# A block inertial Bregman proximal algorithm for nonsmooth nonconvex problems with application to symmetric nonnegative matrix tri-factorization

MASOUD AHOOKHOSH[1], LE THI KHANH HIEN[2]. NICOLAS GILLIS[2], AND PANAGIOTIS PATRINOS[1]

ABSTRACT. We propose BIBPA, a block inertial Bregman proximal algorithm for minimizing the sum of a block relatively smooth function (that is, relatively smooth concerning each block) and block separable nonsmooth nonconvex functions. We prove that the sequence generated by BIBPA subsequentially converges to critical points of the objective under standard assumptions, and globally converges when the objective function is additionally assumed to satisfy the Kurdyka-Łojasiewicz (KŁ) property. We also provide the convergence rate when the objective satisfies the Łojasiewicz inequality. We apply BIBPA to the symmetric nonnegative matrix tri-factorization (SymTriNMF) problem, where we propose kernel functions for SymTriNMF and provide closed-form solutions for subproblems of BIBPA.

## 1. INTRODUCTION

This paper is concerned with the minimization of the sum of a block relatively smooth (see Definition 2.2), and a block separable (nonsmooth) nonconvex function. Although this problem has a simple structure, it covers a broad range of optimization problems arising in signal and image processing, machine learning, and inverse problems. In our block-structured nonconvex setting, the most common class of methodologies is *first-order* ones, where the central to their convergence analysis is the so-called *descent lemma* in both the Euclidean setting (e.g., [1, 11, 12, 36, 38]) and the non-Euclidean one (e.g., [10, 35, 48]). While for the Euclidean case, the descent lemma is guaranteed if the function has Lipschitz continuous gradients, in the non-Euclidean setting it holds for *relatively smooth* functions encompassing the class of smooth functions with Lipschitz gradients.

In the *Euclidean setting*, there are large number alternating minimization algorithms for handling our structured problem such as block coordinate methods [13, 14, 37, 46, 47] and Gauss-Seidel methods [8, 15, 28], proximal alternating minimization [5, 7], and proximal alternating linearized minimization [19, 41, 43]. In the *non-Euclidean setting*, several algorithms have been proposed, namely, Bregman forward-backward splitting [3, 9, 10, 20, 35, 45], accelerated Bregman forward-backward splitting [29, 31], stochastic mirror descent methods [30], Bregman proximal alternating linearized minimization [2].

In order to establish the global convergence of generic algorithms for (nonsmooth) non-convex problems, one needs to assume that the celebrated Kurdyka-Łojasiewicz inequality (see Definition 3.10) is satisfied as a feature of the underlying problem's class. The earliest abstract convergence theorem was introduced by Attouch et al. [6] and by Bolte et al. [19], relying on the following conditions that an algorithm should satisfy: (i) *sufficient decrease condition of the cost function*; (ii) *subgradient lower bound of iterations gap*; (iii) *subsequential convergence*. These conditions are shown to be satisfied by many algorithms [6]. In [19], these conditions were extended for proximal alternating linearized minimization. In the case of inertial proximal point algorithms [40, 41], it was shown that some Lyapunov function satisfies the sufficient decrease condition, which leads to a generalization of the abstract convergence theorem. A generalization of this theorem was introduced for variable metric algorithms in [26], which has been recently extended for inertial variable metric algorithms [39]. In this paper, we show that the results of [39] can cover the global convergence of algorithms in non-Euclidean settings.

1.1. **Contribution.** Our contribution is twofold:

1) (*Block inertial Bregman proximal algorithm*) We introduce `BIBPA`, a block generalization of the Bregman proximal gradient method [19] *with an inertial force*. We extend the notion of relative smoothness [10, 35, 48] to its block version (with different kernels for each block) to support our structured nonconvex problems. Notably, these kernel functions are block-wise convex, a property that does not necessarily imply their joint convexity for all blocks. Unlike the global convergence theorem in [6, 19] that verifies the *sufficient decrease condition* and *subgradient lower bound of iterations gap* on the cost function, for `BIBPA` these properties hold for a *Lyapunov function* including Bregman terms (see the equation (3.8)). Then, the global convergence of `BIBPA` is studied under the KŁ property, and its convergence rate is studied for Łojasiewicz-type KŁ functions.

