

Convex-Concave Backtracking for Inertial Bregman Proximal Gradient Algorithms in Non-Convex Optimization

Mahesh Chandra Mukkamala* Peter Ochs† Thomas Pock‡ Shoham Sabach§

Abstract

Backtracking line-search is an old yet powerful strategy for finding better step size to be used in proximal gradient algorithms. The main principle is to locally find a simple convex upper bound of the objective function, which in turn controls the step size that is used. In case of inertial proximal gradient algorithms, the situation becomes much more difficult and usually leads to very restrictive rules on the extrapolation parameter. In this paper, we show that the extrapolation parameter can be controlled by locally finding also a simple concave lower bound of the objective function. This gives rise to a double convex-concave backtracking procedure which allows for an adaptive and optimal choice of both the step size and extrapolation parameters. We apply this procedure to the class of inertial Bregman proximal gradient methods, and prove that any sequence generated converges globally to critical points of the function at hand. Numerical experiments on a number of challenging non-convex problems in image processing and machine learning were conducted and show the power of combining inertial step and double backtracking strategy in achieving improved performances.

2010 Mathematics Subject Classification: Primary 90C25; Secondary 26B25, 49M27, 52A41, 65K05.

Keywords: Composite minimization, proximal gradient algorithms, inertial methods, convex-concave backtracking, non-Euclidean distances, Bregman distance, global convergence, Kurdyka-Łojasiewicz property.

1 Introduction

In this work we are interested in tackling non-convex additive composite minimization problems, which include the sum of two extended-valued functions: a non-smooth function denoted by f (possibly non-convex) and a smooth function denoted by g (possibly non-convex). More precisely, we consider problems of the following form

$$(\mathcal{P}) \quad \inf \{ \Psi(x) \equiv f(x) + g(x) : x \in \bar{C} \},$$

where \bar{C} is a nonempty, closed and convex set in \mathbb{R}^d . We will give a more precise statement in Section 2 about the involved functions. There is a tremendous number of applications in Machine Learning, Computer Vision, Statistics, and many more, that can be formulated in this framework.

Motivated by challenging applications as illustrated in Section 6, we consider here an instance of problem (\mathcal{P}) , where the smooth function g has a gradient that is not necessarily globally Lipschitz continuous. The restrictive assumption of having Lipschitz continuous gradient can be replaced with a certain convexity condition, which was proposed and developed first in [5] for problems (\mathcal{P}) with

*Faculty of Mathematics and Computer Science, Saarland University, 66123 Saarbrücken, Germany, E-mail: mukkamala@math.uni-sb.de

†Faculty of Mathematics and Computer Science, Saarland University, 66123 Saarbrücken, Germany, E-mail: ochs@math.uni-sb.de

‡Institute of Computer Graphics and Vision, Graz University of Technology, 8010 Graz, Austria. E-mail: pock@icg.tugraz.at

§Faculty of Industrial Engineering, The Technion, Haifa, 3200003, Israel. E-mail: ssabach@ie.technion.ac.il.

convex data, and recently extended to the non-convex setting in [12]. More details on these recent developments will be given below in Section 2.

This convexity condition easily yields an approximation of the objective function at hand by a convex function from above (majorant) and a concave function from below (minorant). In the traditional setting, where the gradient of the smooth function g is Lipschitz continuous, the majorant and the minorant are quadratic functions. In this case, it is well-known that the tightness of the quadratic approximations is directly related to restrictions on the step size to be used in the algorithm. The same relation is true for the convexity condition. In addition to their global existence, these approximations can be locally improved by backtracking (line search) strategies and it is well-known that tight approximations are advantageous, as we explain below in more details.

Interestingly, while the step size is usually restricted by the quality of the majorant, extrapolation (also known as inertia or over-relaxation) parameter is also affected by the quality of the minorant. This observation suggests to adapt the majorant and the minorant independently. In this paper we propose an efficient backtracking strategy that locally determines a tight majorant and minorant to exploit as much information as possible from the objective function, to be used in the proposed algorithm. This leads to a highly efficient algorithm, which is able to detect “the degree of local convexity” of the objective function (see Section 3 for details). As the backtracking procedure seeks for tight convex majorants and concave minorants, our idea is to combine it with an inertial step. We propose an inertial version of the Bregman Proximal Gradient (BPG) algorithm, which uses a convex-concave backtracking procedure to dynamically adjust the step size and the extrapolation parameter. Therefore, we call our algorithm *Convex-Concave Inertial BPG* (CoCaIn BPG in short). We prove a global convergence result of this algorithm (see Section 3.2 for an overview of the results and Section 5 for the details) to critical points of the objective function. The efficiency, which we demonstrate on several practical applications, comes from combining inertial step with the novel *convex-concave backtracking strategy*, which fully exploits the power of tight local approximations in achieving large step sizes and large extrapolation parameters that can be used at the same time.

Before concluding this section, we would like to give the reader a first intuition about the convex-concave backtracking strategy on a simple instance of problem (\mathcal{P}) . The setting is developed below.

A simple illustrative example. In the following, we consider the following particular instance of problem (\mathcal{P}) : $C = \mathbb{R}^d$, $f \equiv 0$ and the gradient of g is L -Lipschitz continuous. Even in this simpler setting, the convex-concave backtracking strategy is novel.

In this smooth and non-convex setting, an update step of a classical inertial based gradient method, starting with some $x^0 \in \mathbb{R}^d$, reads as follows

$$\begin{aligned} y^k &= x^k + \gamma_k (x^k - x^{k-1}), \\ x^{k+1} &= y^k - \frac{1}{\bar{L}_k} \nabla g (y^k), \end{aligned}$$

where $\gamma_k \in [0, 1]$, $k \in \mathbb{N}$, is an extrapolation parameter and $\bar{L}_k > 0$ is a step size. If g is convex and the extrapolation parameter γ_k is carefully chosen, this recovers the popular Nesterov Accelerated Gradient method [36] (for $f \neq 0$, again in the convex setting, see [7]). It is well-known that the gradient step above, can be equivalently written as follows

$$x^{k+1} = \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ g(y^k) + \langle \nabla g(y^k), x - y^k \rangle + \frac{\bar{L}_k}{2} \|x - y^k\|^2 \right\}.$$

For sufficiently large \bar{L}_k , the function to be minimized above is a convex quadratic majorant of the function g (due to the classical Descent Lemma), which is a property that is also crucial for the convergence analysis of the algorithm. Classically, $\bar{L}_k \geq L$, $k \in \mathbb{N}$, is a sufficient condition to guarantee the existence of a quadratic majorant. However, locally, i.e., between the points y^k and x^{k+1} , the

parameter \bar{L}_k may be significantly smaller than the global Lipschitz constant L (which will immediately affect the step size of the algorithm). More precisely, note that the Descent Lemma,

$$\left| g(x) - g(y^k) - \langle \nabla g(y^k), x - y^k \rangle \right| \leq \frac{L}{2} \|x - y^k\|^2, \quad \forall x \in \mathbb{R}^d, \quad (1.1)$$

actually guarantees the existence of a quadratic minorant and a quadratic majorant that are determined by the same (global) parameter L . However, only the majorant limits the step size that is used in the algorithm. As shown in Figure 1, tighter approximations can be computed if the parameters of the minorant and the majorant are allowed to differ:

$$-\frac{\underline{L}_k}{2} \|x - y^k\|^2 \leq g(x) - g(y^k) - \langle \nabla g(y^k), x - y^k \rangle \leq \frac{\bar{L}_k}{2} \|x - y^k\|^2, \quad (1.2)$$

i.e., the minorant parameter \underline{L}_k could be different from the majorant parameter \bar{L}_k .

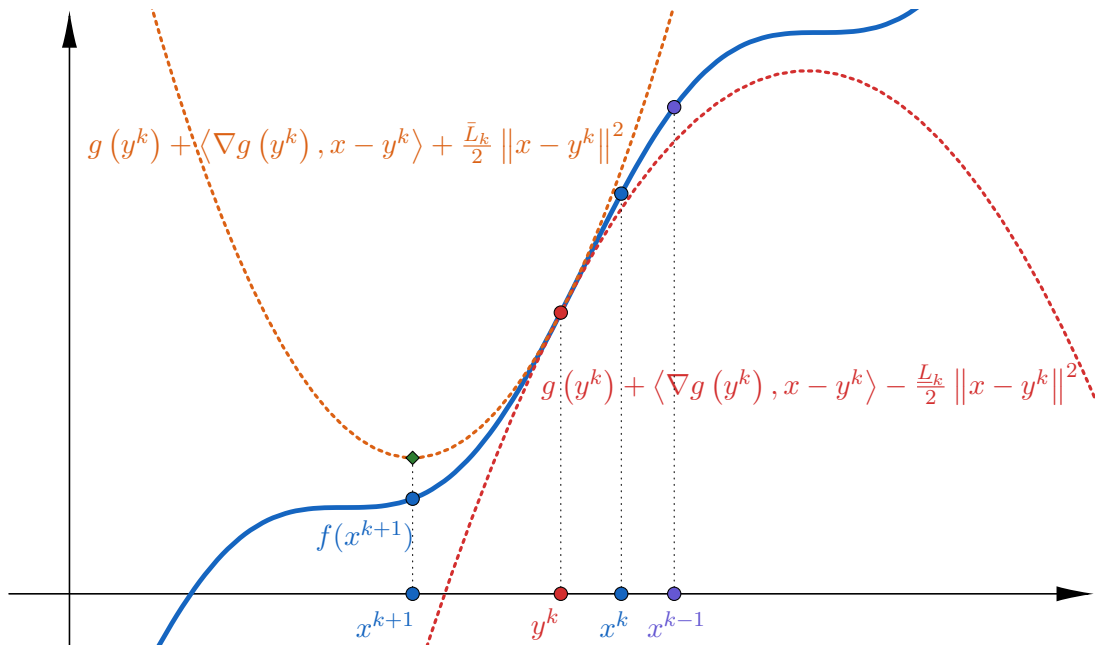


Figure 1: The inequalities in (1.2) guarantee that the objective function has a quadratic concave minorant and a quadratic convex majorant. The proposed convex-concave backtracking strategy locally estimates both the lower and the upper approximations using a double backtracking procedure.

While the step size of the algorithm only depends on the majorant parameter \bar{L}_k , the extrapolation parameter γ_k also depends on the minorant parameter \underline{L}_k . When $\bar{L}_k = \bar{L}$ and $\underline{L}_k = \underline{L}$, for all $k \in \mathbb{N}$, it was established in [54] that for any $0 \leq \gamma_k \leq \bar{\gamma}$, when

$$\bar{\gamma} < \sqrt{\frac{\bar{L}}{\underline{L} + \bar{L}}} \quad \left(= \frac{1}{\sqrt{2}} \text{ for } \bar{L} = \underline{L} \right),$$

the generated sequence converges linearly (under certain error bound condition).

If the minorant parameter \underline{L}_k is close to 0, which means that the function g is “locally convex”, the extrapolation parameter γ_k can be taken close to 1, which makes it “similar” to an Accelerated Gradient method in the non-convex setting.

Below, we will show that using the minorant and the majorant in a local fashion (instead of their global counterparts) is very useful in developing the inertial Bregman Proximal Gradient method.

Notation. We use standard notation and concepts which, unless otherwise specified, can all be found in [45].

2 The Bregman Framework

In this section we will first recall the definition of Bregman distance, which stands at the heart of our developments. Based on that we will shortly review the recent concept of smooth adaptable functions, which in some sense extends and generalizes the class of smooth functions with globally Lipschitz continuous gradient. Then, we will provide the basic and essential ingredients to deal with the Bregman Proximal Gradient method.

We begin with the notion of kernel generating distance functions, which was recently stated in [12] (in this respect see also [4]).

Definition 2.1. (Kernel Generating Distance) Let C be a nonempty, convex and open subset of \mathbb{R}^d . Associated with C , a function $h : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ is called a *kernel generating distance* if it satisfies the following:

- (i) h is proper, lower semicontinuous and convex, with $\text{dom } h \subset \overline{C}$ and $\text{dom } \partial h = C$.
- (ii) h is C^1 on $\text{int } \text{dom } h \equiv C$.

We denote the class of kernel generating distances by $\mathcal{G}(C)$.

Given $h \in \mathcal{G}(C)$, the *Bregman distance* that is associated to h , is a proximity measure $D_h : \text{dom } h \times \text{int } \text{dom } h \rightarrow \mathbb{R}_+$ which is defined by

$$D_h(x, y) := h(x) - [h(y) + \langle \nabla h(y), x - y \rangle].$$

This object is not a distance according to the classical definition (for example, it is not symmetric in general). However, the Bregman distance between two points is nonnegative if and only if the function h is convex. If h is known to be strictly convex, we have that $D_h(x, y) = 0$ if and only if $x = y$. The classic example of a Bregman distance is the squared Euclidean distance, which is generated by $h(x) = \|x\|^2$. For more examples, results and applications of Bregman distances, see [18, 49, 25, 6, 50] and references therein.

An important property that is always crucial when dealing with Bregman distances is the well-known *three-points identity* [21, Lemma 3.1]: for any $y, z \in \text{int } \text{dom } h$ and $x \in \text{dom } h$,

$$D_h(x, z) - D_h(x, y) - D_h(y, z) = \langle \nabla h(y) - \nabla h(z), x - y \rangle. \quad (2.1)$$

We conclude this part by restating *our optimization model*

$$(\mathcal{P}) \quad \inf \{ \Psi \equiv f(x) + g(x) : x \in \overline{C} \},$$

and making the first connection to the Bregman framework. One important feature of using Bregman distances in optimization algorithms is the ability of relate the constraint set C to a certain kernel generating distances function $h \in \mathcal{G}(C)$. From now on, we make the following assumption.

Assumption A. (i) $h \in \mathcal{G}(C)$ with $\overline{C} = \overline{\text{dom } h}$.

(ii) $f : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ is a proper and lower semicontinuous function (possibly non-convex) with $\text{dom } f \cap C \neq \emptyset$.

(iii) $g : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ is a proper and lower semicontinuous function (possibly non-convex) with $\text{dom } h \subset \text{dom } g$, which is continuously differentiable on C .

(iv) $v(\mathcal{P}) := \inf \{ \Psi(x) : x \in \overline{C} \} > -\infty$.

2.1 Smooth Adaptable Functions

One goal of this work is to deal with the non-convex optimization model (\mathcal{P}) where the gradient of the smooth function g is not globally Lipschitz. Recently, Bauschke, Bolte and Teboulle [5], observed that the property of having Lipschitz continuous gradient can be interpreted equivalently as a certain convexity condition on the function itself. This opens the gate for generalizing known results in the convex setting. It was extended to the non-convex setting in [12] with the concept of smooth adaptable functions given below.

