\Upsilon_1 = \mathcal{X}$ . Using equation (45) and the property that  $\nabla\phi \circ \nabla\phi^* = \operatorname{id}$  we have that

$$x_0 = x_1 - \nabla\phi^*(-\nabla f(x_1) - \nu(x_1)), \quad (47)$$

for some  $\nu(x_1) \in \partial J^\infty(x_1)$ . Hence, for any  $x_1 = y \in \mathcal{X}$ , we can pick  $\mathbf{x} \in \mathcal{P}_x^\infty$  where  $x = x_0$  defined according to (47), which shows that  $\Upsilon_1 \supseteq \mathcal{X}$ , as desired. We therefore have that  $\nabla\phi(\Delta x_t) = -\nabla J^{T-1}(x_t)$  for all  $\mathbf{x} \in \mathcal{P}^T$ . Lastly, it is easy to see that the recursion follows from (46).  $\blacksquare$

## Appendix E. Proofs for Section 4

Over the course of this section, we assume without loss of generality that  $f_\star = 0$  since we may simply consider the function  $\tilde{f}(x) = f(x) - f_\star$ , which satisfies this property.

### E.1. Proof of Lemma 8

**Proof** Consider  $\mathbf{x}, \mathbf{y} \in \mathcal{A}_x^T$  and  $\rho \in (0, 1)$ . Then we have

$$\mathcal{R}_T(\mathbf{x} + \rho(\mathbf{y} - \mathbf{x})) = \sum_{t=1}^T f(x_t + \rho(y_t - x_t)) + \phi(\Delta x_t + \rho(\Delta x_t - \Delta y_t)) .$$

Noting that by the convexity of  $f$  and the strict convexity of  $\phi$ , we have

$$\begin{aligned} f(x_t + \rho(y_t - x_t)) &\leq (1 - \rho)f(x_t) + \rho f(y_t) , \\ \phi(\Delta x_t + \rho(\Delta x_t - \Delta y_t)) &< (1 - \rho)\phi(\Delta x_t) + \rho\phi(\Delta y_t) , \end{aligned}$$

and hence, we find that

$$\mathcal{R}_T(\mathbf{x} + \rho(\mathbf{y} - \mathbf{x})) < (1 - \rho)\mathcal{R}_T(\mathbf{x}) + \rho\mathcal{R}_T(\mathbf{y}) ,$$

showing that  $\mathcal{R}_T$  is strictly convex, and hence has a unique minimum. Applying (Ekeland and Temam, 1999, Proposition 1.2), the solution must be unique and by (Ekeland and Temam, 1999, Proposition 2.1) is the unique critical point of  $\mathcal{R}_T$ .  $\blacksquare$

### E.2. Proof of Lemma 9

**Proof** Let us first assume that  $T \in \mathbb{N} \cup \{\infty\}$ . In order to show that  $J^T$  is convex and differentiable we leverage convex analysis tools from Rockafellar (1970). Let us introduce the (abuse of) notation

$$\mathcal{R}_T(x; \mathbf{x}_{1:T}) = \mathcal{R}_T(\mathbf{x}) ,$$

for  $\mathbf{x} \in \mathcal{A}_x^T$  such that  $x_0 = x$  and  $\mathbf{x}_{1:T} = \{x_t\}_{t=1}^T$ , which allows us to separate the initial value and the remainder of the path of the optimizer. Note that by the convexity of  $f$  and  $\phi$  that  $\mathcal{R}_T$  is convex in both variables and that by definition we have  $J^T(x) = \min_{\mathbf{y}_{1:T} \in \mathcal{X}^{\otimes T}} \mathcal{R}_T(x; \mathbf{y}_{1:T})$ . Now for  $x, y \in \mathcal{X}$ ,  $\mathbf{x}_{1:T}, \mathbf{y}_{1:T} \in \mathcal{X}^{\otimes T}$ ,  $\rho \in (0, 1)$