2) (*Globally convergent scheme for solving the SymTriNMF problem*) With appropriate selection of kernel functions for Bregman distances, it turns out that the objective of the *symmetric nonnegative matrix tri-factorization* (SymTriNMF) problem is block relatively smooth, and the corresponding subproblems can be solved in closed form, an important property when dealing with machine learning problems that include a large number of variables. To the best of our knowledge, `BIBPA` is the first scheme with a rigorous theoretical guarantee of convergence for the SymTriNMF problem.

1.2. **Related works.** There are three papers [2, 50, 51] that are closely related to this paper. In [2], we introduced a multi-block relative smoothness condition that exploits a single kernel function for all blocks, while in the current paper we assume a block relative smoothness condition allowing a different kernel function for each block. Moreover, our algorithm `BIBPA` involves dynamic step-sizes and inertial terms for each block that makes our derivation and analysis different from those of [2]. Beside of the algorithmic differences with [50], we use nonseparable (nonconvex) kernels as apposed to the separable convex kernel used in [50] for each block. An inertial Bregman proximal gradient algorithm was presented in [51] for composite minimization that does not support our block structure nonconvex problems and therefore is different in derivation and analysis concerning our work.

1.3. **Organization.** The remainder of this paper is organized as follows. While Section 2 discusses the problem statement and the block relative smoothness, Section 3 introduces and analyzes a block inertial Bregman proximal algorithm (`BIBPA`). In Section 4, it is shown the `BIBPA`'s subproblems are solved in closed form. Some conclusion are delivered in Section 5.

## 2. Problem statement and block relative smoothness

We consider the structured nonsmooth nonconvex minimization problem

$$\underset{x \in \overline{C}}{\textbf{minimize}} \quad \Phi(\boldsymbol{x}) \equiv f(\boldsymbol{x}) + \sum_{i=1}^{N} g_i(x_i), \tag{2.1}$$

where $C$ is a nonempty, convex, and open set in $\mathbb{R}^n$ and $\overline{C}$ denotes its closure. Setting $n = \sum_{i=1}^{N} n_i$ and $i = 1, \ldots, N$, we assume the following hypotheses:

**Assumption I** (requirements for composite minimization (2.1))**.**

A1 $g_i : \mathbb{R}^{n_i} \to \overline{\mathbb{R}} := \mathbb{R} \cup \{\infty\}$ *is proper and lower semicontinuous (lsc);*

A2 $h_i : \mathbb{R}^n \to \overline{\mathbb{R}}$ *is $i$-th block Legandre,* $\textbf{int dom}\, h_1 = \ldots = \textbf{int dom}\, h_N$, $\overline{C} \subseteq \overline{\textbf{dom}\, h_1}$, *and* $\textbf{dom}\, g \cap C \neq \emptyset$ *with* $g := \sum_{i=1}^{N} g_i$;

A3 $f : \mathbb{R}^n \to \overline{\mathbb{R}}$ *is* $C^1(\textbf{int dom}\, h_1)$ *and* $(L_1, \ldots, L_N)$*-smooth relative to* $(h_1, \ldots, h_N)$;

A4 $\textbf{arg min}\, \big\{\Phi(x) \mid \boldsymbol{x} \in \overline{C}\big\} \neq \emptyset.$

2.1. **Notation.** We denote by $\overline{\mathbb{R}} := \mathbb{R} \cup \{\infty\}$ the extended-real line. We use boldface lower-case letters (e.g., $\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{z}$) for vectors in $\mathbb{R}^n$ and use normal lower-case letters (e.g., $z_i$, $x_i$, $y_i$) for vectors in $\mathbb{R}^{n_i}$, for $n_i \in \mathbb{N}$. For the identity matrix $I_n$, we set $U_i \in \mathbb{R}^{n \times n_i}$ such that $I_n = (U_1, \ldots, U_N) \in \mathbb{R}^{n \times n}$. The set of cluster points of $(\boldsymbol{x}^k)_{k \in \mathbb{N}}$ is denoted as $\omega(\boldsymbol{x}^0)$. A function $f : \mathbb{R}^n \to \overline{\mathbb{R}}$ is *proper* if $f > -\infty$ and $f \not\equiv \infty$, in which case its *domain* is defined as the set $\textbf{dom}\, f := \{\boldsymbol{x} \in \mathbb{R}^n \mid f(\boldsymbol{x}) < \infty\}$. A vector $\boldsymbol{v} \in \partial f(\boldsymbol{x})$ is a *subgradient* of $f$ at $\boldsymbol{x}$, and the set of all such vectors is called the *subdifferential* $\partial f(\boldsymbol{x})$ [42, Definition 8.3], i.e.