Definition 2.2 (*L-smooth Adaptable*). A pair (g, h) is called *L-smooth adaptable* (**L-smad**) on C if there exists $L > 0$ such that $Lh - g$ and $Lh + g$ are convex on C .

The convexity requirement of $Lh + g$ can be written with respect to a different parameter $\ell \leq L$, which is key to the proposed double backtracking procedure to be developed in Section 3.1. In this section, for the sake of simplicity, we use $\ell = L$ for simplicity.

The optimization model (\mathcal{P}) appears with a smooth term in the objective function which is very common in many fields of applications. A crucial pillar in designing and analyzing algorithms for tackling this model, is usually based on the fact that the smooth part in the objective function has a Lipschitz continuous gradient. This property, via the well-known Descent Lemma, guarantees us that a lower and an upper quadratic approximations exist. For *L-smooth adaptable* functions, we will use the following extended version of the Descent Lemma (see [12, Lemma 2.1, p. 2134]).

Lemma 2.1 (*Extended Descent Lemma*). *The pair of functions (g, h) is L-smooth adaptable on C if and only if:*

$$|g(x) - g(y) - \langle \nabla g(y), x - y \rangle| \leq LD_h(x, y), \quad \forall x, y \in \text{int dom } h. \quad (2.2)$$

Remark 2.1 (*Invariance to Strong Convexity*). We would like to note that the *L-smooth adaptable* property is *invariant* when h is additionally assumed to be σ -strongly convex. Indeed, as described in [12], since convexity of g is not needed, we can define $\omega(x) := (\sigma_1/2) \|x\|^2$, and then for any $0 < \sigma_1 < \sigma$, we have

$$Lh - g = L(h - \omega) - (g - L\omega) := L\bar{h} - \bar{g},$$

namely, the new pair (\bar{g}, \bar{h}) satisfies the **L-smad** property on C .

2.2 The Bregman Proximal Gradient Algorithm

In this section we review the basic notations and results needed to study Bregman based optimization methods. We first recall the definition of the Bregman proximal mapping [49], which associated with a proper and lower semi-continuous function $f : \mathbb{R}^d \rightarrow (-\infty, +\infty]$, and defined by

$$\text{prox}_f^h(x) \in \text{argmin} \left\{ f(u) + D_h(u, x) : u \in \mathbb{R}^d \right\}, \quad \forall x \in \text{int dom } h.$$

With $h \equiv (1/2) \|\cdot\|^2$, the above boils down to the classical set-valued *Moreau proximal mapping* introduced in [34]. We refer the reader to the recent survey paper [50], and references therein. Here, we will focus on the Bregman proximal gradient mapping, which will take a central role in the algorithm to be developed in the next section. Given $x \in \text{int dom } h$ and a step size $\tau > 0$, the *Bregman proximal gradient* mapping is defined by

$$\begin{aligned} T_\tau(x) &\in \text{argmin} \left\{ f(u) + \langle \nabla g(x), u - x \rangle + \frac{1}{\tau} D_h(u, x) : u \in \bar{C} \right\} \\ &= \text{argmin} \left\{ f(u) + \langle \nabla g(x), u - x \rangle + \frac{1}{\tau} D_h(u, x) : u \in \mathbb{R}^d \right\}, \end{aligned} \quad (2.3)$$

where the second equality follows from the fact that $\text{dom } h \subset \bar{C}$. Note that here with $h \equiv (1/2) \|\cdot\|^2$, the above recovers the classical proximal gradient mapping. Since f could be non-convex, the mapping

T_τ is not, in general, single-valued. This mapping emerges from the usual approach, which consists of linearizing the differentiable function g around a point x and regularizing it with a proximal distance from that point. Similar to [12], the following assumption guarantees that the Bregman proximal gradient mapping is well-defined.

Assumption B. (i) The function $h + \tau f$ is supercoercive for all $\tau > 0$, that is,

$$\lim_{\|u\| \rightarrow \infty} \frac{h(u) + \tau f(u)}{\|u\|} = \infty.$$

(ii) For all $x \in C$, we have $T_\tau(x) \subset C$.

Assumption B(i) is a standard coercivity condition, which is for instance automatically satisfied when \overline{C} is compact. On the other hand, Assumption B(ii) can be shown to hold under a classical constraint qualification condition. It also holds automatically when f is convex or when $C = \mathbb{R}^d$. The following result from [12], ensures that the Bregman proximal gradient mapping is well-defined.

Lemma 2.2 (Well-Posedness of T_τ). *Suppose that Assumptions A and B hold, and let $x \in \text{int dom } h$. Then, the set $T_\tau(x)$ is a nonempty and compact subset of $\text{int dom } h$.*

3 The Inertial Bregman Proximal Gradient Method

Our proposed algorithm belongs to the class of inertial based optimization methods. The most well-known method in this class is the so-called Heavy-ball method, which was introduced by Polyak [44] to minimize convex and smooth functions. A popular variant of the method, which applied to the additive composite model (\mathcal{P}) with $C = \mathbb{R}^d$, takes the following form. Start with any $x^0 = x^1 \in \mathbb{R}^d$, and generate iteratively a sequence $\{x^k\}_{k \in \mathbb{N}}$ via

$$y^k = x^k + \gamma_k (x^k - x^{k-1}), \quad (3.1)$$

$$x^{k+1} \in \operatorname{argmin}_u \left\{ f(u) + \langle \nabla g(y^k), u - y^k \rangle + \frac{1}{2\tau_k} \|u - y^k\|^2 \right\}, \quad (3.2)$$

where $\gamma_k \in [0, 1]$ is an *extrapolation* parameter and $\tau_k > 0$ is a *step size*. In [41], an inertial proximal gradient algorithm, called *iPiano*, was proposed¹. It was shown that under Assumption A, if f is convex and g has a globally Lipschitz continuous gradient, the sequence $\{x^k\}_{k \in \mathbb{N}}$ converges globally to a critical point (in this setting, under additional error-bound condition, a linear rate of convergence was proved in [54]). The case where also the function f is not necessarily convex was treated in [13, 38]. Two years later, in [43] a block version of the method, called *iPALM* was proposed and analyzed in the fully non-convex setting, i.e., both f and g are non-convex. In this case, a global convergence result to critical points was also established. A unified analysis was presented in [40].

In this work we propose a Bregman variant of the method mentioned above (see steps (3.1) and (3.2)), which also handle the two involved parameter γ_k and τ_k , $k \in \mathbb{N}$, in a dynamic fashion. To this end we incorporate into our basic steps two routines aiming at controlling and updating these parameters.

3.1 The Convex-Concave Backtracking Procedure

As we already illustrated on a simple example in the introduction, the origin of this procedure comes from the fact that for smooth adaptable functions we can build lower and upper approximations as given in Lemma 2.1:

$$-\underline{L}D_h(x, y) \leq g(x) - g(y) - \langle \nabla g(y), x - y \rangle \leq \bar{L}D_h(x, y), \quad \forall x, y \in \text{int dom } h. \quad (3.3)$$

¹With a small modification that the proximity term is centered around the extrapolated point y^k , while the gradient of g is evaluated at x^k .

Even though the existence of the parameters \underline{L} and \bar{L} could be globally guaranteed, in practice it is often difficult or computationally expensive to evaluate them. In such cases it is recommended to apply a backtracking procedure that can locally verify the validity of the inequalities given in (3.3). However, in most cases only the upper approximation and the corresponding parameter \bar{L} are used. Here, we will develop a double backtracking procedure that locally verifies both the lower and the upper approximations, in order to better control and update the extrapolation parameter γ_k and the step size τ_k at each iteration $k \in \mathbb{N}$. To the best of our knowledge, this is the first attempt to use the lower approximation in algorithms for tackling non-convex problems. It should be noted that in the case that g is convex we have by definition $\underline{L} = 0$, or even a convex quadratic lower approximation can be found when g is strongly convex (see [50] for a discussion and references about a strong convexity property with respect to a Bregman distance). Based on the concepts described above, we will make the following additional assumptions on the involved functions.

Assumption C. (i) The function $h : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ is σ -strongly convex on \mathbb{R}^d .

(ii) The pair of functions (g, h) is L -smooth adaptable on C .

(iii) There exists $\alpha \in \mathbb{R}$ such that $f(\cdot) - (\alpha/2) \|\cdot\|^2$ is convex².

Few comments on the assumption above are now in order. The first item is related to Remark 2.1, which says that the smooth adaptable property is invariant to strongly convex kernel generating distance functions h . The third assumption allows us to deal with non-convex functions f since α could be negative. See Section 6 for examples of functions that satisfy all these assumptions. Now we are ready to present our algorithm, which is called Convex-Concave Inertial (CoCaIn) Bregman Proximal Gradient.

Convex-Concave Inertial BPG

Input. $\delta, \varepsilon > 0$ with $1 > \delta > \varepsilon$.

Initialization. $x^0 = x^1 \in \text{int dom } h \cap \text{dom } f$, $\bar{L}_0 > \frac{-\alpha}{(1-\delta)\sigma}$ and $\tau_0 \leq \bar{L}_0^{-1}$.

General Step. For $k = 1, 2, \dots$, compute

$$y^k = x^k + \gamma_k (x^k - x^{k-1}) \in \text{int dom } h, \quad (3.4)$$

where γ_k is chosen such that

$$(\delta - \varepsilon) D_h(x^{k-1}, x^k) \geq (1 + \underline{L}_k \tau_{k-1}) D_h(x^k, y^k) \quad (3.5)$$

holds and such that \underline{L}_k satisfies

$$g(x^k) \geq g(y^k) + \langle \nabla g(y^k), x^k - y^k \rangle - \underline{L}_k D_h(x^k, y^k). \quad (3.6)$$

Now, choose $\bar{L}_k \geq \bar{L}_{k-1}$, set $\tau_k \leq \min\{\tau_{k-1}, \bar{L}_k^{-1}\}$ and compute

$$x^{k+1} \in \operatorname{argmin}_u \left\{ f(u) + \langle \nabla g(y^k), u - y^k \rangle + \frac{1}{\tau_k} D_h(u, y^k) \right\} \quad (3.7)$$

with \bar{L}_k fulfilling

$$g(x^{k+1}) \leq g(y^k) + \langle \nabla g(y^k), x^{k+1} - y^k \rangle + \underline{L}_k D_h(x^{k+1}, y^k). \quad (3.8)$$

²Such functions are called semi-convex with modulus α (see [38, 39]).

The two input parameters δ and ε are free to be chosen by the user. As we will see later the parameter ε measures the descent to be achieved at each iteration of the algorithm. While δ bounds from below the parameter \bar{L}_k , which affects the step size of the algorithm according to the relation to τ_k (more details about these aspects will be given below in Section 6). The steps (3.4) and (3.7) are the classical steps of the inertial proximal gradient method, while here since we are dealing with the Bregman variant, it must be guaranteed that the auxiliary vector y^k as defined in (3.4) belongs to $\text{int dom } h$. Otherwise the Bregman proximal gradient step (3.7) is not defined (see Section 2.2). Even though, in general, it is not easy to guarantee that, in our case this will not be an issue. Indeed, in order to derive global convergence results of Bregman based algorithms in the non-convex setting an essential assumption seems to be that the kernel generating distance function h has a full domain, i.e., $\text{dom } h = \mathbb{R}^d$ (see, for instance, [12] for more details about this limitation). The steps (3.6) and (3.8) implement the double backtracking procedure (see Section 5.4). The step (3.5) is designed to control the extrapolation parameter γ_k , $k \in \mathbb{N}$, and should be validate at each iteration. However, a natural question would be if such a parameter always exists? We postpone the positive answer to this question, to Section 4, and conclude this section with a list of our theoretical contributions.

3.2 Summary of the Convergence Results

Before we proceed with the well-posedness of CoCaIn BPG and the convergence analysis, we provide here a brief summary of our results.

- We show the *well-posedness of CoCaIn BPG*, in the sense that, one can always find γ_k such that (3.5) is satisfied for all $k \in \mathbb{N}$ (see Lemma 4.1). Moreover, we show that it suffices to know the Bregman symmetric coefficient $\alpha(h)$ (Definition 4.1), in order to estimate the extrapolation parameter γ_k , $k \in \mathbb{N}$.
- In the Euclidean setting, i.e., when $h = (1/2) \|\cdot\|^2$, we provide an *explicit formula for the maximal extrapolation parameter*

$$0 \leq \gamma_k \leq \bar{\gamma}, \quad \bar{\gamma} < \sqrt{\frac{\bar{L}_{k-1}}{\bar{L}_{k-1} + \underline{L}_k}},$$

which uses the majorant parameter \bar{L}_{k-1} from the previous iterate, which is a key for the efficient implementation of the proposed convex-concave backtracking procedure. When $\bar{L}_{k-1} = \underline{L}_k$, we easily recover that $\bar{\gamma} < 1/\sqrt{2}$.

- *Stability and convergence of the objective function values of CoCaIn BPG*, which relies on finding an appropriate sequence of Lyapunov functions $\{\Phi_\delta^k(x^k, x^{k-1})\}_{k \in \mathbb{N}}$ that enjoys a sufficient descent property (see Proposition 5.1).
- *Global convergence of a sequence generated by the CoCaIn BPG method* to critical points of the objective function Ψ (see Theorem 5.2). This result relies on the concept of Gradient-like Descent Sequences (see Definition 5.1), see [11].

4 Well-Posedness of CoCaIn BPG

Now, we would like to verify the well-posedness of the CoCaIn BPG algorithm. An important tool in achieving our goal is the recently introduced symmetry coefficient of a Bregman distance, which measures the lack of symmetry in $D_h(\cdot, \cdot)$, see [5].

Definition 4.1 (Symmetry Coefficient). Given $h \in \mathcal{G}(C)$, its *symmetry coefficient* is defined by

$$\alpha(h) := \inf \left\{ \frac{D_h(x, y)}{D_h(y, x)} : x, y \in \text{int dom } h, x \neq y \right\} \in [0, 1].$$

An important and immediate consequence of this definition is the fact that for all $x, y \in \text{int dom } h$ we have

$$\alpha(h) D_h(x, y) \leq D_h(y, x) \leq \alpha(h)^{-1} D_h(x, y), \quad (4.1)$$

where we have adopted the convention that $0^{-1} = +\infty$ and $+\infty \times r = +\infty$ for all $r \geq 0$. Clearly, the closer is $\alpha(h)$ to 1, the more symmetric D_h is with perfect symmetry when $\alpha(h) = 1$ (which holds if and only if $h = \|\cdot\|^2$).