$$\begin{aligned} J^T((1 - \rho)x + \rho y) &= \min_{\mathbf{z}_{1:T} \in \mathcal{X}^{\otimes T}} \mathcal{R}_T((1 - \rho)x + \rho y; \mathbf{z}_{1:T}) \\ &\leq \mathcal{R}_T((1 - \rho)x + \rho y; (1 - \rho)\mathbf{x}_{1:T} + \rho\mathbf{y}_{1:T}) \\ &\leq (1 - \rho) \min_{\mathbf{x}_{1:T} \in \mathcal{X}^{\otimes T}} \mathcal{R}_T(x; \mathbf{x}_{1:T}) + \rho \min_{\mathbf{y}_{1:T} \in \mathcal{X}^{\otimes T}} \mathcal{R}_T(y; \mathbf{y}_{1:T}) \\ &= (1 - \rho) J^T(x) + \rho J^T(y) . \end{aligned}$$

Taking the minimum over  $\mathbf{x}_{1:t}$  and  $\mathbf{y}_{1:t}$ , we obtain

$$\begin{aligned} J^T((1 - \rho)x + \rho y) &\leq (1 - \rho) \min_{\mathbf{x}_{1:T} \in \mathcal{X}^{\otimes T}} \mathcal{R}_T(x; \mathbf{x}_{1:T}) + \rho \min_{\mathbf{y}_{1:T} \in \mathcal{X}^{\otimes T}} \mathcal{R}_T(y; \mathbf{y}_{1:T}) \\ &= (1 - \rho) J^T(x) + \rho J^T(y) , \end{aligned}$$

demonstrating the claim that  $J^T$  is convex for all  $T \in \mathbb{N} \cup \{\infty\}$ .

Next, we show that  $J^T$  is differentiable for  $T \in \mathbb{N}$ . Note that for each fixed  $\mathbf{y}_{1:T}$ , that the function  $x \mapsto \mathcal{R}_T(x; \mathbf{y}_{1:T})$  is both convex and differentiable, where the differentiability of  $x$  follows from the differentiability of  $\phi$ . Fix  $x \in \mathcal{X}$ , and define a sequence  $\{\mathbf{y}_{1:T}^i\}_{i \in \mathbb{N}} \subseteq \mathcal{X}^{\otimes T}$  such that

$$f_i(x) = \mathcal{R}_T(x; \mathbf{y}_{1:T}^i) \longrightarrow \min_{\mathbf{y}_{1:T} \in \mathcal{X}^{\otimes T}} \mathcal{R}_T(x; \mathbf{y}_{1:T}) = J^T(x) .$$

since each  $f_i$  is convex and differentiable over  $\mathcal{X}$ , we can apply (Rockafellar, 1970, Theorem 25.7) to claim that  $\nabla f_i(x) \rightarrow \nabla J^T(x)$ , and hence  $\nabla J^T(x)$  is differentiable.

Next, we show that both of the convergence statements of Lemma 9 hold, which in turn imply the differentiability for  $T = \infty$ . Once more, we will leverage the results of (Rockafellar, 1970, Theorem 25.7). Notice that for each  $T \in \mathbb{N}$ ,  $J^T$  is convex and differentiable. Moreover, note that  $J^T$  is pointwise non-decreasing and bounded above due to Theorem 2, and hence  $J^T \rightarrow J^\infty$  pointwise. Applying (Rockafellar, 1970, Theorem 25.7), we get that these properties imply that  $J^T \rightarrow J^\infty$  and  $\nabla J^T \rightarrow \nabla J^\infty$  uniformly on compact sets, which also show that  $J^\infty$  is differentiable. ■

### E.3. Proof of Lemma 10

**Proof** We note that Lemma 10 is a special case of Theorem 7, where for any  $T \in \mathbb{N} \cup \{\infty\}$ ,  $J^T$  is convex and differentiable. In this case, we find that the necessary conditions of Theorem 7 are satisfied. Moreover, we have that since  $J^T$  is differentiable,  $\partial J^T(x) = \{\nabla J^T(x)\}$ . Applying this to the result of Theorem 7, we obtain the desired result. ■

### E.4. Proof of Lemma 11

**Proof** We split the proof according to the individual properties listed in the statement of the Lemma.

**Proof of Property i.** We recall the result from (Rockafellar, 1970, Theorem 26.5) which states that a function is Legendre convex if and only if its dual is Legendre convex. Hence, it is sufficient for us to show that  $(J^T)^*$  is Legendre convex. What remains to be shown are that  $(J^T)^*$  is convex, differentiable and satisfies the property that  $\lim_{\|x\| \rightarrow \infty} \|\nabla J^T(x)\| = \infty$ .