$$\partial f(\boldsymbol{x}) = \Big\{\boldsymbol{v} \in \mathbb{R}^n \mid \exists (\boldsymbol{x}^k, \boldsymbol{v}^k)_{k \in \mathbb{N}} \text{ s.t. } \boldsymbol{x}^k \to \boldsymbol{x}, \ f(\boldsymbol{x}^k) \to f(\boldsymbol{x}), \ \widehat{\partial} f(\boldsymbol{x}^k) \ni \boldsymbol{v}^k \to \boldsymbol{v}\Big\},$$

where $\widehat{\partial} f(\boldsymbol{x})$ is the set of *regular subgradients* of $f$ at $\boldsymbol{x}$, namely

$$\widehat{\partial} f(\boldsymbol{x}) = \Big\{\boldsymbol{v} \in \mathbb{R}^n \mid f(\boldsymbol{z}) \geq f(\boldsymbol{x}) + \langle \boldsymbol{v}, \boldsymbol{z} - \boldsymbol{x}\rangle + o(\|\boldsymbol{z} - \boldsymbol{x}\|), \ \forall \boldsymbol{z} \in \mathbb{R}^n\Big\}.$$

2.2. **Block relative smoothness.** We first describe the notion of *block relative smoothness*, which is an extension of the relative smoothness [10, 35]. To this end, we introduce the notion of *block kernel* functions, which coincides with the classical one (cf. [3, Definition 2.1]) for $N = 1$.

**Definition 2.1** (*$i$-th block convexity and kernel function*)**.** *Let* $h : \mathbb{R}^n \to \overline{\mathbb{R}}$ *be a proper and lower semicontinuous (lsc) function with* $\textbf{int dom}\, h \neq \emptyset$ *and such that* $h \in C^1(\textbf{int dom}\, h)$. *For a fixed vector* $\boldsymbol{x} \in \mathbb{R}^n$ *and* $i \in \{1, \ldots, N\}$, *we say that $h$ is*

  *(i)* *$i$-th block (strongly/strictly) convex if the function* $h(\boldsymbol{x} + U_i(\cdot - x_i))$ *is (strongly/strictly) convex for all* $\boldsymbol{x} \in \textbf{dom}\, h$;

 *(ii)* *a $i$-th block kernel function if $h$ is $i$-th block convex and* $h(\boldsymbol{x} + U_i(\cdot - x_i))$ *is 1-coercive for all* $\boldsymbol{x} \in \textbf{dom}\, h$, *i.e.,* $\lim_{\|z\| \to \infty} \frac{h(\boldsymbol{x} + U_i(z - x_i))}{\|z\|} = \infty$;

*(iii)* *$i$-th block essentially smooth, if for every sequence* $(\boldsymbol{x}^k)_{k \in \mathbb{N}} \subseteq \textbf{int dom}\, h$ *converging to a boundary point of* $\textbf{dom}\, h$, *we have* $\|\nabla_i h(\boldsymbol{x}^k)\| \to \infty$;

*(iv)* *$i$-th block Legendre if it is $i$-th block essentially smooth and $i$-th block strictly convex.*

Let $h : \mathbb{R}^n \to \overline{\mathbb{R}}$ be a Legendre function. Then, the classical definition of *Bregman distances* (cf. [23]) leads to the function $\mathbf{D}_h : \mathbb{R}^n \times \mathbb{R}^n \to \overline{\mathbb{R}}$ given by

$$\mathbf{D}_h(\boldsymbol{y}, \boldsymbol{x}) := \begin{cases} h(\boldsymbol{y}) - h(\boldsymbol{x}) - \langle \nabla h(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x}\rangle & \text{if } \boldsymbol{y} \in \textbf{dom}\, h, \boldsymbol{x} \in \textbf{int dom}\, h, \\ \infty & \text{otherwise.} \end{cases} \tag{2.2}$$