To this end, we need to convince the reader about the existence of γ_k , $k \in \mathbb{N}$, which satisfies (3.5), i.e., that

$$(\delta - \varepsilon) D_h(x^{k-1}, x^k) \geq (1 + \underline{L}_k \tau_{k-1}) D_h(x^k, y^k),$$

holds true. The following result provides a positive answer.

Lemma 4.1 (General Extrapolation Behavior). *Given $h \in \mathcal{G}(C)$ with $\alpha(h) > 0$. Let $x_1, x_2, y \in \text{int dom } h$ and $y := x_1 + \gamma(x_1 - x_2)$ with $\gamma \geq 0$. Then, there exist $\kappa > 0$ and γ^* such that*

$$D_h(x_1, y) \leq \kappa D_h(x_2, x_1), \quad \forall \gamma \in [0, \gamma^*]. \quad (4.2)$$

Proof. From the three points identity (see (2.1)) we have

$$\begin{aligned} D_h(y, x_2) &= D_h(y, x_1) + D_h(x_1, x_2) + \langle \nabla h(x_1) - \nabla h(x_2), y - x_1 \rangle \\ &= D_h(y, x_1) + D_h(x_1, x_2) + \gamma \langle \nabla h(x_1) - \nabla h(x_2), x_1 - x_2 \rangle \\ &= D_h(y, x_1) + D_h(x_1, x_2) + \gamma (D_h(x_1, x_2) + D_h(x_2, x_1)). \end{aligned}$$

Now, from (4.1), we obtain that

$$D_h(y, x_2) \leq \frac{1}{\alpha(h)} [D_h(x_1, y) + (\gamma\alpha(h) + 1 + \gamma) D_h(x_2, x_1)].$$

On the other hand, since $x_1 = (y + \gamma x_2) / (1 + \gamma)$, we can use the fact that $u \rightarrow D_h(u, v)$, for a fixed $v \in \text{int dom } h$, is a convex function and therefore

$$D_h(x_1, y) \leq \frac{\gamma}{1 + \gamma} D_h(x_2, y) \leq \frac{\gamma}{\alpha(h)(1 + \gamma)} D_h(y, x_2),$$

where the last inequality follows from (4.1). By combining the last two inequalities we derive that

$$D_h(x_1, y) \leq \frac{\gamma}{\alpha(h)^2(1 + \gamma)} [D_h(x_1, y) + (\gamma\alpha(h) + 1 + \gamma) D_h(x_2, x_1)],$$

and, by re-arranging we have

$$D_h(x_1, y) \leq \frac{\gamma(\gamma\alpha(h) + 1 + \gamma)}{\alpha(h)^2(1 + \gamma) - \gamma} D_h(x_2, x_1).$$

First, it is easy to verify that for $\gamma < \alpha(h)^2 / (1 - \alpha(h)^2)$, the denominator is positive. In addition, to find γ such that

$$\frac{\gamma(\gamma\alpha(h) + 1 + \gamma)}{\alpha(h)^2(1 + \gamma) - \gamma} \leq \kappa,$$

we will use simple algebraic manipulations. Indeed, by re-arranging we have

$$\underbrace{\gamma^2(\alpha(h) + 1)}_a + \underbrace{\gamma(1 + \kappa - \alpha(h)^2 \kappa)}_b - \alpha(h)^2 \kappa \leq 0.$$

Since $\alpha(h)^2 \leq 1$, it follows that $b > 0$. We also have that $\Delta = b^2 + 4a\alpha(h)^2 \kappa > 0$, and thus there exists a positive root denoted by γ^* . Therefore, for any $\gamma \in [0, \gamma^*]$, the desired result follows. \square

Remark 4.1. Note that in the above lemma, γ^* depends only on the symmetry coefficient $\alpha(h)$. Therefore, for the Euclidean distance with $\alpha(h) = 1$, this implies that,

$$\gamma^* = \frac{-1 + \sqrt{1 + 8\kappa}}{4}.$$

However, for the Euclidean distance, the expression in (4.2), can be simplified significantly. Indeed, since we take $h = (1/2) \|\cdot\|^2$, then using the fact that $y^k - x^k = \gamma_k (x^k - x^{k-1})$ we obtain that $\gamma_k \leq \sqrt{\kappa}$. In the case of CoCaIn BPG, we have the following restriction on the maximal extrapolation parameter that can be used

$$\gamma_k \leq \sqrt{\frac{\delta - \varepsilon}{1 + \underline{L}_k \tau_{k-1}}} \leq \sqrt{\frac{(\delta - \varepsilon) \bar{L}_{k-1}}{\bar{L}_{k-1} + \underline{L}_k}}.$$

A related bound also appeared in [54] as we discussed in the introduction. When, the values of \underline{L}_k and \bar{L}_{k-1} are almost equal and $\delta - \varepsilon \approx 1$, then it is possible to choose the inertial parameter γ_k such that $\gamma_k \approx 1/\sqrt{2}$. We discuss more about bounds of γ_k , $k \in \mathbb{N}$, in Section 5.3.

5 Convergence Analysis of CoCaIn BPG

Before, we proceed to the convergence analysis, we need the following technical lemma.

Lemma 5.1 (Function Descent Property). *Let $\{x^k\}_{k \in \mathbb{N}}$ be a sequence generated by CoCaIn BPG. Then, for all $k \in \mathbb{N}$, we have*

$$\Psi(x^k) \geq \Psi(x^{k+1}) + \frac{1}{\tau_k} D_h(x^k, x^{k+1}) + \frac{\alpha}{2} \|x^{k+1} - x^k\|^2 - \left(\frac{1}{\tau_k} + \underline{L}_k\right) D_h(x^k, y^k). \quad (5.1)$$

Proof. Fix $k \geq 1$. From the convexity of $f(\cdot) - (\alpha/2) \|\cdot\|^2$, which holds thanks to Assumption C(iii), we obtain from the sub-gradient inequality [45, Example 8.8 and Proposition 8.12] that

$$f(x^k) - \frac{\alpha}{2} \|x^k\|^2 \geq f(x^{k+1}) - \frac{\alpha}{2} \|x^{k+1}\|^2 + \langle \xi^{k+1} - \alpha x^{k+1}, x^k - x^{k+1} \rangle,$$

where $\xi^{k+1} \in \partial f(x^{k+1})$. By rearranging the inequality we obtain

$$f(x^k) \geq f(x^{k+1}) + \frac{\alpha}{2} \|x^{k+1} - x^k\|^2 + \langle \xi^{k+1}, x^k - x^{k+1} \rangle. \quad (5.2)$$

From the optimality condition of step (3.7), we have that

$$\xi^{k+1} + \nabla g(y^k) + \frac{1}{\tau_k} (\nabla h(x^{k+1}) - \nabla h(y^k)) = \mathbf{0},$$

which combined with (5.2) yields that

$$\begin{aligned} f(x^k) &\geq f(x^{k+1}) + \frac{\alpha}{2} \|x^{k+1} - x^k\|^2 - \langle \nabla g(y^k), x^k - x^{k+1} \rangle \\ &\quad + \frac{1}{\tau_k} \langle \nabla h(y^k) - \nabla h(x^{k+1}), x^k - x^{k+1} \rangle \\ &= f(x^{k+1}) + \frac{\alpha}{2} \|x^{k+1} - x^k\|^2 - \langle \nabla g(y^k), x^k - x^{k+1} \rangle \\ &\quad + \frac{1}{\tau_k} (D_h(x^k, x^{k+1}) + D_h(x^{k+1}, y^k) - D_h(x^k, y^k)), \end{aligned}$$

where the last equality follows from the three-points identity (see (2.1)). On the other hand, using the lower approximation given in (3.6) and the upper approximation given in (3.8), we have that

$$g(x^k) \geq g(x^{k+1}) + \langle \nabla g(y^k), x^k - x^{k+1} \rangle - \underline{L}_k D_h(x^k, y^k) - \bar{L}_k D_h(x^{k+1}, y^k).$$

Combining the last two inequalities and using the fact that $\tau_k^{-1} \geq \bar{L}_k$, implies that

$$\Psi(x^k) \geq \Psi(x^{k+1}) + \frac{\alpha}{2} \|x^{k+1} - x^k\|^2 + \frac{1}{\tau_k} D_h(x^k, x^{k+1}) - \left(\frac{1}{\tau_k} + \underline{L}_k\right) D_h(x^k, y^k),$$

which completes the proof. \square

Since we are dealing with inertial based methods, which belongs to the class of non-descent methods, we can not expect to use classical convergence techniques for non-convex problems (see below for more information about it). In order to overcome the lack of descent, we will use the Lyapunov technique, which involves the construction of a sequence of new functions, which will be used to “better” measure the progress of the algorithm, where by progress we mean a decrement in the Lyapunov function values. In several cases a trivial Lyapunov function would be to use the function itself, however in the case of non-descent methods, it is not a good choice, since it does not capture well the behavior of the iterates. The behavior of two subsequent iterates must be taken into consideration along with the function, as observed in [41, 48].

5.1 Lyapunov Function Descent Property of CoCaIn BPG

Let $\{x^k\}_{k \in \mathbb{N}}$ be a sequence generated by CoCaIn BPG. We define, at iterate $k \in \mathbb{N}$, the following Lyapunov function

$$\Phi_\delta^k(x^k, x^{k-1}) = \tau_{k-1} \left(\Psi(x^k) - v(\mathcal{P}) \right) + \delta D_h(x^{k-1}, x^k). \quad (5.3)$$

This Lyapunov function involves two terms: (i) the term $\tau_{k-1} (\Psi(x^k) - v(\mathcal{P}))$, which measures the progress in original function values Ψ with respect to the global optimal value of problem (\mathcal{P}) and (ii) the term given by $\delta D_h(x^{k-1}, x^k)$, which ensures that the iterates stay close enough, with respect to the Bregman distance. Before, we motivate further the usage of this Lyapunov function, we show its descent property.

Proposition 5.1. *Let $\{x^k\}_{k \in \mathbb{N}}$ be a sequence generated by CoCaIn BPG. Then, for all $k \in \mathbb{N}$, we have*

$$\Phi_\delta^k(x^k, x^{k-1}) \geq \Phi_\delta^{k+1}(x^{k+1}, x^k) + \varepsilon D_h(x^{k-1}, x^k). \quad (5.4)$$

Proof. Multiplying (5.2) with τ_k , we obtain

$$\begin{aligned} \tau_k \left(\Psi(x^k) - v(\mathcal{P}) \right) &\geq \tau_k \left(\Psi(x^{k+1}) - v(\mathcal{P}) \right) + \frac{\alpha \tau_k}{2} \|x^{k+1} - x^k\|^2 + D_h(x^k, x^{k+1}) \\ &\quad - (1 + \underline{L}_k \tau_k) D_h(x^k, y^k). \end{aligned}$$

By the definition of the Lyapunov function Φ_δ^k and the fact that $\tau_k \leq \tau_{k-1}$ we have

$$\begin{aligned} \Phi_\delta^k(x^k, x^{k-1}) &\geq \Phi_\delta^{k+1}(x^{k+1}, x^k) + \frac{\alpha \tau_k}{2} \|x^{k+1} - x^k\|^2 + (1 - \delta) D_h(x^k, x^{k+1}) \\ &\quad + \delta D_h(x^{k-1}, x^k) - (1 + \underline{L}_k \tau_k) D_h(x^k, y^k). \end{aligned}$$

With $1 - \delta > 0$ and the strong convexity of $h(\cdot)$, that follows from Assumption C(i), we obtain

$$\frac{\alpha \tau_k}{2} \|x^{k+1} - x^k\|^2 + (1 - \delta) D_h(x^k, x^{k+1}) \geq \left(\frac{\alpha \tau_k}{2} + (1 - \delta) \frac{\sigma}{2} \right) \|x^{k+1} - x^k\|^2 \geq 0,$$

where the last inequality holds, since $\tau_k^{-1} \geq \bar{L}_k$ and $\bar{L}_k \geq -\alpha / (1 - \delta) \sigma$. Next, we observe that

$$D_h(x^k, y^k) \leq \frac{\delta - \varepsilon}{(1 + \underline{L}_k \tau_{k-1})} D_h(x^{k-1}, x^k) \leq \frac{\delta - \varepsilon}{(1 + \underline{L}_k \tau_k)} D_h(x^{k-1}, x^k),$$

where the first inequality is due to the step (3.5) of the algorithm and the second inequality is due to fact that $\tau_k \leq \tau_{k-1}$. By rearranging we obtain,

$$\delta D_h(x^{k-1}, x^k) - (1 + \underline{L}_k \tau_k) D_h(x^k, y^k) \geq \varepsilon D_h(x^{k-1}, x^k)$$

thus completing the proof. \square

Proposition 5.2. *Let $\{x^k\}_{k \in \mathbb{N}}$ be a sequence generated by CoCaIn BPG. Then, the following assertions hold:*

- (i) *The sequence $\left\{ \Phi_\delta^{k+1}(x^{k+1}, x^k) \right\}_{k \in \mathbb{N}}$ is nonincreasing.*
- (ii) *$\sum_{k=1}^{\infty} D_h(x^{k-1}, x^k) < \infty$, and hence the sequence $\{D_h(x^{k-1}, x^k)\}_{k \in \mathbb{N}}$ converges to zero.*
- (iii) *$\min_{1 \leq k \leq n} D_h(x^{k-1}, x^k) \leq \Phi_\delta^1(x^1, x^0) / (\varepsilon n)$.*

Proof. (i) This follows trivially from Proposition 5.1, since $\varepsilon > 0$.

(ii) Let n be a positive integer. Summing (5.4) from $k = 1$ to n we get

$$\sum_{k=1}^n D_h(x^{k-1}, x^k) \leq \frac{1}{\varepsilon} (\Phi_\delta^1(x^1, x^0) - \Phi_\delta^{n+1}(x^{n+1}, x^n)) \leq \frac{1}{\varepsilon} \Phi_\delta^1(x^1, x^0), \quad (5.5)$$

since $\Phi_\delta^{n+1}(x^{n+1}, x^n) \geq 0$. Taking the limit as $n \rightarrow \infty$, we obtain the first desired assertion, from which we immediately deduce that $\{D_h(x^{k-1}, x^k)\}_{k \in \mathbb{N}}$ converges to zero.

(iii) From (5.5) we also obtain,

$$n \min_{1 \leq k \leq n} D_h(x^{k-1}, x^k) \leq \sum_{k=1}^n D_h(x^{k-1}, x^k) \leq \frac{1}{\varepsilon} \Phi_\delta^1(x^1, x^0),$$

which after division by n yields the desired result. \square

In order to proceed with the global convergence analysis of CoCaIn BPG, we will need throughout the rest of this section, to additionally assume the following.