We first show that  $(J^T)^*$  is strictly convex. Recall the recursion on  $(J^T)^*$  from Lemma 21,

$$(J^T)^*(q) = \tilde{\phi}^*(q) + (J^{T-1} + f)^*(q) . \quad (48)$$

Since  $\phi$  is Legendre convex,  $\tilde{\phi}^*$  must also be Legendre convex and hence strictly convex. Since  $J^{T-1}$  is convex (property i) and  $f$  is strictly convex, we therefore have that  $(J^{T-1} + f)^*$  is convex. Since  $(J^T)^* = \tilde{\phi}^* + (J^{T-1} + f)^*$  is the sum of a strictly convex and a convex function, it is strictly convex and hence  $(J^T)^*$  is strictly convex.

Next, we show that  $(J^T)^*$  is differentiable. First note that  $\tilde{\phi}^*$  is differentiable since it is Legendre convex. Next, recall that  $J^{T-1} + f$  is strictly convex, and hence by (Rockafellar, 1970, Theorem 26.3)  $(J^{T-1} + f)^*$  is differentiable. Hence, by (48) we have that  $(J^T)^*$  is strictly convex.

Now, note that by (Rockafellar, 1970, Lemma 26.7), a convex function  $g : \mathcal{X} \rightarrow \mathbb{R}$  satisfies  $\lim_{\|x\| \rightarrow \infty} \|\nabla g(x)\| = \infty$  if and only if  $g$  is co-finite, that is,  $g$  satisfies

$$\lim_{\lambda \rightarrow \infty} g(\lambda y) / \lambda = \infty \quad \forall 0 \neq y \in \mathcal{X} .$$



By Fenchel's inequality, we have that  $(J^{T-1} + f)^*(q) \geq -J^{T-1}(0) + f(0) = -\alpha > -\infty$  and applying Lemma 21 with  $p = 0$ , we obtain the bound

$$(J^T)^*(q) = \phi^*(q) + (J^{T-1} + f)^*(q) \geq \phi^*(q) - \alpha .$$

Since  $\phi$  is assumed to be Legendre convex (Rockafellar, 1970, Theorem 26.5) implies that  $\tilde{\phi}^*$  is also Legendre convex and hence co-finite. Hence, we have that

$$\lim_{\lambda \rightarrow \infty} \frac{(J^T)^*(\lambda y)}{\lambda} \geq \lim_{\lambda \rightarrow \infty} \frac{\tilde{\phi}^*(\lambda y) - \alpha}{\lambda} = \infty ,$$

showing that  $(J^T)^*$  is also co-finite and hence Legendre convex, as desired.

**Proof of Property ii.** Let  $x^* \in \operatorname{argmin}_{x \in \mathcal{X}} f(x)$ , and consider  $\mathbf{x} \in \cup_{T \in \mathbb{N}} \mathcal{A}_{x^*}^T$  defined by  $x_t = x^*$  for all  $t \in \mathbb{N}$ . Note that under this definition,  $0 = \mathcal{R}(\mathbf{x}) \geq \min_{\mathbf{y} \in \mathcal{A}_{x^*}^T} \mathcal{R}_T(\mathbf{y}) = J^T(x^*) \geq \min_x J^T(x) \geq 0$ . By property *i*, we have that  $J^T$  is Legendre convex and hence strictly convex, so we have that this minimum is unique.