However, in the remainder of this paper, we extend this definition for the cases that $h$ is only an $i$-th block Legendre function. Fixing all blocks except the $i$-th one, the Bregman distance (2.2) will reduce to $\mathbf{D}_h(\boldsymbol{x}+U_i(y_i-x_i),\boldsymbol{x}) = h(\boldsymbol{x}+U_i(y_i-x_i))-h(\boldsymbol{x})-\langle\nabla_i h(\boldsymbol{x}), y_i-x_i\rangle$, which measures the proximity between $\boldsymbol{x}+U_i(y_i-x_i)$ and $\boldsymbol{x}$ with respect to the $i$-th block of variables. Moreover, the kernel $h$ is $i$-th block convex if and only if $\mathbf{D}_h(\boldsymbol{x}+U_i(y_i-x_i),\boldsymbol{x}) \geq 0$ for all $\boldsymbol{x}+U_i(y_i-x_i) \in \mathbf{dom}\,h$ and $\boldsymbol{x} \in \mathbf{int\,dom}\,h$. Note that if $h$ is $i$-th block strictly convex, then $\mathbf{D}_h(\boldsymbol{x}+U_i(y_i-x_i),\boldsymbol{x}) = 0$ if and only if $x_i = y_i$.

We are now in a position to present the notion of *block relative smoothness*, which is the central tool for our analysis in the next section.

**Definition 2.2** (block relative smoothness). *For $i \in [N]$, let $h_i : \mathbb{R}^n \to \overline{\mathbb{R}}$ be $i$-th block kernel functions and let $f : \mathbb{R}^n \to \overline{\mathbb{R}}$ be a proper and lsc function. If there exists $L_i > 0$, $i \in [N]$, such that $L_i h_i(\boldsymbol{x}+U_i(z-x_i)) - f(\boldsymbol{x}+U_i(z-x_i))$ are convex for all $\boldsymbol{x}, \boldsymbol{x}+U_i(z-x_i) \in \mathbf{int\,dom}\,h_i$, then, $f$ is called $(L_1,\ldots,L_N)$-smooth relative to $(h_1,\ldots,h_N)$.*

Note that if $N = 1$, the block relative smoothness is reduced to standard relative smoothness [10, 35]. If $f$ is $L$-Lipschitz continuous, then both $L/2\|\cdot\|^2 - f$ and $L/2\|\cdot\|^2 + f$ are convex, i.e., the relative smoothness of $f$ generalizes the notions of Lipschitz continuity.

**Proposition 2.3** (characterization of block relative smoothness). *For $i = 1,\ldots,N$, let $h_i : \mathbb{R}^n \to \overline{\mathbb{R}}$ be $i$-th block kernels and let $f : \mathbb{R}^n \to \overline{\mathbb{R}}$ be a proper lsc function and $f \in C^1$. Then, the following statements are equivalent:*

*(a) $(L_1,\ldots,L_N)$-smooth relative to $(h_1,\ldots,h_N)$;*

*(b) for all $(\boldsymbol{x},\boldsymbol{y}) \in \mathbf{int\,dom}\,h_i \times \mathbf{int\,dom}\,h_i$ and $i = 1,\ldots,N$,*

$$f(\boldsymbol{x}+U_i(y_i-x_i)) \leq f(\boldsymbol{x}) + \langle\nabla_i f(\boldsymbol{x}), y_i-x_i\rangle + L_i\,\mathbf{D}_{h_i}(\boldsymbol{x}+U_i(y_i-x_i),\boldsymbol{x}); \quad (2.3)$$

*(c) for all $(\boldsymbol{x},\boldsymbol{y}) \in \mathbf{int\,dom}\,h_i \times \mathbf{int\,dom}\,h_i$ and $i = 1,\ldots,N$,*

$$\langle\nabla_i f(\boldsymbol{x}) - \nabla_i f(\boldsymbol{y}), x_i - y_i\rangle \leq L_i\langle\nabla_i h_i(\boldsymbol{x}) - \nabla_i h_i(\boldsymbol{y}), x_i - y_i\rangle; \quad (2.4)$$