Assumption D. (i) $\text{dom } h = \mathbb{R}^d$.

(ii) ∇h and ∇g are Lipschitz continuous on any bounded subset of \mathbb{R}^d .

5.2 Global Convergence for CoCaIn BPG

In this subsection we show the global convergence result of CoCaIn BPG. The goal is to show that the whole sequence $\{x^k\}_{k \in \mathbb{N}}$, that is generated by CoCaIn BPG, converges to a critical point. To this end, we denote the set of critical points by

$$\text{crit } \Psi = \left\{ x \in \mathbb{R}^d : 0 \in \partial \Psi(x) \equiv \partial f(x) + \nabla g(x) \right\}.$$

Note that, such a set is well-defined due to Fermat's rule [45, Theorem 10.1, p. 422] and due to the concept of limiting subdifferential.

From now on we will make the following assumption regarding the sequence of majorant parameters $\{\bar{L}_k\}_{k \in \mathbb{N}}$: there exists an integer $K \in \mathbb{N}$ such that $\bar{L}_k = \bar{L}$ for all $k \geq K$ (K can be as large as the user wishes). It should be noted that thanks to Assumption C(ii) and Lemma 2.1, there exists a global majorant parameter \bar{L} such that (3.8) holds true for all $k \in \mathbb{N}$. On the other hand, since in anyway we require that the parameters do not decrease between two successive iterations, it makes sense that at some point we will stop changing them and continue with a fixed value. However, it

is very important not using the global parameter \bar{L} right from the beginning since in practice the parameter \bar{L}_k determined by (3.8) might be much smaller (especially in early stages of the algorithm).

In the second phase of the algorithm, i.e., when $k \geq K$, it also makes sense to assume that $\tau_k = \tau$ for all $k \geq K$ where $\tau \leq \bar{L}^{-1}$. This immediately suggest that our Lyapunov function can also be simplified. More precisely, we define the following new Lyapunov function:

$$\Psi_{\delta_1}(x^k, x^{k-1}) = \begin{cases} \Phi_{\delta}^k(x^k, x^{k-1}), & k < K, \\ \Psi(x^k) + \delta_1 D_h(x^{k-1}, x^k), & k \geq K, \end{cases} \quad (5.6)$$

where $\delta_1 = \delta/\tau$.

The global convergence result is based on showing that CoCaIn BPG generates a gradient-like descent sequence according to Definition 5.1 (see below). This involves three properties which need to be verified: “sufficient descent condition”, “relative error condition” and “continuity condition”. Such a convergence analysis is based on a recent technique, which was initiated by Attouch and Bolte [1], and later on was simplified and unified in [11]. A more general framework was proposed in [40].

The main tool that stands behind this technique is the *Kurdyka-Lojasiewicz* (KL) property [31, 32] (see [8] for the non-smooth case), which is properly defined in the appendix. This property has been used in several recent works that deal with non-convex optimization problems (see [1, 3, 11] for early foundational works). For more details and information on the KL property, we refer the reader to the following papers [8, 1, 10, 2, 3, 11, 40] and references therein.

Verifying that a given function satisfies the KL property could be difficult, however in their seminal work [8], Bolte, Daniilidis and Lewis prove that any proper, lower semicontinuous and semi-algebraic function satisfies the KL property on its domain. This important result makes this proof technique very powerful, since we are familiar with many semi-algebraic functions that appear very often in applications. In fact, the same result holds for (possibly non-smooth) functions that are definable in an o-minimal structure [8, 9]. For examples and more details about the relations between KL and other important notions, see [8, 10] and references therein.

In order to derive the global convergence of our algorithm we follow this proof technique that we shortly recall now. For the interested readers we refer to [12, Appendix 6, p. 2147], where a short and self-contained summary of this proof methodology can be found. It should be noted again that here we consider a modification, which fits non-descent methods like CoCaIn BPG.

Definition 5.1 (Gradient-like Descent Sequence). A sequence $\{x^k\}_{k \in \mathbb{N}}$ is called a *gradient-like descent sequence* for minimizing Ψ_{δ_1} if the following three conditions hold:

(C1) *Sufficient decrease condition.* There exists a positive scalar ρ_1 such that

$$\rho_1 \left\| x^k - x^{k-1} \right\|^2 \leq \Psi_{\delta_1}(x^k, x^{k-1}) - \Psi_{\delta_1}(x^{k+1}, x^k), \quad \forall k \in \mathbb{N}.$$

(C2) *Relative error condition.* There exist an integer $K \in \mathbb{N}$ and a positive scalar ρ_2 such that

$$\left\| w^{k+1} \right\| \leq \rho_2 \left(\left\| x^k - x^{k-1} \right\| + \left\| x^{k+1} - x^k \right\| \right), \quad w^{k+1} \in \partial \Psi_{\delta_1}(x^{k+1}, x^k), \quad \forall k \geq K.$$

(C3) *Continuity condition.* Let \bar{x} be a limit point of a subsequence $\{x^k\}_{k \in \mathcal{K}}$, then $\limsup_{k \in \mathcal{K} \subset \mathbb{N}} \Psi(x^k) \leq \Psi(\bar{x})$.

Based on Definition 5.1 and the KL property, the following global convergence result holds true. We provide its proof in the appendix.

Theorem 5.1 (Global Convergence). *Let $\{x^k\}_{k \in \mathbb{N}}$ be a bounded gradient-like descent sequence for minimizing Ψ_{δ_1} . If Ψ satisfy the KL property, then the sequence $\{x^k\}_{k \in \mathbb{N}}$ has finite length, i.e., $\sum_{k=1}^{\infty} \|x^{k+1} - x^k\| < \infty$ and it converges to $x^* \in \text{crit } \Psi$.*

Now, in a sequence of lemmas, we prove that CoCaIn BPG generates a gradient-like descent sequence for minimizing Ψ_{δ_1} . In order to prove condition (C1), we first note that Proposition 5.2 is also valid for the new Lyapunov function Ψ_{δ_1} as recorded now (for the sake of simplicity we omit the exact details of the proof, which is almost identical to the proof above).

Proposition 5.3. *Let $\{x^k\}_{k \in \mathbb{N}}$ be a sequence generated by CoCaIn BPG. Then, the following assertions hold:*

- (i) *The sequence $\{\Psi_{\delta_1}(x^{k+1}, x^k)\}_{k \in \mathbb{N}}$ is nonincreasing, converging and condition (C1) of Definition 5.1 holds true.*
- (ii) *$\sum_{k=1}^{\infty} D_h(x^{k-1}, x^k) < \infty$, and hence the sequence $\{D_h(x^{k-1}, x^k)\}_{k \in \mathbb{N}}$ converges to zero.*
- (iii) *$\min_{1 \leq k \leq n} D_h(x^{k-1}, x^k) \leq (\Psi_{\delta_1}(x^1, x^0) - \Psi_*) / (\varepsilon n)$ where $\Psi_* = v(\mathcal{P}) > -\infty$ (by Assumption A(iv)).*

Now we can prove the following result, which means that condition (C2) holds true.

Proposition 5.4. *Let $\{x^k\}_{k \in \mathbb{N}}$ be a bounded sequence generated by CoCaIn BPG. Then, there exist $w^{k+1} \in \partial \Psi_{\delta_1}(x^{k+1}, x^k)$ and a positive scalar ρ_2 such that*

$$\|w^{k+1}\| \leq \rho_2 \left(\|x^k - x^{k-1}\| + \|x^{k+1} - x^k\| \right), \quad \forall k \geq K.$$

Proof. Fix $k \geq K$. By the definition of the Lyapunov function $\Psi_{\delta_1}(\cdot, \cdot)$ we obtain that

$$\partial \Psi_{\delta_1}(x^{k+1}, x^k) = \left(\partial \Psi(x^{k+1}) + \delta_1 \nabla^2 h(x^{k+1})(x^{k+1} - x^k), \delta_1 (\nabla h(x^k) - \nabla h(x^{k+1})) \right).$$

Writing the optimality condition of the optimization problem which defines x^{k+1} (see (3.7) and recall that for $k \geq K$, we have that $\tau_k = \tau$) yields that

$$0 \in \partial f(x^{k+1}) + \nabla g(y^k) + \frac{1}{\tau} (\nabla h(x^{k+1}) - \nabla h(y^k)).$$

Therefore

$$\nabla g(x^{k+1}) - \nabla g(y^k) + \frac{1}{\tau} (\nabla h(y^k) - \nabla h(x^{k+1})) \in \partial \Psi(x^{k+1}),$$

and by defining

$$w_1^{k+1} \equiv \nabla g(x^{k+1}) - \nabla g(y^k) + \frac{1}{\tau} (\nabla h(y^k) - \nabla h(x^{k+1})) + \delta_1 \nabla^2 h(x^{k+1})(x^{k+1} - x^k),$$

and $w_2^{k+1} \equiv \delta_1 (\nabla h(x^k) - \nabla h(x^{k+1}))$ we obviously obtain that $w^{k+1} \in \partial \Psi_{\delta_1}(x^{k+1}, x^k)$ where $w^{k+1} = (w_1^{k+1}, w_2^{k+1})$. Since $\{x^k\}_{k \in \mathbb{N}}$ is a bounded sequence and both ∇h and ∇g are Lipschitz continuous on bounded subsets of \mathbb{R}^d (see Assumption D(ii)), there exists $M > 0$ such that

$$\begin{aligned} \|w_1^{k+1}\| &\leq \left\| \nabla g(x^{k+1}) - \nabla g(y^k) \right\| + \frac{1}{\tau} \left\| \nabla h(y^k) - \nabla h(x^{k+1}) \right\| + \delta_1 \left\| \nabla^2 h(x^{k+1}) \right\| \cdot \|x^{k+1} - x^k\| \\ &\leq M \left(1 + \frac{1}{\tau} \right) \|x^{k+1} - y^k\| + \delta_1 M \|x^{k+1} - x^k\|, \end{aligned}$$

where the last inequality follows also from the fact that $\|\nabla^2 h(x^{k+1})\| \leq M$, since ∇h is Lipschitz continuous on bounded subsets of \mathbb{R}^d . Using step (3.4) we obtain that

$$\begin{aligned} \|w_1^{k+1}\| &\leq M \left(1 + \frac{1}{\tau}\right) \left(\|x^{k+1} - x^k\| + \gamma_k \|x^k - x^{k-1}\|\right) + \delta_1 M \|x^{k+1} - x^k\| \\ &\leq M \left(1 + \delta_1 + \frac{1}{\tau}\right) \|x^{k+1} - x^k\| + M \left(1 + \frac{1}{\tau}\right) \|x^k - x^{k-1}\|, \end{aligned}$$

where we have used the fact that $\gamma_k \leq 1$, $k \in \mathbb{N}$. Since, we also have that

$$\|w_2^{k+1}\| = \delta_1 \|\nabla h(x^k) - \nabla h(x^{k+1})\| \leq \delta_1 M \|x^{k+1} - x^k\|,$$

the desired result is proved and condition (C2) also holds true. \square

Now we are left with showing that CoCaIn BPG generates a sequence that satisfies condition (C3).

Proposition 5.5. *Let $\{x^k\}_{k \in \mathbb{N}}$ be a bounded sequence generated by CoCaIn BPG. Let x^* be a limit point of a subsequence $\{x^k\}_{k \in \mathcal{K}}$, then $\limsup_{k \in \mathcal{K} \subset \mathbb{N}} \Psi(x^k) \leq \Psi(x^*)$.*

Proof. Consider a subsequence $\{x^{n_k}\}_{k \in \mathbb{N}}$ which converges to x^* (there exists such a subsequence since the sequence $\{x^k\}_{k \in \mathbb{N}}$ is assumed to be bounded). Using Proposition 5.3(ii) and the strong convexity of $h(\cdot)$, we obtain that $\lim_{k \rightarrow \infty} \|x^k - x^{k-1}\| = 0$. Therefore, the sequence $\{x^{n_k-1}\}_{k \in \mathbb{N}}$ also converges to x^* . From the definition of y^k , see (3.4), it also follows that $\{y^{n_k-1}\}_{k \in \mathbb{N}}$ also converges to x^* . In addition, since h is continuously differentiable on \mathbb{R}^d we have that $\lim_{k \rightarrow \infty} D_h(x^*, y^{n_k-1}) = 0$. Now, from (3.7), it follows (after some simplifications), for all $k \geq K$, that

$$f(x^k) \leq f(x^*) + \langle x^* - x^k, \nabla g(y^{k-1}) \rangle + \frac{1}{\tau} D_h(x^*, y^{k-1}) - \frac{1}{\tau} D_h(x^k, y^{k-1}).$$

Substituting k by n_k and letting $k \rightarrow \infty$, we obtain from the fact that g is continuously differentiable on \mathbb{R}^d , that

$$\limsup_{k \rightarrow \infty} f(x^{n_k}) \leq f(x^*).$$

Using this, and recalling that here g is continuous, we obtain that $\limsup_{k \in \mathcal{K} \subset \mathbb{N}} \Psi(x^{n_k}) \leq \Psi(x^*)$, where $\mathcal{K} = \{n_k : k \geq K\}$. \square

The global convergence of CoCaIn BPG now easily follows from our general result on gradient-like descent sequences (see Theorem 5.1)

Theorem 5.2 (Global Convergence of CoCaIn BPG). *Let $\{x^k\}_{k \in \mathbb{N}}$ be a bounded sequence generated by CoCaIn BPG. If f and g satisfy the KL property, then the sequence $\{x^k\}_{k \in \mathbb{N}}$ has finite length, i.e., $\sum_{k=1}^{\infty} \|x^{k+1} - x^k\| < \infty$ and it converges to $x^* \in \text{crit } \Psi$.*

Before, we conclude this section, we provide a simplified variant of CoCaIn BPG.

5.3 CoCaIn BPG Without Backtracking

Note that CoCaIn BPG uses a local estimate of the minorant and majorant parameters \underline{L}_k and \bar{L}_k , $k \in \mathbb{N}$, determined by the backtracking steps (3.6) and (3.8), respectively. However, when the global parameter L is known (guaranteed in Assumption C(ii)), we can skip the backtracking steps, and provide a simplified variant of CoCaIn BPG.

For the inertial step (5.9), when $h = (1/2) \|\cdot\|^2$ we can obtain that

$$\gamma_k \leq \sqrt{\frac{\delta - \epsilon}{2}},$$

with $\bar{L} = \underline{L}$. Using Remark 4.1, if $\delta - \epsilon \approx 1$, one could choose the extrapolation parameter as follows $\gamma_k \approx 1/\sqrt{2}$. However, in general, the closed form expression for γ_k is difficult to obtain, for which backtracking line-search strategy can be used.