**Proof of Property iii.** Note that by property *i*,  $(J^T)^*$  is Legendre convex. Hence the relative convex with respect to  $\tilde{\phi}^*$  follows directly from Lemma 27, since we have

$$\begin{aligned} D_{(J^T)^*}(q, p) &= D_{\tilde{\phi}^*}(q, p) + D_{(J^{T-1}+f)^*}(q, p) \\ &\geq D_{\tilde{\phi}^*}(q, p) \end{aligned}$$

where the inequality follows from the positivity of the Bregman divergence. Hence, we obtain one of the claims for *Property i* of the theorem. For the second claim, we apply the result of Lemma 22 to get the desired result. ■

## E.5. Proof of Lemma 12

**Proof** We first note here that the assumption that  $\phi$  is quadratic implies that  $\tilde{\phi}(x) = \phi(-x) = \phi(x)$ . Now, assume that  $f$  is  $\lambda$ -relatively-smooth with respect to  $\phi$ . Let us define the set

$$\Gamma = \{ \gamma \in [0, 1] : D_{J^\infty}(x, y) \leq \gamma D_{\phi^*}(x, y) \quad \forall x, y \in \mathcal{X} \} ,$$

as well as its infimum  $\underline{\gamma} = \inf \Gamma$ . Note that by Theorem 11-iii,  $\Gamma$  is non-empty and hence  $\underline{\gamma}$  is well-defined. Furthermore, it is easy to see that  $\Gamma$  is closed and bounded and hence compact, therefore  $\underline{\gamma} \in \Gamma$ .

Now since  $f$  is  $\lambda$ -relatively-smooth, by the linearity of the Bregman divergence we obtain that  $D_{J^\infty+f}(x, y) \leq (\underline{\gamma} + \lambda) D_{\phi^*}(x, y)$ , and hence, applying Lemma 22, we have that

$$D_{(J^\infty+f)^*}(p, q) \geq (\underline{\gamma} + \lambda)^{-1} D_{\phi^*}(p, q) . \quad (49)$$

Recalling Lemma 21, and taking the limit as  $T \rightarrow \infty$ , we have

$$D_{(J^\infty)^*}(q, p) = D_{\phi^*}(q, p) + D_{(J^\infty+f)^*}(q, p) ,$$

hence combining with (49), we obtain

$$D_{(J^\infty)^*}(q, p) \geq (1 + (\underline{\gamma} + \lambda)^{-1})D_{\phi^*}(q, p) . \quad (50)$$

Applying Lemma 22 once more to (50), we find that

$$D_{J^\infty}(x, y) \leq (1 + (\underline{\gamma} + \lambda)^{-1})^{-1}D_{\phi^*}(x, y) ,$$

and hence we have that  $(1 + (\underline{\gamma} + \lambda)^{-1})^{-1} \in \Gamma$ . By the definition of  $\underline{\gamma}$ , however, we have that

$$\underline{\gamma} \leq (1 + (\underline{\gamma} + \lambda)^{-1})^{-1} . \quad (51)$$

Noting that equation (51) can be re-arranged into a quadratic inequality in terms of  $\underline{\gamma}$  and that  $\underline{\gamma} \in [0, 1]$ , we can solve the inequality to obtain that

$$\underline{\gamma} \leq \frac{1}{2} \left( \sqrt{\lambda^2 + 4\lambda} - \lambda \right) \in (0, 1) ,$$

as desired. In order to obtain the converse result, we begin by assuming that  $f$  is  $\mu$ -relatively-convex, and repeat the same sequence of steps with the inequalities reversed and modifying the definition of the set  $\Gamma$  and of  $\underline{\gamma}$  accordingly (as a sup).  $\blacksquare$

## E.6. Proof of Theorem 13

**Proof** Over the course of this proof, we use the short-hand notation  $\nabla J^\infty(x_t) = \nabla J_t^\infty$ .

We begin by noting that Lemma 24-2 and the strict convexity of  $\phi$  implies that  $\tilde{\phi}^*(\nabla J_{t+1}^\infty) \leq \tilde{\phi}^*(\nabla J_t^\infty)$ , and hence  $\{\tilde{\phi}^*(\nabla J_{t+1}^\infty)\}$  is decreasing. Next, recalling Lemma 24-3, we have that

$$\tilde{\phi}^*(\nabla J_{t+1}^\infty) \leq -D_{\tilde{\phi}^*}(\nabla J^T(x^*), \nabla J^T(x_t)) + D_{J^T}(x_t, x^*) - D_{J^T}(x_{t+1}, x^*) . \quad (52)$$