*(d) if $f \in C^2(\mathbf{int\,dom}\,f)$ and $h \in C^2(\mathbf{int\,dom}\,h_i)$, then*

$$L_i\nabla^2_{x_i x_i}h_i(\boldsymbol{x}) - \nabla^2_{x_i x_i}f(\boldsymbol{x}) \geq 0, \quad \forall\boldsymbol{x} \in \mathbf{int\,dom}\,h_i, \ i = 1,\ldots,N. \quad (2.5)$$

*Proof.* The proof is a straightforward extension of those given in [35, Proposition 1.1], by fixing all the blocks except one of them. $\square$

### 2.3. **Motivating example: symmetric nonnegative matrix tri-factorization.**
We consider a symmetric matrix $X \in \mathbb{R}^{m\times m}$ and aim to decompose it in the form $X = UVU^T$, where $U \in \mathbb{R}^{m\times r}_+$ and $V \in \mathbb{R}^{r\times r}_+$. This translates to the minimization of $\frac{1}{2}\|X - UVU^T\|^2_F$ for $U, V \geq 0$, leading to the unconstrained problem

$$\min_{U\in\mathbb{R}^{m\times r},V\in\mathbb{R}^{r\times r}} \frac{1}{2}\|X - UVU^T\|^2_F + \delta_{U\geq 0} + \delta_{V\geq 0}. \quad (2.6)$$

**Proposition 2.4** (block relative smoothness of SymTriNMF objective). *Let functions $h_1 : \mathbb{R}^{m\times r} \times \mathbb{R}^{r\times r} \to \overline{\mathbb{R}}$ and $h_2 : \mathbb{R}^{m\times r} \times \mathbb{R}^{r\times r} \to \overline{\mathbb{R}}$ be strongly convex kernel functions as*

$$h_1(U,V) := \tfrac{a_1}{4}\|V\|^2_F\|U\|^4_F + \tfrac{b_1}{2}(\|X\|_F\|V\|_F + \varepsilon_1)\|U\|^2_F, \quad (2.7)$$

$$h_2(U,V) := \tfrac{a_2}{2}\big(\|U\|^4_F + \varepsilon_2\big)\|V\|^2_F. \quad (2.8)$$

*with $\varepsilon_1,\varepsilon_2 > 0$. Then the function $f : \mathbb{R}^{m\times r}\times\mathbb{R}^{r\times r} \to \overline{\mathbb{R}}$ given by $f(U,V) := \frac{1}{2}\|X-UVU^T\|^2_F$ is $(L_1,L_2)$-smooth relative to $(h_1,h_2)$ with*

$$L_1 \geq \mathbf{max}\left\{\tfrac{6}{a_1}, \tfrac{2}{b_1}\right\}, \quad and \quad L_2 \geq \tfrac{1}{a_2}. \quad (2.9)$$

*Proof.* Plugging the partial derivative $\nabla_U f(U, V) = -XUV^T - X^T UV + UVU^T UV^T + UV^T U^T UV$ into the definition of directional derivative, we obtain

$$\nabla^2_{UU} f(U, V)Z = -2XZV^T + UVU^T ZV + UVZ^T UV + ZVU^T UV^T$$
$$+ UV^T U^T ZV + UV^T Z^T UV + ZV^T U^T UV,$$

which consequently leads to

$$\langle Z, \nabla^2_{UU} f(U, V)Z \rangle = -2\langle Z, XZV^T \rangle + \langle Z, UVU^T ZV \rangle + \langle Z, UVZ^T UV \rangle + \langle Z, ZVU^T UV^T \rangle$$
$$+ \langle Z, UV^T U^T ZV \rangle + \langle Z, UV^T Z^T UV \rangle + \langle Z, ZV^T U^T UV \rangle$$
$$\leq \left(2\|X\| \|V\| + 6\|U\|^2 \|V\|^2\right)\|Z\|^2_F.$$