CoCaIn BPG Without Backtracking**Input.** $\delta, \varepsilon > 0$ with $1 > \delta > \varepsilon$.**Initialization.** $x^0 = x^1 \in \text{int dom } h \cap \text{dom } f$, $L \geq \max\{\frac{-\alpha}{(1-\delta)\sigma}, L\}$ and $\tau_0 \leq L^{-1}$.**General Step.** For $k = 1, 2, \dots$, compute

$$y^k = x^k + \gamma_k (x^k - x^{k-1}) \in \text{int dom } h, \quad (5.7)$$

$$x^{k+1} \in \operatorname{argmin}_u \left\{ f(u) + \langle \nabla g(y^k), u - y^k \rangle + \frac{1}{\tau_k} D_h(u, y^k) \right\}, \quad (5.8)$$

where $\tau_k \leq \min\{\tau_{k-1}, L^{-1}\}$ and $\gamma_k \geq 0$ satisfies

$$(\delta - \varepsilon) D_h(x^{k-1}, x^k) \geq 2D_h(x^k, y^k). \quad (5.9)$$

5.4 Implementing the Double Backtracking Procedure

The update steps of CoCaIn BPG based on the double backtracking strategy (see steps (3.6) and (3.8)). Here, we describe some implementation details of these two steps. Note that the inner loops for finding the minorant and the majorant parameters \underline{L}_k and \bar{L}_k , $k \in \mathbb{N}$, are implemented in a sequential fashion. By this, we mean that at iteration $k \in \mathbb{N}$ we first execute the steps (3.4), (3.5) and (3.6) in order to compute an appropriate y^k , only then we proceed to steps (3.7) and (3.8) in order to compute x^{k+1} . Note that the fact that the sequence $\{\bar{L}_k\}_{k \in \mathbb{N}}$ does not decrease is crucial in order to decouple the steps (3.4) and (3.7). More precisely, we now describe the backtracking procedure to find \underline{L}_k . Let $\underline{\nu} > 1$ be a scaling parameter and arbitrarily initialize $\underline{L}_{k,0} > 0$. Then, we find the smallest $\underline{L}_k \in \{\underline{\nu}^0 \underline{L}_{k,0}, \underline{\nu}^1 \underline{L}_{k,0}, \underline{\nu}^2 \underline{L}_{k,0}, \dots\}$ that satisfies (3.6) and such that $\gamma_k \geq 0$ satisfies

$$D_h(x^k, y^k) \leq \frac{\delta - \varepsilon}{\underline{L}_k \tau_{k-1} + 1} D_h(x^{k-1}, x^k).$$

We can now describe the procedure to find \bar{L}_k . Let $\bar{\nu} > 1$ and initialize $\bar{L}_{k,0} := \bar{L}_{k-1}$, then we take the smallest $\bar{L}_k \in \{\bar{\nu}^0 \bar{L}_{k,0}, \bar{\nu}^1 \bar{L}_{k,0}, \bar{\nu}^2 \bar{L}_{k,0}, \dots\}$ that satisfies (3.8). Therefore, $\{\bar{L}_k\}_{k \in \mathbb{N}}$ is monotonically non-decreasing. Note, however, we do not require any monotonicity of the sequence $\{\underline{L}_k\}_{k \in \mathbb{N}}$.

The double backtracking strategy preserves the sign of \underline{L}_k , however, only $-\underline{L}_k \leq \bar{L}_k$ is required. Changing the sign of \underline{L}_k when the function is locally strongly convex might lead to additional acceleration. However, we leave this kind of adaptation for future work.

6 Numerical Experiments

Our goal in this section is to illustrate the performance of CoCaIn BPG in various situations. We start with a minimization of univariate functions, which emphasizes the power of incorporating inertial terms into the BPG algorithm and using the double backtracking procedure. Then we provide some insights on the following practical applications: Quadratic Inverse Problems in Phase Retrieval, Structured Matrix Factorization, Non-convex Robust Denoising with Non-convex Total Variation Regularization, and Non-convex Quadratic Problems with Cubic Regularization.

6.1 Finding Global Minima of Univariate Functions

We begin with two examples of minimizing univariate non-convex functions, which shed some light on the two main features of our algorithm: (i) inertial term, and (ii) double backtracking procedure. We consider unconstrained minimization of functions $g : \mathbb{R} \rightarrow \mathbb{R}$, which have Lipschitz continuous gradient, i.e., model (\mathcal{P}) with $d = 1$, $f \equiv 0$ and $C = \mathbb{R}$. The two functions are: $g(x) = \log(1 + x^2)$

and $g(x) = (1 + e^x)^{-1}$. We compare three methods: CoCaIn BPG with $h = (1/2) \|\cdot\|^2$ and refer to it as the *CoCaIn with Euclidean distance*, the classical *Gradient Descent* (GD) method with backtracking (which is actually CoCaIn with Euclidean distance and with $\gamma_k = 0$ for all $k \in \mathbb{N}$). We also use *iPiano*³ of [41] (with the inertial parameter set to 0.7). When using a backtracking procedure in the GD and iPiano methods, we mean that only the majorant parameter is used. We use the same initialization for all the algorithms and report the performance in Figure 2.

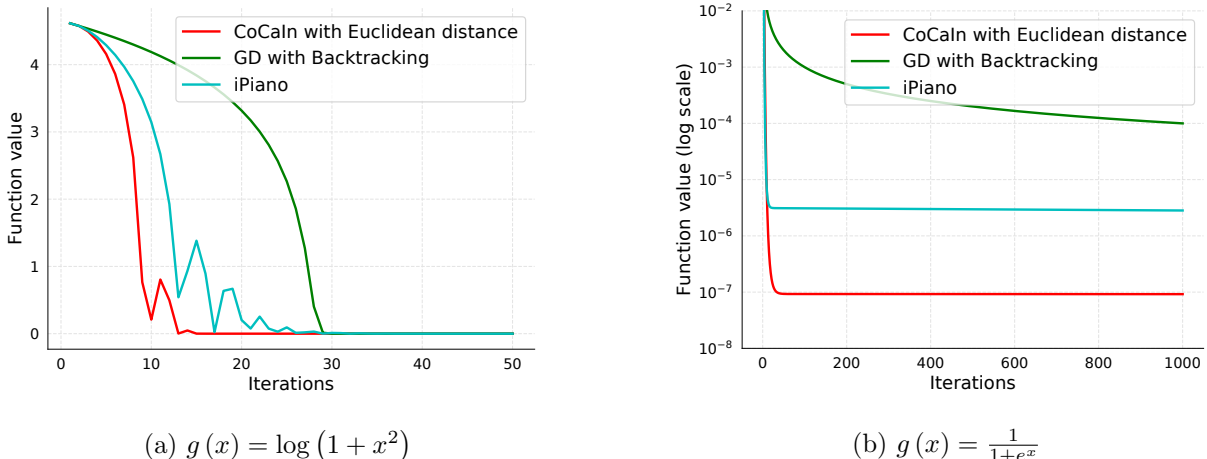


Figure 2: **Better performance by CoCaIn.** In the left-hand side plot, the function has a unique critical point. The CoCaIn BPG finds it faster than the other two methods. In the right-hand side plot, the function has a very small gradient and CoCaIn BPG reaches a significantly lower function value than the two other methods. These plots hint that CoCaIn BPG can significantly accelerate the convergence speed with comparison to GD and iPiano which use only simple backtracking procedure.

We provide now another example of minimizing a univariate function that is non-smooth and non-convex, with the purpose of illustrating an important feature of the CoCaIn BPG algorithm: sensitivity to local minima and critical points. In this experiment, we consider the non-smooth and non-convex function $\Psi(x) = |x| + \sin(x) + \cos(x)$, with many critical points as shown in the center plot of Figure 3. We take here $f(x) = |x|$ and $g(x) = \sin(x) + \cos(x)$ (which is obviously a non-convex function with Lipschitz continuous gradient). Here again we take $h = (1/2) \|\cdot\|^2$. In order to apply CoCaIn BPG, the main computational step is of the following form:

$$x^{k+1} \in \operatorname{argmin}_x \left\{ |x| + \left\langle x - y^k, \cos(y^k) - \sin(y^k) \right\rangle + \frac{1}{2\tau_k} (x - y^k)^2 \right\}, \quad (6.1)$$

which results in the following update step

$$x^{k+1} = \max \left\{ 0, \left| y^k - \tau_k \nabla g(y^k) \right| - \tau_k \right\} \operatorname{sgn} \left((y^k - \tau_k \nabla g(y^k)) \right). \quad (6.2)$$

We compare CoCaIn BPG with Euclidean distance to the classical *Proximal Gradient* (PG) method with backtracking (CoCaIn BPG with Euclidean distance and $\gamma_k = 0$, $k \in \mathbb{N}$), and also to *iPiano* as before. As mentioned in the first experiment, when using a backtracking procedure in the PG and iPiano methods we mean that only the majorant parameter is used.

As shown in Figure 3, CoCaIn BPG achieves the global minimum, whereas the Proximal Gradient method with backtracking gets stuck in a local minimum. We performed the same experiment starting at 100 equidistant points sampled from the interval $[-15, 15]$. The average final function value for CoCaIn was 2.75, whereas for the Proximal Gradient method with backtracking it was 3.21 and for

³In this particular case, the method coincides with the Heavy-ball method [44].

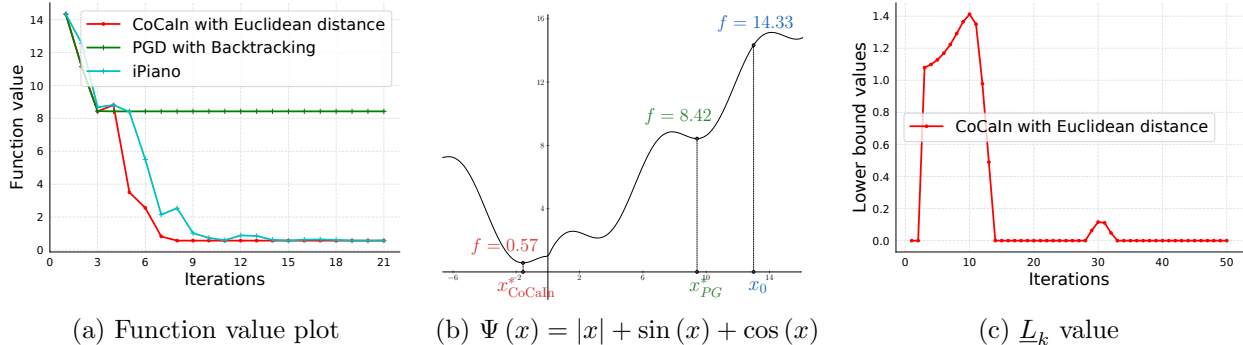


Figure 3: **CoCaIn can find the global minimum.** The left-hand side plot explicitly shows the behaviour in terms of function values versus the iterations counter. In the center plot, we use x_{PG}^* as a short hand notation for the critical point achieved by the Proximal Gradient method with backtracking, and for CoCaIn BPG method we use x_{CoCaIn}^* . The iPiano method achieves the same critical point as the CoCaIn BPG method but slower. In the right-hand side plot, we plot \underline{L}_k (the minorant parameter) obtained by CoCaIn BPG method versus the iterations counter. The hilly structures represent that CoCaIn BPG can bypass local maxima and eventually converge to zero. Meaning that CoCaIn BPG adapts to the “local convexity” of the function.

the iPiano it was 3.37. This means that CoCaIn BPG reaches the global minimum from 52 points, whereas the Proximal Gradient method with backtracking achieves the global minimum only from 27 points and the iPiano from 39 points. Hence, the behavior illustrated in Figure 3 is not due to the choice of initialization, but rather due to additional features of the CoCaIn BPG algorithm in comparison to the two other algorithms. This illustrates the great power of using double backtracking procedure in minimizing univariate non-convex functions.

6.2 Escaping Spurious Stationary Points

Here, we provide an evidence that CoCaIn BPG can escape spurious stationary points in minimizing non-convex functions of two variables. Let $b_i \in \mathbb{R}$, $i = 1, 2, \dots, m$, be samples of a noisy signal with additive Gaussian noise. A very common task in signal processing is to recover the true data. However, due to the noise, data can be prone to several outliers. In such cases, a robust loss [26] is used. Moreover, prior information about the data, can be embedded through a regularizing term (for instance, a sparsity promoting regularizer). Given $\lambda, \rho > 0$, we consider a minimization of the following form

$$\Psi(x) = \lambda \sum_{i=1}^m \log \left(1 + \rho (x_i - b_i)^2 \right) + \sum_{i=1}^m \log (1 + |x_i|), \quad (6.3)$$

with

$$f(x) := \sum_{i=1}^m \log (1 + |x_i|) \quad \text{and} \quad g(x) := \lambda \sum_{i=1}^m \log \left(1 + \rho (x_i - b_i)^2 \right).$$

Our goal is to recover the true data x . The function f is a non-convex sparsity promoting regularizer (also known as the log-sum penalty term [16, 37]) and the function g is a robust loss. For illustration purposes, we consider a simple instance of problem (6.3) where $m = 2$, $\lambda = 0.5$ and $\rho = 100$. For minimizing this function we set $C = \mathbb{R}^2$ and $h(x) := (1/2)(x_1^2 + x_2^2)$ to be used in the CoCaIn BPG method.

Before presenting the numerical results, we would like to note that in this example, the function $f(x) - (\alpha/2)h(x)$ is convex for any $\alpha \leq -1$ and $Lh - g$ is convex for all $L \geq 100$. Each iteration of CoCaIn BPG would require to compute the Bregman proximal gradient mapping, which in this case reduces to the classical proximal gradient mapping (due to the choice of h). Note that due to the separability of the functions f and g , the needed minimization problem can be split into

two individual minimizations with respect to x_1 and x_2 . These two optimization problems (after simple manipulations) reduces to the computation of the proximal mapping of the univariate function $\log(1 + |x|)$. A closed form formula can be found in [27] and reads as follows:

$$\text{prox}_{\tau \log(1+|x|)}(y) = \begin{cases} \text{sgn}(y) \operatorname{argmin}_{x \in E} \left\{ \log(1 + |x|) + \frac{1}{2\tau} (x - |y|)^2 \right\}, & \text{if } (|y| - 1)^2 - 4(\tau - |y|) \geq 0, \\ 0, & \text{otherwise,} \end{cases}$$

where

$$E = \left\{ 0, \left[\frac{|y| - 1 + \sqrt{(|y| - 1)^2 - 4(\tau - |y|)}}{2} \right]_+, \left[\frac{|y| - 1 - \sqrt{(|y| - 1)^2 - 4(\tau - |y|)}}{2} \right]_+ \right\},$$

with $[x]_+ := \max\{0, x\}$.