Now, noting that  $\tilde{\phi}^*(\nabla J_{t+1}^\infty)$  is decreasing, we get that

$$\begin{aligned} t \tilde{\phi}^*(\nabla J_t^\infty) &\leq \sum_{u=1}^t \tilde{\phi}^*(\nabla J_u^\infty) \\ &\leq \sum_{u=0}^{t-1} \left\{ -D_{\tilde{\phi}^*}(\nabla J^\infty(x^*), \nabla J^\infty(x_u)) + D_{J^\infty}(x_u, x^*) - D_{J^\infty}(x_{u+1}, x^*) \right\} \\ &\leq \sum_{u=0}^{t-1} \left\{ -D_{\tilde{\phi}^*}(\nabla J^\infty(x^*), \nabla J^\infty(x_{u-1})) \right\} + D_{J^\infty}(x_0, x^*) - D_{J^\infty}(x_t, x^*) \\ &\leq D_{J^\infty}(x_0, x^*) , \end{aligned}$$

and hence, dividing both sides by  $t$ , we obtain the bound in the statement of the theorem.  $\blacksquare$

### E.7. Proof of Theorem 14

**Proof** Over the course of this proof, we will use the short-hand notation  $J^\infty(x_t) = J_t^\infty$ . We recall once more that if  $\phi$  is quadratic, then  $\nabla\phi$  is linear and hence we have that  $D_\phi(x, y) = \phi(x - y)$ . This proof follows closely the proof of Theorem 13, but where we replace the use of Lemma 24 with Lemma 23. We begin by noting that Lemma 23-2 and the convexity of  $\phi$  imply that  $J_{t+1}^\infty \leq J_t^\infty$ , and hence  $\{J_t^\infty\}$  is non-increasing. Next, recalling Lemma 23-3, we have that

$$D_{J^\infty}(x_{t+1}, x^*) \leq -D_{J^\infty}(x^*, x_t) + D_\phi(x^*, x_t) - D_\phi(x^*, x_{t+1}). \quad (53)$$

Since  $\{J_t^\infty\}$  is non-increasing and since  $D_{J^\infty}(x_t, x^*) = J_t^\infty$  we get that by using (53),

$$\begin{aligned} t J_t^\infty &\leq \sum_{u=1}^t J_u^\infty \\ &\leq \sum_{u=0}^{t-1} \{-D_{J^\infty}(x^*, x_u) + D_\phi(x^*, x_u) - D_\phi(x^*, x_{u+1})\} \\ &\leq \sum_{u=0}^{t-1} \{D_\phi(x^*, x_u) - D_\phi(x^*, x_{u+1})\} \\ &= D_\phi(x^*, x_0) - D_\phi(x^*, x_t) \\ &\leq D_\phi(x^*, x_0) = \phi(x^* - x_0), \end{aligned}$$

hence, dividing both sides by  $t$ , we obtain the first bound in the statement of the theorem.

To obtain the second, we note that if  $J^\infty$  is also  $\mu$ -relatively convex with respect to  $\phi^*$ , we have that

$$\mu D_\phi(x, y) \leq D_{J^\infty}(y, x),$$

and hence, applying this fact along with the bound Lemma 24-3, we obtain that

$$\begin{aligned} D_\phi(x^*, x_t) &\geq D_\phi(x^*, x_{t+1}) + D_{J^\infty}(x_{t+1}, x^*) + D_{J^\infty}(x^*, x_t) \\ &\geq D_\phi(x^*, x_{t+1}) + \mu D_\phi(x_{t+1}, x^*) + \mu D_\phi(x^*, x_t), \end{aligned}$$

Now, noting that  $D_\phi(x, y) = \phi(x - y)$ , we have that

$$\begin{aligned} \phi(x_{t+1} - x^*) &\leq \left( \frac{1 - \mu}{1 + \mu} \right) \phi(x_t - x^*) \\ &\leq \left( 1 - \frac{2\mu}{1 + \mu} \right) \phi(x_t - x^*), \end{aligned}$$

cascading this inequality, and noting that  $\mu \phi(x - y) \leq D_{J^\infty}(x, y) \leq \lambda \phi(x - y)$ , we obtain the second bound in the statement of the theorem.  $\blacksquare$