On the other hand, from $\nabla_U h_1(U, V) = \left(a_1 \|U\|^2_F \|V\|^2_F + b_1(\|X\|_F \|V\|_F + \varepsilon_1)\right)U$, we have

$$\nabla^2_{UU} h_1(U, V)Z = \left(2a_1 \|V\|^2 \langle U, Z \rangle\right)U + \left(a_1 \|V\|^2_F \|U\|^2_F + b_1(\|X\|_F \|V\|_F + \varepsilon_1)\right)Z,$$

implying that

$$\langle Z, \nabla^2_{UU} h_1(U, V)Z \rangle \geq \left(a_1 \|V\|^2_F \|U\|^2_F + b_1(\|X\|_F \|V\|_F + \varepsilon_1)\right)\|Z\|^2_F$$
$$\geq \left(a_1 \|V\|^2 \|U\|^2 + b_1(\|X\| \|V\| + \varepsilon_1)\right)\|Z\|^2_F.$$

Therefore, the inequality

$$\langle Z, (L_1 \nabla^2_{UU} h_1(U, V) - \nabla^2_{UU} f(U, V))Z \rangle$$
$$\geq \left((L_1 a_1 - 6)\|V\|^2 \|U\|^2 + (L_1 b_1 - 2)\|X\| \|V\| + \varepsilon_1 L_1\right)\|Z\|^2_F \geq 0$$

holds if $L_1 a_1 - 6 \geq 0$ and $L_1 b_1 - 2 \geq 0$, as claimed.

It follows from $\nabla_V f(U, V) = U^T XU + U^T UVU^T U$ that

$$\nabla^2_{VV} f(U, V)Z = \lim \frac{\nabla_U f(U + tZ, V) - \nabla_U f(U, V)}{t} = U^T UZU^T U,$$

leading to the inequality $\langle Z, \nabla^2_{VV} f(U, V)Z \rangle = \langle Z, U^T UZU^T U \rangle \leq \|U\|^4 \|Z\|^2_F$. Now, using $\nabla_V h_2(U, V) = a_2\left(\|U\|^4 + \varepsilon_2\right)V$, we get $\langle Z, \nabla^2_{VV} h_2(U, V)Z \rangle = a_2\left(\|U\|^4 + \varepsilon_2\right)\|Z\|^2_F$, i.e.,

$$\langle Z, (L_2 \nabla^2_{VV} h_2(U, V) - \nabla^2_{VV} f(U, V))Z \rangle = \left((L_2 a_2 - 1)\|U\|^4 + \varepsilon_2 L_2\right)\|Z\|^2_F \geq 0$$

if $L_2 a_2 - 1 \geq 0$, giving our desired results. $\qquad\square$

## 3. Block inertial Bregman proximal algorithm

This section discusses our algorithm, starting from the prox-boundedness extension [42].

**Definition 3.1** (block prox-boundedness). *A function* $g : \mathbb{R}^n \to \overline{\mathbb{R}}$ *is* block prox-bounded *if for each* $i \in \{1, \ldots, N\}$ *there exists* $\gamma_i > 0$ *and* $x \in \mathbb{R}^n$ *such that*

$$\inf_{z \in \mathbb{R}^{n_i}} \left\{ g(x + U_i(z - x_i)) + \tfrac{1}{\gamma_i} \mathbf{D}_{h_i}(x + U_i(z - x_i), x) \right\} > -\infty.$$

*The supremum of the set of all such* $\gamma_i$ *is the threshold* $\gamma^h_{i,g}$ *of the block prox-boundedness,*

$$\gamma^{h_i}_{i,g} := \sup_{\gamma_i > 0} \left\{ \gamma_i : \exists x \in \mathbb{R}^n, \inf_{z \in \mathbb{R}^{n_i}} \left\{ g(x + U_i(z - x_i)) + \tfrac{1}{\gamma_i} \mathbf{D}_{h_i}(x + U_i(z - x_i), x) \right\} > -\infty \right\}.$$
(3.1)

**Proposition 3.2** (characteristics of block prox-boundedness). *For* $h_i : \mathbb{R}^n \to \overline{\mathbb{R}}$ *and proper and lsc functions* $g_i : \mathbb{R}^{n_i} \to \overline{\mathbb{R}}$ $(i = 1, \ldots, N)$, *the following statements are equivalent:*

*(a)* $g = \sum_{i=1}^N g_i$ *is block prox-bounded;*

*(b)* *for all* $i = 1, \ldots, N$, $g_i + r_i h_i(x + U_i(z - x_i))$ *is bounded below on* $\mathbb{R}^{n_i}$ *for some* $r_i \in \mathbb{R}$;