Now we can apply CoCaIn BPG method and the function behavior is described in Figure 4.

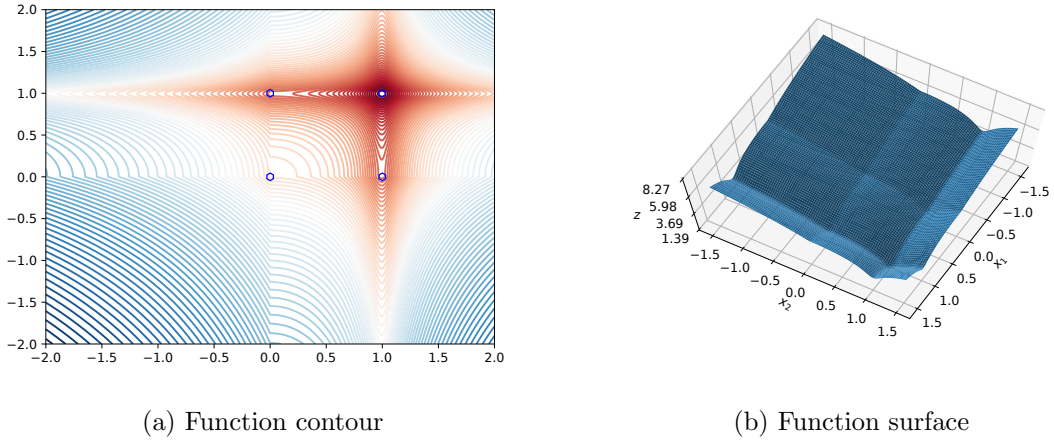


Figure 4: **Function with spurious stationary points.** The left-hand side plot shows the contours of the objective function, and the four critical points (denoted with blue diamond). In the right-hand side plot, we show the objective function, where the z -axis represents the function value. Here, the critical points appear as downward kink.

The performance of CoCaIn BPG is illustrated in Figure 5, which shows that CoCaIn BPG can indeed escape spurious critical points to reach the global minimum.

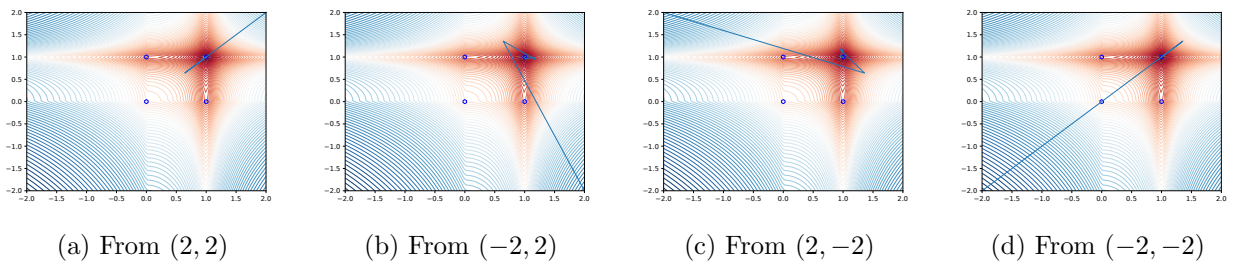


Figure 5: **CoCaIn can find the global minimum.** The CoCaIn BPG algorithm finds the global minimum at (1, 1), from various initialization points.

6.3 Quadratic Inverse Problems in Phase retrieval

Phase retrieval has been an active research topic for several years [15, 52, 24, 33]. It gained lot of attention from the optimization community, due to resulting hard non-convex problems [12, 24, 19]. The phase retrieval problem can be described as follows. Given sampling vectors $a_i \in \mathbb{R}^d$, $i = 1, 2, \dots, m$, and the measurements $b_i > 0$, we seek to find a vector $x \in \mathbb{R}^d$ such that the following system of quadratic equations is satisfied,

$$|\langle a_i, x \rangle|^2 \approx b_i^2, \quad \forall i = 1, 2, \dots, m. \quad (6.4)$$

One typical way to tackle this system of quadratic equations is by solving an optimization problem that seeks to minimize a certain error/noise measure in accomodating the equations. This usually leads to non-convex objectives (see [19]). The objective function also depends on the type of noise in the system (for instance, an objective function that is used for Gaussian noise can not be appropriate also for Poisson noise). We consider here, the additive Gaussian noise, and the related popular non-convex squared error measure, which is given by

$$\Psi(x) = f(x) + \frac{1}{4} \sum_{i=1}^m \left(\langle a_i, x \rangle^2 - b_i^2 \right)^2, \quad (6.5)$$

with

$$g(x) = \frac{1}{4} \sum_{i=1}^m \left(\langle a_i, x \rangle^2 - b_i^2 \right)^2.$$

The function f acts as a regularizing term and is used to incorporate certain prior information on the wished solution. We conduct experiments with two options of regularizing functions: (i) squared ℓ_2 -norm, that is, $f(x) = (\lambda/2) \|x\|^2$ and (ii) ℓ_1 -norm, that is, $f(x) = \lambda \|x\|_1$. When applying here the CoCaIn BPG method we use the following kernel generating distance function

$$h(x) = \frac{1}{4} \|x\|_2^4 + \frac{1}{2} \|x\|_2^2. \quad (6.6)$$

We obviously have that $\text{dom } h = \mathbb{R}^d$ and we record below a result [12, Lemma 5.1, p. 2143], which shows that the pair (g, h) satisfies the L-smad property (see Definition 2.2).

Lemma 6.1. *Let g and h be as defined above. Then, for any L satisfying*

$$L \geq \sum_{i=1}^m \left(3 \|a_i a_i^T\|^2 + \|a_i a_i^T\| |b_i^2| \right),$$

the function $Lh - g$ is convex on \mathbb{R}^d .

Therefore, in this case, all the Assumptions A, B, C and D are valid. We now discuss the update step of CoCaIn BPG, which requires the solution of the following subproblem

$$x^{k+1} \in \operatorname{argmin}_x \left\{ f(x) + \langle \nabla g(y^k), x - y^k \rangle + \frac{1}{\tau_k} D_h(x, y^k) \right\}. \quad (6.7)$$

Following [12], we provide closed form formulas for these optimization problems when f is either the squared ℓ_2 -norm or the ℓ_1 -norm.

ℓ_1 -norm. Here we use the following closed form solution, derived in [12, Proposition 5.1, p. 2145]. First, we define the soft-thresholding operator with respect to the parameter $\theta > 0$, as follows

$$\mathcal{S}_\theta(y) = \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \theta \|x\|_1 + \frac{1}{2} \|x - y\|^2 \right\} = \max\{|y| - \theta, 0\} \operatorname{sgn}(y), \quad (6.8)$$

where all operations are applied coordinate-wise. Then the closed form solution of problem (6.7) is given by

$$x^{k+1} = t^* \mathcal{S}_{\lambda\tau_k} \left(\nabla h \left(y^k \right) - \tau_k \nabla g \left(y^k \right) \right),$$

where t^* is the unique positive real root of the following cubic equation

$$t^3 \left\| \mathcal{S}_{\lambda\tau_k} \left(\nabla h \left(y^k \right) - \tau_k \nabla g \left(y^k \right) \right) \right\|_2^2 + t - 1 = 0.$$

Squared ℓ_2 -norm. Using similar arguments as of [12, Proposition 5.1, p. 2145], we can easily derive that the solution of problem (6.7) is given by

$$x^{k+1} = t^* \left(\tau_k \nabla g \left(y^k \right) - \nabla h \left(y^k \right) \right),$$

where t^* is the unique real root of the following cubic equation

$$t^3 \left\| \tau_k \nabla g \left(y^k \right) - \nabla h \left(y^k \right) \right\|^2 + (2\lambda\tau_k + 1)t + 1 = 0.$$

We illustrate, in Figure 6, the performance of CoCaIn BPG compared with two other algorithms: (i) the Bregman Proximal Gradient Method with backtracking (denoted by BPG-WB) using the same kernel generating distance function (which is exactly CoCaIn BPG with $\gamma_k = 0$ for all $k \in \mathbb{N}$) and (ii) the Inexact Bregman Proximal Minimization Line Search Algorithm (denoted by IBPM-LS) of [42]. We also compare with the Bregman Proximal Gradient (BPG) method of [12] without backtracking and with the parameter L as derived in Lemma 6.1.

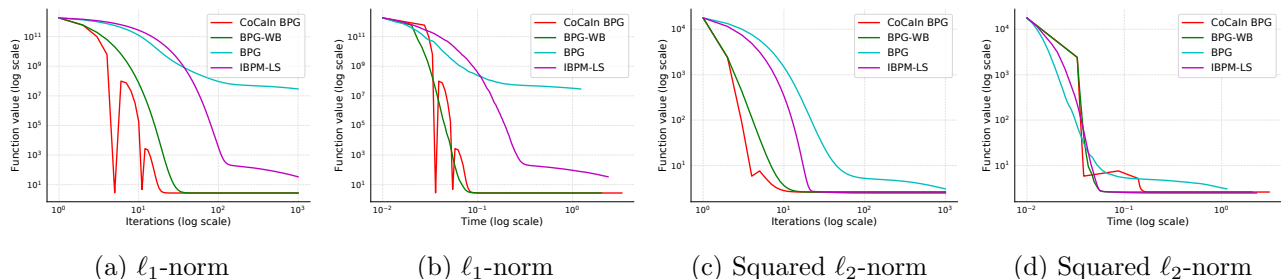


Figure 6: **CoCaIn BPG for Phase Retrieval.** All plots illustrate that CoCaIn BPG and BPG with Backtracking performances are much better. The difference is very significant when compared with BPG (without backtracking). This is due to the large L used in the algorithm, thus resulting in smaller steps. On the other hand, CoCaIn BPG uses the local parameters \underline{L}_k and \bar{L}_k , thus enjoys larger steps. The function values versus the time plots reveal that CoCaIn BPG rapidly attains a lower function value in a very early stage.

6.4 Non-convex Robust Denoising with Non-convex TV Regularization

We consider the problem of denoising a 1D signal, where the d measurements are provided via a vector $b \in \mathbb{R}^d$. The goal is to obtain the true signal, denoted by $x \in \mathbb{R}^d$. However, in real world applications, it is possible that the measurements are noisy with outliers. The standard routine to deal with outliers is to use a loss function, which is robust to outliers. Some examples include, ℓ_1 -norm loss and Huber loss [26]. The basic idea is to heavily penalize small errors and reasonably penalize large errors. This is done to ensure that the predicted data x , is not influenced significantly by outliers. The above mentioned loss functions are convex, however it is rarely the case that non-convex losses are used. We consider a fully non-convex formulation of the problem, which includes a non-convex loss function along with a non-convex regularization. We illustrate two ways in which one could make use of the CoCaIn BPG algorithm, such that the required assumptions hold true and the involved steps of algorithms are easy to compute.

Non-smooth data fitting term. Let

$$f(x) := \sum_{i=1}^d \log(1 + |x_i - b_i|) \quad \text{and} \quad g(x) = \lambda \sum_{i=1}^{d-1} \log\left(1 + \rho(x_{i+1} - x_i)^2\right),$$

where $\lambda, \rho > 0$. It is easy to see that both, the loss function f and the regularizing function g are non-convex. Observe that here f is also non-smooth while g is smooth. Here, the function g is the non-convex variant of the popular Total Variation (TV) regularizer, which is used to prefer smooth signals while preserving sharp changes in the signal (such as edges of images). For an overview on non-convex regularizations we refer the reader to [37, 55]. We verify now that the functions f and g satisfy our assumptions with respect to $h(x) = (1/2) \|x\|^2$.

Lemma 6.2. *The pair (g, h) satisfies the L -smad property with $L \geq 8\lambda\rho$ and $f(x) - (\alpha/2) \|x\|^2$ is convex for all $\alpha \leq -1$.*

Proof. First, it is easy to prove the convexity of $f(x) - (\alpha/2) \|x\|^2$, by checking that its right-hand side derivative is monotonically increasing [28, Theorem 6.4], for all $\alpha \leq -1$. Secondly, the convexity of the function $Lh - g$, can be proven by showing the existence of $L_i > 0$, for all i , such that

$$\frac{L_i}{2} (x_i^2 + x_{i+1}^2) - \lambda \log\left(1 + \rho(x_{i+1} - x_i)^2\right), \quad (6.9)$$

is convex. The maximum eigenvalue of the hessian of the function $\lambda \log\left(1 + \rho(x_{i+1} - x_i)^2\right)$, is bounded from above by $4\lambda\rho$. Therefore, for any $L_i \geq 8\lambda\rho$, the convexity of the function defined in (6.9) is proven. This proves the desired result. \square

Due to separability of the function f , we can split the computation of the corresponding Bregman Proximal Gradient mapping, into the following separable subproblems

$$x_i^{k+1} \in \operatorname{argmin}_{x_i \in \mathbb{R}} \left\{ \log(1 + |x_i - b_i|) + \left\langle x_i - y_i^k, \frac{2\lambda(y_{i+1}^k - y_i^k)}{1 + \rho(y_{i+1}^k - y_i^k)^2} \right\rangle + \frac{1}{2\tau_k} (x_i - y_i^k)^2 \right\},$$

which as discussed in Section 6.2, can be reduced to the computation of the proximal mapping of the function $\log(1 + |x - b|)$.

Smooth data fitting term. Here, we consider the same optimization problem but with a smooth data fitting function

$$f(x) := \sum_{i=1}^d \log\left(1 + (x_i - b_i)^2\right).$$

In this case we still have that the function $f(x) - (\alpha/2) \|x\|^2$ is convex for all $\alpha \leq -d/4$. In addition, the main step of the CoCaIn BPG algorithm, can be split into separable subproblems with respect to each variable x_i . This subproblem is a minimization of a quasi-convex function, hence the negative gradient always points towards the global minimizer. This means that we can find global minimizer, via setting the derivative to zero, which reduces to a cubic equation that can be explicitly solved.