### E.8. Proof of Lemma 15

**Proof** We separate the proof into two parts, the first proving (13) and the second proving (14).

**Proof of (13):** Note that  $a_t = f(x_t) - f_* + \phi(\Delta x_t) \geq 0$  is non-increasing by Lemma 25 and that  $J_t^\infty = \sum_{s=t+1}^\infty a_s \leq \frac{C}{t}$  for  $C = \lambda \phi(x_0 - x^*)$  by equation (11). Hence,

$$t^2 a_{2t} \leq t \sum_{u=t+1}^{2t} a_u \leq t J_t^\infty \leq C .$$

Therefore, we have that

$$\limsup_{t \rightarrow \infty} t^2 a_t \leq 4C , \quad (54)$$

and hence by definition of the lim sup,  $a_t > \frac{4C}{t^2}$  for at most finitely many  $t$ , and we have the desired result.

**Proof of (14):** Here, we use a version of the *reverse Stolz-Cesàro* theorem. From equation (12), we have that

$$J_t^\infty = J^\infty(x_t) \leq c_0 e^{-c_1 t} = b_t$$

for  $c_0 = \lambda \phi(x_0 - x^*)$  and  $c_1 = -\log(1 - \frac{2\gamma}{1+\gamma})$ . Hence, we have that

$$\limsup_{t \rightarrow \infty} \frac{J_t^\infty}{b_t} \in [0, 1] \quad \text{and} \quad \lim_{t \rightarrow \infty} \frac{b_t}{b_{t+1}} = e^{-c_1} = B \neq 1 .$$

Hence, noting that

$$\frac{J_t^\infty - J_{t+1}^\infty}{b_{t+1} - b_t} = \frac{\frac{J_t^\infty}{b_t} \frac{b_t}{b_{t+1}} - \frac{J_{t+1}^\infty}{b_{t+1}}}{1 - \frac{b_t}{b_{t+1}}}$$

we compute

$$\begin{aligned} \limsup_{t \rightarrow \infty} \frac{f(x_{t+1}) + \phi(\Delta x_t)}{b_{t+1} - b_t} &= \limsup_{t \rightarrow \infty} \frac{J_t^\infty - J_{t+1}^\infty}{b_{t+1} - b_t} \\ &= \limsup_{t \rightarrow \infty} \frac{\frac{J_t^\infty}{b_t} \frac{b_t}{b_{t+1}} - \frac{J_{t+1}^\infty}{b_{t+1}}}{1 - \frac{b_t}{b_{t+1}}} \\ &\leq \limsup_{t \rightarrow \infty} \frac{\frac{J_t^\infty}{b_t} \frac{b_t}{b_{t+1}}}{1 - \frac{b_t}{b_{t+1}}} \leq \frac{B}{1-B} . \end{aligned}$$

Hence, we have that  $\limsup_{t \rightarrow \infty} \frac{f(x_{t+1}) + \phi(\Delta x_t)}{b_{t+1} - b_t} \leq \frac{B}{1-B}$ , which by definition implies that  $\frac{f(x_{t+1}) + \phi(\Delta x_t)}{b_{t+1} - b_t} > \frac{B}{1-B}$  at most finitely many times, giving the desired result.  $\blacksquare$

## E.9. Proof of Lemma 17

**Proof** Recall the definition of the Gâteaux derivative from equation (4),

$$\mathcal{R}'_T(\mathbf{x})(\delta \mathbf{x}) = \sum_{t=1}^T \langle \nabla \phi(\Delta x_t) - \nabla \phi(\Delta x_{t-1}) + \nabla f(x_t) , \delta x_{t-1} \rangle .$$

Computing the dual norm using the above expression, we find that

$$\begin{aligned} \left\| \mathcal{R}'_T(\boldsymbol{x}^\theta) \right\|_{2,*}^2 &= \sum_{t=1}^T \left\| \nabla \phi(\Delta x_t) - \nabla \phi(\Delta x_{t-1}) + \nabla f(x_t) \right\|^2 \\ &= \sum_{t=1}^T \mathcal{L}(\theta; x_{t-1}) , \end{aligned}$$

as desired. ■