*(c) for all $i = 1, \ldots, N$, $\liminf_{\|z\| \to \infty} g_i(z)/h_i(x + U_i(z - x_i)) > -\infty$.*

*Proof.* The proof is a straightforward adaptation of [2, Proposition 2.7]. $\qquad\square$

For a given points $x^k, x^{k-1} \in \mathbb{R}^n$ and $\alpha_i^k \geq 0$, let us define the function $\mathcal{M}_{h_i/\gamma_i^k} : \mathbf{dom}\, h_i \times \mathbf{int\, dom}\, h_i \times \mathbf{int\, dom}\, h_i \to \overline{\mathbb{R}}$ given by

$$\mathcal{M}_{h_i/\gamma_i^k}(x, x^k, x^{k-1}) := \langle \nabla f(x^k) - \tfrac{\alpha_i^k}{\gamma_i^k}(x^k - x^{k-1}), x - x^k \rangle + \tfrac{1}{\gamma_i^k} \mathbf{D}_{h_i}(x, x^k) + \sum_{i=1}^{N} g_i(x_i) \quad (3.2)$$

and the *block inertial Bregman proximal* mapping $\mathbf{T}_{h_i/\gamma_i^k} : \mathbf{int\, dom}\, h_i \times \mathbf{int\, dom}\, h_i \rightrightarrows \mathbb{R}^{n_i}$ as

$$\begin{aligned}
\mathbf{T}_{h_i/\gamma_i^k}(x^k, x^{k-1}) &:= \underset{z \in \mathbb{R}^{n_i}}{\arg\min}\ \mathcal{M}_{h_i/\gamma_i^k}(x^k + U_i(z - x_i^k), x^k, x^{k-1}) \\
&= \underset{z \in \mathbb{R}^{n_i}}{\arg\min}\ \langle \nabla_i f(x^k) - \tfrac{\alpha_i^k}{\gamma_i^k}(x_i^k - x_i^{k-1}), z - x_i^k \rangle + \tfrac{1}{\gamma_i^k} \mathbf{D}_{h_i}(x^k + U_i(z - x_i^k), x^k) + g_i(z),
\end{aligned} \quad (3.3)$$

which is set-valued by nonconvexity of $g_i$ $(i = 1, \ldots, N)$, and it reduces to the inertial Bregman forward-backward mapping for $N = 1$; cf. [21]. For a given sequence $(x^k)_{k \in \mathbb{N}}$, we introduce the following notation

$$x^{k,i} := (x_1^{k+1}, \ldots, x_i^{k+1}, x_{i+1}^k, \ldots, x_N^k), \quad (3.4)$$

i.e., $x^{k,0} = x^k$ and $x^{k,N} = x^{k+1}$. Using this notation and the mapping (3.3), we next introduce the *block inertial Bregman proximal algorithm* (BIBPA); see Algorithm 1.

---

**Algorithm 1 (BIBPA)** Block Inertial Bregman Proximal Algorithm

---

INPUT $\quad x^0 \in \mathbf{int\, dom}\, h_1$, $I_n = (U_1, \ldots, U_N) \in \mathbb{R}^{n \times n}$ with $U_i \in \mathbb{R}^{n \times n_i}$ and the identity matrix $I_n$, $k = 0$.

1: **while** some stopping criterion is not met **do**

2: $\quad x^{k,0} = x^k$;

3: $\quad$ **for** $i = 1, \ldots, N$ **do** choose $\gamma_i^k$ and $\alpha_i^k$ as Prop. 3.5 and compute

$$x_i^{k,i} \in \mathbf{T}_{h_i/\gamma_i^k}(x^{k,i-1}, x^{k-1}), \quad x^{k,i} = x^{k,i-1} + U_i(x_i^{k,i} - x_i^{k,i-1}); \quad (3.5)$$

4: $\quad x^{k+1} = x^{k,N}$, $k = k + 1$;

OUTPUT $\quad$ A vector $x^k$.

---

In order to verify the well-definedness of the iterations generated by BIBPA, we next investigate some important properties of the mapping $\mathbf{T}_{h_i/\gamma_i^k}$.