6.5 Non-convex Quadratic Problems with Cubic Regularization

Here we consider the broad class of problems, known as Non-convex Quadratic Problems with Cubic Regularization which take the following form

$$\min_{x \in \mathbb{R}^d} \left\{ \frac{1}{2} x^T A x + \langle b, x \rangle + \frac{\rho}{3} \|x\|^3 \right\}$$

where $A \in \mathbb{R}^{d \times d}$, $b \in \mathbb{R}^d$ and $\rho > 0$. For simplicity we assume that A is symmetric, even though it is straightforward to extend our discussion to non-symmetric matrices. These problems typically arise in trust-region methods and cubic regularized Newton methods, which is currently a very active topic of research [35, 17, 23, 46, 51, 29, 57, 53]. Gradient descent is a popular method to solve such problems [51], but the theoretical guarantees usually limited to objective functions with Lipschitz continuous gradient. However, it is clear that this is not the case here. Motivated by this limitation, we show below two ways, in which CoCaIn BPG can handle this class of challenging problems with theoretical guarantees. We do not aim here to provide the best algorithm to solve such problem, rather to illustrate the general applicability of CoCaIn BPG. We start with the first option, where we set the following functions:

$$g(x) := \frac{1}{2}x^T Ax + \langle b, x \rangle, \quad f(x) := \frac{\rho}{3} \|x\|^3 \quad \text{and} \quad h(x) = \frac{1}{2} \|x\|^2.$$

Since f is a full-domain convex function and g is quadratic, all the required assumptions hold true. In order to apply CoCaIn BPG algorithm, the following subproblem (ignoring constant terms) must be solved

$$x^{k+1} \in \operatorname{argmin}_x \left\{ \frac{\rho}{3} \|x\|^3 + \langle Ay^k + b, x - y^k \rangle + \frac{1}{2\tau_k} \|x - y^k\|^2 \right\}.$$

By writing the optimality condition at the point x^{k+1} , we obtain

$$\tau_k \rho \left\| x^{k+1} \right\| x^{k+1} + x^{k+1} = x^k - \tau_k Ax^k - \tau_k b.$$

Denote $z^k := y^k - \tau_k Ax^k - \tau_k b$, then we obviously have that $x^{k+1} = tz^k$ where t can be obtained by solving the following quadratic equation

$$\tau_k \rho t^2 \left\| z^k \right\| + t - 1 = 0.$$

Since $1 + 4\tau_k \left\| z^k \right\| > 0$, there exists a solution for t , and the solution for x^{k+1} follows trivially.

We present another way in which CoCaIn BPG can be applied. Set the following functions:

$$g(x) = \frac{1}{2}x^T Ax + \langle b, x \rangle + \frac{\rho}{3} \|x\|^3, \quad f \equiv 0 \quad \text{and} \quad h(x) = \frac{\rho}{3} \|x\|^3 + \frac{\|A\|}{2} \|x\|^2.$$

It is easy to show that the function $Lh - g$ is convex for all $L \geq 1$ and all the required assumptions hold true. Similar arguments as above show that the main step of CoCaIn BPG can also be computed explicitly.

6.6 Structured Matrix Factorization

We now illustrate the performance of CoCaIn BPG on Structured Matrix Factorization problems. We show that with proper initialization, CoCaIn BPG can result in a superior performance, even though L -smooth adaptable property does not hold. This hints at possible generalizations of CoCaIn BPG. We leave such extensions and the corresponding theoretical justifications for future work.

Structured matrix factorization⁴ is an important optimization problem. Its applications include, for example, dictionary learning, deconvolution and source separation. The main goal is to obtain the unknown matrix factors $U \in \mathbb{R}^{M \times K}$ and $Z \in \mathbb{R}^{K \times N}$ of the matrix $A \in \mathbb{R}^{M \times N}$ such that $A = UZ + Q$ where $Q \in \mathbb{R}^{M \times N}$ is the error term. The factors are obtained typically through an optimization problems of the following form

$$\min_{U, Z} \left\{ \Psi(U, Z) \equiv \frac{1}{2} \|A - UZ\|_F^2 + \lambda F_1(U) + \lambda F_2(Z) \right\}, \quad (6.10)$$

⁴We follow the tutorial at <http://www.albertaueyung.com/post/python-matrix-factorization/> and also we use the Nimfa python library (see <http://nimfa.biolab.si/>)

where F_1 and F_2 are regularization terms, $\lambda > 0$ a regularizing parameter. The data-fidelity function G is the Frobenius norm of the approximation error matrix Q . Here, we set $g = G$ and $f(U, Z) = \lambda F_1(U) + \lambda F_2(Z)$, with the obvious matrix to vector transformations. In applications like Collaborative Filtering [30] where A is sparse, the Frobenius norm is restricted to the non-zero indices of A . We refer to [22, 20, 47, 56, 42] for detailed summaries of matrix factorization formulations, algorithms and applications.

We consider experiments with two settings of regularization terms, (i) squared ℓ_2 -norm, where we use $F_1(U) = (1/2) \|U\|_F^2$ and $F_2(Z) = (1/2) \|Z\|_F^2$ and (ii) ℓ_1 -norm, where we use $F_1(U) = \|U\|_1$ and $F_2(Z) = \|Z\|_1$. We compare CoCaIn BPG with Euclidean distance as the Bregman distance, Proximal Gradient Descent (PGD) with backtracking (can be obtained by setting $\gamma_k = 0$, $k \in \mathbb{N}$, in CoCaIn BPG with Euclidean distance) and iPiano of [41] with extrapolation parameter set to 0.7. We use a toy example, where we generate a random sparse matrix $A \in \mathbb{R}^{200 \times 200}$ and report the performance in Figures 7a and 7b. We also use Medulloblastoma dataset⁵ [14], which deals with microarray data of child tumors, to obtain a dense matrix $A \in \mathbb{R}^{5893 \times 34}$. We report the performance of the algorithms in Figures 7c and 7d.

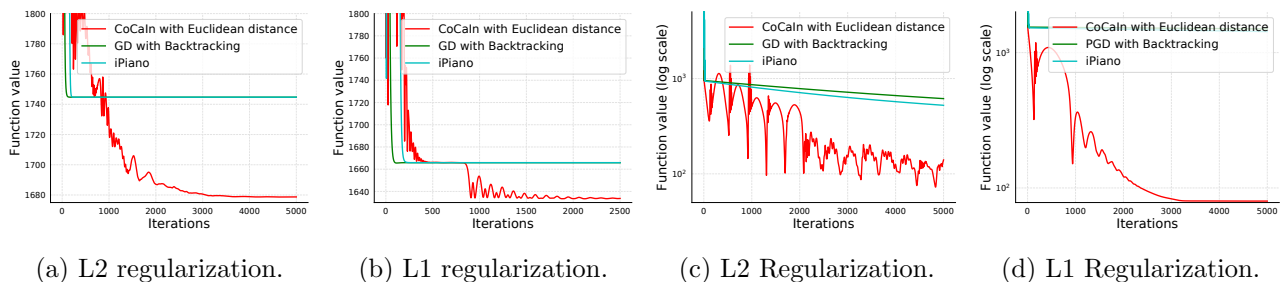


Figure 7: **CoCaIn for Matrix Factorization.** Figures 7a and 7b show the performance on the toy matrix $A \in \mathbb{R}^{200 \times 200}$, with $K = 10$. Here, iPiano and PGD with backtracking perform very similar, however CoCaIn BPG achieves a lower function value. Figures 7c and 7d are for the matrix $A \in \mathbb{R}^{5893 \times 34}$ (obtained from the Medulloblastoma dataset [14]), with $K = 2$. Here, CoCaIn BPG’s performance is an order of magnitude better in terms of function values in log scale compared with both iPiano and PGD with backtracking.

7 Acknowledgments

Mahesh Chandra Mukkamala and Peter Ochs acknowledge funding by the German Research Foundation (DFG Grant OC 150/1-1). Thomas Pock acknowledges support by the ERC starting grant HOMOVIS, No. 640156.

8 Appendix: Proof of Theorem 5.1

The set of all limit points of $\{x^k\}_{k \in \mathbb{N}}$ is defined by

$$\omega(x^0) := \left\{ \bar{x} \in \mathbb{R}^d : \exists \text{ an increasing sequence of integers } \{k_l\}_{l \in \mathbb{N}} \text{ such that } x^{k_l} \rightarrow \bar{x} \text{ as } l \rightarrow \infty \right\}.$$

We first prove the following result.

⁵More information regarding the dataset and its importance to the matrix factorization can be found at <http://nimfa.biolaab.si/nimfa.examples.medulloblastoma.html>

Lemma 8.1. *Let $\{x^k\}_{k \in \mathbb{N}}$ be a bounded gradient-like descent sequence for minimizing Ψ_{δ_1} . Then, $\omega(x^0)$ is a nonempty and compact subset of $\text{crit } \Psi$, and we have*

$$\lim_{k \rightarrow \infty} \text{dist}(x^k, \omega(x^0)) = 0. \quad (8.1)$$

In addition, the objective function Ψ is finite and constant on $\omega(x^0)$.

Proof. Since $\{x^k\}_{k \in \mathbb{N}}$ is bounded there is $x^* \in \mathbb{R}^d$ and a subsequence $\{x^{k_q}\}_{q \in \mathbb{N}}$ such that $x^{k_q} \rightarrow x^*$ as $q \rightarrow \infty$ and hence $\omega(x^0)$ is nonempty. Moreover, the set $\omega(x^0)$ is compact since it can be viewed as an intersection of compact sets. Now, from condition (C3), Proposition 5.3(ii) and the lower semicontinuity of Ψ (which follows from the lower semi-continuity of f and g , see Assumption A), we obtain

$$\lim_{q \rightarrow \infty} \Psi_{\delta_1}(x^{k_q+1}, x^{k_q}) = \lim_{q \rightarrow \infty} \Psi(x^{k_q}) = \Psi(x^*). \quad (8.2)$$

On the other hand, from conditions (C1) and (C2), we know that there is $w^k \in \partial \Psi_{\delta_1}(x^{k+1}, x^k)$, $k \in \mathbb{N}$, such that $w^k \rightarrow \mathbf{0}$ as $k \rightarrow \infty$. The closedness property of $\partial \Psi_{\delta_1}$ implies thus that $\mathbf{0} \in \partial \Psi_{\delta_1}(x^*, x^*) = \partial \Psi(x^*)$. This proves that x^* is a critical point of Ψ , and hence (8.1) is valid.

To complete the proof, let $\lim_{k \rightarrow \infty} \Psi_{\delta_1}(x^{k+1}, x^k) = l \in \mathbb{R}$. Then $\{\Psi_{\delta_1}(x^{k_q+1}, x^{k_q})\}_{q \in \mathbb{N}}$ converges to l and from (8.2) we have $\Psi(x^*) = l$. Hence the restriction of Ψ_{δ_1} to $\omega(x^0)$ equals l . \square

We recall now the definition of the Kurdyka-Łojasiewicz (KL) property [31, 32] and [8] (for the non-smooth case). Denote $[\alpha < \Psi < \beta] := \{x \in \mathbb{R}^d : \alpha < \Psi(x) < \beta\}$. Let $\eta > 0$, and set

$$\Phi_\eta = \{\varphi \in C^0[0, \eta) \cap C^1(0, \eta) : \varphi(0) = 0, \varphi \text{ concave and } \varphi' > 0\}.$$

Definition 8.1 (The Non-smooth KL Property). A proper and lower semicontinuous function $\Psi : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ has the Kurdyka-Łojasiewicz (KL) property locally at $\bar{u} \in \text{dom } \Psi$ if there exist $\eta > 0$, $\varphi \in \Phi_\eta$, and a neighborhood $U(\bar{u})$ such that

$$\varphi'(\Psi(u) - \Psi(\bar{u})) \text{dist}(0, \partial \Psi(u)) \geq 1,$$

for all $u \in U(\bar{u}) \cap [\Psi(\bar{u}) < \Psi(u) < \Psi(\bar{u}) + \eta]$.

Our last ingredient is a key uniformization of the KL property proven in [11, Lemma 6, p. 478], which we record below.

Lemma 8.2 (Uniformized KL Property). *Let Ω be a compact set and let $\Psi : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be a proper and lower semicontinuous function. Assume that Ψ is constant on Ω and satisfies the KL property at each point of Ω . Then, there exist $\tilde{\varepsilon} > 0$, $\eta > 0$ and $\varphi \in \Phi_\eta$ such that for all \bar{x} in Ω one has,*

$$\varphi'(\Psi(x) - \Psi(\bar{x})) \text{dist}(0, \partial \Psi(x)) \geq 1, \quad (8.3)$$

and all $x \in \{x \in \mathbb{R}^d : \text{dist}(x, \Omega) < \tilde{\varepsilon}\} \cap [\Psi(\bar{x}) < \Psi(x) < \Psi(\bar{x}) + \eta]$.

We can now restate and prove Theorem 5.1

Theorem 8.1. *Let $\{x^k\}_{k \in \mathbb{N}}$ be a bounded gradient-like descent sequence for minimizing Ψ_{δ_1} . If Ψ and h satisfy the KL property, then the sequence $\{x^k\}_{k \in \mathbb{N}}$ has finite length, i.e., $\sum_{k=1}^{\infty} \|x^{k+1} - x^k\| < \infty$ and it converges to $x^* \in \text{crit } \Psi$.*

Proof. Since $\{x^k\}_{k \in \mathbb{N}}$ is bounded there exists a subsequence $\{x^{k_q}\}_{q \in \mathbb{N}}$ such that $x^{k_q} \rightarrow \bar{x}$ as $q \rightarrow \infty$. In a similar way as in Lemma 8.1 (using also Proposition 5.3(ii)) we get that

$$\lim_{k \rightarrow \infty} \Psi_{\delta_1}(x^{k+1}, x^k) = \lim_{k \rightarrow \infty} \Psi(x^k) = \Psi(\bar{x}). \quad (8.4)$$

If there exists an integer \bar{k} for which $\Psi_{\delta_1}(x^{k+1}, x^k) = \Psi(\bar{x})$ then condition (C1) would imply that $x^{\bar{k}+1} = x^{\bar{k}}$. A trivial induction show then that the sequence $\{x^k\}_{k \in \mathbb{N}}$ is stationary and the announced results are obvious. Since $\{\Psi_{\delta_1}(x^{k+1}, x^k)\}_{k \in \mathbb{N}}$ is a nonincreasing sequence, it is clear from (8.4) that $\Psi(\bar{x}) < \Psi_{\delta_1}(x^{k+1}, x^k)$ for all $k > 0$. Again from (8.4) for any $\eta > 0$ there exists a nonnegative integer k_0 such that $\Psi_{\delta_1}(x^{k+1}, x^k) < \Psi(\bar{x}) + \eta$ for all $k > k_0$. From Lemma 8.1 we know that $\lim_{k \rightarrow \infty} \text{dist}(x^k, \omega(x^0)) = 0$. This means that for any $\tilde{\varepsilon} > 0$ there exists a positive integer k_1 such that $\text{dist}(x^k, \omega(x^0)) < \tilde{\varepsilon}$ for all $k > k_1$.

From Lemma 8.1 applied of Ψ_{δ_1} , we know that $\omega(x^0)$ is nonempty and compact and that the function Ψ is finite and constant on $\omega(x^0)$. Hence, we can apply the Uniformization Lemma 8.2 with $\Omega = \omega(x^0)$. Therefore, for any $k > l := \max\{k_0, k_1\}$, we have

$$\varphi' \left(\Psi_{\delta_1}(x^k, x^{k-1}) - \Psi(\bar{x}) \right) \text{dist} \left(0, \partial \Psi_{\delta_1}(x^k, x^{k-1}) \right) \geq 1. \quad (8.5)$$

This makes sense since we know that $\Psi_{\delta_1}(x^k, x^{k-1}) > \Psi(\bar{x})$ for any $k > l$. Combining (8.5) with condition (C2), see Proposition 5.4, we get that

$$\varphi' \left(\Psi_{\delta_1}(x^k, x^{k-1}) - \Psi(\bar{x}) \right) \geq \rho_2^{-1} \left(\|x^{k-1} - x^{k-2}\| + \|x^k - x^{k-1}\| \right)^{-1}. \quad (8.6)$$

For convenience, we define for all $p, q \in \mathbb{N}$ and \bar{x} the following quantity

$$\Delta_{k,k+1} := \varphi \left(\Psi_{\delta_1}(x^k, x^{k-1}) - \Psi(\bar{x}) \right) - \varphi \left(\Psi_{\delta_1}(x^{k+1}, x^k) - \Psi(\bar{x}) \right).$$

From the concavity of φ we get that

$$\Delta_{k,k+1} \geq \varphi' \left(\Psi_{\delta_1}(x^k, x^{k-1}) - \Psi(\bar{x}) \right) \left(\Psi_{\delta_1}(x^k, x^{k-1}) - \Psi_{\delta_1}(x^{k+1}, x^k) \right). \quad (8.7)$$

Combining condition (C1) with (8.6) and (8.7) yields, for any $k > l$, that

$$\Delta_{k,k+1} \geq \frac{\|x^{k+1} - x^k\|^2}{\rho (\|x^{k-1} - x^{k-2}\| + \|x^k - x^{k-1}\|)}, \quad \text{where } \rho := \rho_2 / \rho_1.$$

Using the fact that $2\sqrt{\alpha\beta} \leq \alpha + \beta$ for all $\alpha, \beta \geq 0$, we infer from the later inequality that

$$4 \|x^{k+1} - x^k\| \leq \|x^{k-1} - x^{k-2}\| + \|x^k - x^{k-1}\| + 4\rho \Delta_{k,k+1}. \quad (8.8)$$

Summing up (8.8) for $i = l+1, \dots, k$ yields

$$\begin{aligned} 4 \sum_{i=l+1}^k \|x^{i+1} - x^i\| &\leq \sum_{i=l+1}^k \|x^{i-1} - x^{i-2}\| + \sum_{i=l+1}^k \|x^i - x^{i-1}\| + 4\rho \sum_{i=l+1}^k \Delta_{i,i+1} \\ &\leq \sum_{i=l+1}^k \|x^{i+1} - x^i\| + \|x^l - x^{l-1}\| + \|x^{l+1} - x^l\| \\ &\quad + \sum_{i=l+1}^k \|x^{i+1} - x^i\| + \|x^l - x^{l-1}\| + 4\rho \sum_{i=l+1}^k \Delta_{i,i+1} \\ &= 2 \sum_{i=l+1}^k \|x^{i+1} - x^i\| + 2 \|x^l - x^{l-1}\| + \|x^{l+1} - x^l\| + 4\rho \Delta_{l+1,k+1}, \end{aligned}$$

where the last inequality follows from the fact that $\Delta_{p,q} + \Delta_{q,r} = \Delta_{p,r}$ for all $p, q, r \in \mathbb{N}$. Since $\varphi \geq 0$, recalling the definition of $\Delta_{l+1,k+1}$, we thus have for any $k > l$ that

$$2 \sum_{i=l+1}^k \|x^{i+1} - x^i\| \leq 2 \|x^l - x^{l-1}\| + \|x^{l+1} - x^l\| + 4\rho \varphi \left(\Psi_{\delta_1}(x^{l+1}, x^l) - \Psi(\bar{x}) \right),$$

which implies that $\sum_{k=1}^{\infty} \|x^{k+1} - x^k\| < \infty$, i.e., $\{x^k\}_{k \in \mathbb{N}}$ is a Cauchy sequence and hence together with Lemma 8.1, we obtain the global convergence to a critical point. \square

References

- [1] H. Attouch and J. Bolte. On the convergence of the proximal algorithm for nonsmooth functions involving analytic features. *Mathematical Programming*, 116(1-2):5–16, 2009.
- [2] H. Attouch, J. Bolte, P. Redont, and A. Soubeyran. Proximal alternating minimization and projection methods for nonconvex problems: an approach based on the Kurdyka-Łojasiewicz inequality. *Mathematics of Operations Research*, 35(2):438–457, 2010.
- [3] H. Attouch, J. Bolte, and B. F. Svaiter. Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward–backward splitting, and regularized Gauss–Seidel methods. *Mathematical Programming*, 137(1-2):91–129, 2013.
- [4] A. Auslender and M. Teboulle. Interior gradient and proximal methods for convex and conic optimization. *SIAM Journal on Optimization*, 16(3):697–725, 2006.
- [5] H. H. Bauschke, J. Bolte, and M. Teboulle. A descent lemma beyond Lipschitz gradient continuity: first-order methods revisited and applications. *Mathematics of Operations Research*, 42(2):330–348, 2017.
- [6] H. H. Bauschke and J. M. Borwein. Legendre functions and the method of random bregman projections. *Journal of Convex Analysis*, 4(1):27–67, 1997.
- [7] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [8] J. Bolte, A. Daniilidis, and A. Lewis. The Łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. *SIAM Journal on Optimization*, 17:1205–1223, 2006.
- [9] J. Bolte, A. Daniilidis, A.S. Lewis, and M. Shiota. Clarke subgradients of stratifiable functions. *SIAM Journal on Optimization*, 18(2):556–572, 2007.
- [10] J. Bolte, A. Daniilidis, O. Ley, and L. Mazet. Characterizations of Łojasiewicz inequalities: subgradient flows, talweg, convexity. *Transactions of the American Mathematical Society*, 362(6):3319–3363, 2010.
- [11] J. Bolte, S. Sabach, and M. Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1-2):459–494, 2014.
- [12] J. Bolte, S. Sabach, M. Teboulle, and Y. Vaisbourd. First order methods beyond convexity and Lipschitz gradient continuity with applications to quadratic inverse problems. *SIAM Journal on Optimization*, 28(3):2131–2151, 2018.
- [13] R. I. BoT, E. R. Csetnek, and S. C. László. An inertial forward–backward algorithm for the minimization of the sum of two nonconvex functions. *EURO Journal on Computational Optimization*, 4(1):3–25, 2016.
- [14] J.-P. Brunet, P. Tamayo, T. R. Golub, and J. P. Mesirov. Metagenes and molecular pattern discovery using matrix factorization. *Proceedings of the National Academy of Sciences*, 101(12):4164–4169, 2004.
- [15] E. J. Candes, X. Li, and M. Soltanolkotabi. Phase retrieval via wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007, 2015.
- [16] E. J. Candes, M. B. Wakin, and S. Boyd. Enhancing sparsity by reweighted ℓ_1 minimization. *Journal of Fourier analysis and applications*, 14(5-6):877–905, 2008.

- [17] Y. Carmon and J. C. Duchi. Analysis of krylov subspace solutions of regularized non-convex quadratic problems. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 10705–10715. Curran Associates, Inc., 2018.
- [18] Y. Censor and S. A. Zenios. Proximal minimization algorithm with D-functions. *Journal of Optimization Theory and Applications*, 73(3):451–464, 1992.
- [19] H. Chang, S. Marchesini, Y. Lou, and T. Zeng. Variational phase retrieval with globally convergent preconditioned proximal algorithm. *SIAM Journal on Imaging Sciences*, 11(1):56–93, 2018.
- [20] S. Chaudhuri, R. Velmurugan, and R. M. Rameshan. *Blind image deconvolution*. Springer, 2016.
- [21] G. Chen and M. Teboulle. Convergence analysis of a proximal-like minimization algorithm using bregman functions. *SIAM Journal on Optimization*, 3(3):5380–543, 1993.
- [22] A. Cichocki, R. Zdunek, A. H. Phan, and S. Amari. *Nonnegative matrix and tensor factorizations: applications to exploratory multi-way data analysis and blind source separation*. John Wiley & Sons, 2009.
- [23] A. R. Conn, N. I. M. Gould, and P. L. Toint. *Trust region methods*, volume 1. SIAM, 2000.
- [24] J.C. Duchi and F. Ruan. Solving (most) of a set of quadratic equalities: Composite optimization for robust phase retrieval. *ArXiv preprint arXiv:1705.02356*, 2017.
- [25] J. Eckstein. Nonlinear proximal point algorithms using Bregman functions, with applications to convex programming. *Mathematics of Operations Research*, 18(1):202–226, 1993.
- [26] J. Friedman, T. Hastie, and R. Tibshirani. *The elements of statistical learning*. Springer series in statistics New York, 2001.
- [27] P. Gong, C. Zhang, L. Zhaosong, J. Z. Huang, and J. Ye. A general iterative shrinkage and thresholding algorithm for non-convex regularized optimization problems. In S. Dasgupta and D. McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28, pages 37–45. PMLR, 2013.
- [28] J.-B. Hiriart-Urruty and C. Lemarechal. *Fundamentals of Convex Analysis*. Springer Science & Business Media, 2012.
- [29] J. M. Kohler and A. Lucchi. Sub-sampled cubic regularization for non-convex optimization. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 1895–1904. PMLR, 2017.
- [30] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.
- [31] K. Kurdyka. On gradients of functions definable in o-minimal structures. *Université de Grenoble. Annales de l’Institut Fourier*, 48(3):769–783, 1998.
- [32] S. Lojasiewicz. Une propriété topologique des sous-ensembles analytiques réels. In *Les Équations aux Dérivées Partielles (Paris, 1962)*, pages 87–89. Éditions du Centre National de la Recherche Scientifique, Paris, 1963.
- [33] D. R. Luke. Phase retrieval, What’s new? *SIAG/OPT Views and News*, 25(1):1–6, 2017.
- [34] J. J. Moreau. Proximité et dualité dans un espace hilbertien. *Bulletin de la Société mathématique de France*, 93:273–299, 1965.

- [35] Y. Nesterov and B. T. Polyak. Cubic regularization of newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.
- [36] Y. E. Nesterov. A method for solving the convex programming problem with convergence rate $O(1/k^2)$. *Doklady Akademii Nauk SSSR*, 269(3):543–547, 1983.
- [37] M. Nikolova. Analysis of the recovery of edges in images and signals by minimizing nonconvex regularized least-squares. *Multiscale Modeling & Simulation*, 4(3):960–991, 2005.
- [38] P. Ochs. *Long term motion analysis for object level grouping and nonsmooth optimization methods*. PhD thesis, Albert-Ludwigs-Universität Freiburg, Mar 2015.
- [39] P. Ochs. Local convergence of the heavy-ball method and ipiano for non-convex optimization. *Journal of Optimization Theory and Applications*, 177(1):153–180, 2018.
- [40] P. Ochs. Unifying abstract inexact convergence theorems and block coordinate variable metric ipiano. *SIAM Journal on Optimization*, 29(1):541–570, 2019.
- [41] P. Ochs, Y. Chen, T. Brox, and T. Pock. iPiano: inertial proximal algorithm for nonconvex optimization. *SIAM Journal on Imaging Sciences*, 7(2):1388–1419, 2014.
- [42] P. Ochs, J. Fadili, and T. Brox. Non-smooth non-convex bregman minimization: Unification and new algorithms. *Journal of Optimization Theory and Applications*, 181(1):244–278, 2019.
- [43] T. Pock and S. Sabach. Inertial proximal alternating linearized minimization (iPALM) for non-convex and nonsmooth problems. *SIAM Journal on Imaging Sciences*, 9(4):1756–1787, 2016.
- [44] B. T. Polyak. Some methods of speeding up the convergence of iterative methods. *Akademija Nauk SSSR. Žurnal Vyčislitel'noj Matematiki i Matematičeskoj Fiziki*, 4:791–803, 1964.
- [45] R. T. Rockafellar and R. J.-B. Wets. *Variational Analysis*, volume 317 of *Fundamental Principles of Mathematical Sciences*. Springer-Verlag, Berlin, 1998.
- [46] C. Song and J. Liu. Inexact proximal cubic regularized Newton methods for convex optimization. *ArXiv preprint arXiv:1902.02388*, 2019.
- [47] J.-L. Starck, F. Murtagh, and J. Fadili. *Sparse image and signal processing: wavelets, curvelets, morphological diversity*. Cambridge University Press, 2010.
- [48] W. Su, S. Boyd, and E. Candes. A differential equation for modeling Nesterovs accelerated gradient method: Theory and insights. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2510–2518. Curran Associates, Inc., 2014.
- [49] M. Teboulle. Entropic proximal mappings with application to nonlinear programming. *Mathematics of Operations Research*, 17(3):670–690, 1992.
- [50] M. Teboulle. A simplified view of first order methods for optimization. *Mathematical Programming*, 170(1):67–96, 2018.
- [51] N. Tripuraneni, M. Stern, C. Jin, J. Regier, and M. I. Jordan. Stochastic cubic regularization for fast nonconvex optimization. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 2899–2908. Curran Associates, Inc., 2018.
- [52] G. Wang, G. B. Giannakis, and Y. C. Eldar. Solving systems of random quadratic equations via truncated amplitude flow. *IEEE Transactions on Information Theory*, 64(2):773–794, 2018.

- [53] Z. Wang, Y. Zhou, Y. Liang, and G. Lan. Cubic regularization with momentum for nonconvex optimization. *ArXiv preprint arXiv:1810.03763*, 2018.
- [54] B. Wen, X. Chen, and T. K. Pong. Linear convergence of proximal gradient algorithm with extrapolation for a class of nonconvex nonsmooth minimization problems. *SIAM Journal on Optimization*, 27(1):124–145, 2017.
- [55] F. Wenand, L. Chu, P. Liu, and R. C. Qiu. Nonconvex regularization based sparse and low-rank recovery in signal processing, statistics, and machine learning. *ArXiv preprint arXiv:1808.05403*, 2018.
- [56] Y. Xu, Z. Li, J. Yang, and D. Zhang. A survey of dictionary learning algorithms for face recognition. *IEEE access*, 5:8502–8514, 2017.
- [57] J. Zhang, L. Xiao, and S. Zhang. Adaptive stochastic variance reduction for subsampled Newton method with cubic regularization. *ArXiv preprint arXiv:1811.11637*, 2018.