**Assumption II.** *For all $z \in \mathbf{T}_{h_i/\gamma_i^k}(x, y)$ and $\gamma_i^k \in (0, 1/L_i)$, $x + U_i(z - x_i) \in C$ and $i = 1, \ldots, N$.*

**Proposition 3.3** (properties of the mapping $\mathbf{T}_{h_i/\gamma_i^k}$)**.** *Under Assumption I and Assumption II, $\gamma_i^k \in (0, \gamma_{i,g}^{h_i})$ for $i \in [N]$, and $x^k, x^{k-1} \in \mathbf{int\, dom}\, h_i$, the following statements are true:*

*(i) $\mathbf{T}_{h_i/\gamma_i^k}(x^k, x^{k-1})$ is nonempty, compact, and outer semicontinuous;*

*(ii) $\mathbf{dom}\, \mathbf{T}_{h_i/\gamma_i^k} = \mathbf{int\, dom}\, h_i \times \mathbf{int\, dom}\, h_i$;*

*(iii) If $x_i^{k,i} \in \mathbf{T}_{h_i/\gamma_i^k}(x^{k,i-1}, x^{k-1})$ for $\gamma_i^k \in (0, 1/L_i)$, then $x^{k,i} \in \mathbf{int\, dom}\, h_i$.*

*Proof.* The proof follows from [2, Proposition 2.10] and Assumption II. $\qquad\square$

In the subsequent lemma, we show that the cost function $\Phi$ satisfies some necessary inequality that will be needed in the next result.

**Lemma 3.4** (cyclic inequality of the cost). *Let Assumption I and Assumption II hold, and let $(\mathbf{x}^k)_{k\in\mathbb{N}}$ be generated by BIBPA. If $h_i$ ($i \in [N]$) is $\sigma_i$-block strongly convex, then we have*

$$\Phi(\mathbf{x}^{k+1}) - \Phi(\mathbf{x}^k) \le \sum_{i=1}^{N} \left( \left( \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} - \frac{1-\gamma_i^k L_i}{\gamma_i^k} \right) \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) + \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} \mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1}) \right). \quad (3.6)$$

*Proof.* For $i \in \{1, \ldots, N\}$ and $x_i^{k,i} \in \mathbf{T}_{h_i/\gamma_i^k}(\mathbf{x}^{k,i-1}, \mathbf{x}^{k-1})$, it holds that

$$\langle \nabla_i f(\mathbf{x}^{k,i-1}) - \tfrac{\alpha_i^k}{\gamma_i^k}(x_i^k - x_i^{k-1}), x_i^{k,i} - x_i^k \rangle + \tfrac{1}{\gamma_i^k} \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) + \sum_{j=1}^{N} g_j(x_j^{k,i}) \le \sum_{j=1}^{N} g_j(x_j^{k,i-1}).$$

Together with Assumption IA3 and Proposition 2.3(b), this implies

$$f(\mathbf{x}^{k,i}) \le f(\mathbf{x}^{k,i-1}) + \langle \nabla_i f(\mathbf{x}^{k,i-1}), x_i^{k,i} - x_i^k \rangle + L_i \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1})$$
$$\le f(\mathbf{x}^{k,i-1}) + \sum_{j=1}^{N} g_j(x_j^{k,i-1}) - \sum_{j=1}^{N} g_j(x_j^{k,i}) - \tfrac{1-\gamma_i^k L_i}{\gamma_i^k} \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) + \tfrac{\alpha_i^k}{\gamma_i^k}\langle x_i^k - x_i^{k-1}, x_i^{k,i} - x_i^k \rangle$$
$$\le f(\mathbf{x}^{k,i-1}) + \sum_{j=1}^{N} g_j(x_j^{k,i-1}) - \sum_{j=1}^{N} g_j(x_j^{k,i}) - \tfrac{1-\gamma_i^k L_i}{\gamma_i^k} \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1})$$
$$+ \tfrac{|\alpha_i^k|}{2\gamma_i^k}\left( \|x_i^k - x_i^{k-1}\|^2 + \|x_i^{k,i} - x_i^k\|^2 \right)$$
$$\le f(\mathbf{x}^{k,i-1}) + \sum_{j=1}^{N} g_j(x_j^{k,i-1}) - \sum_{j=1}^{N} g_j(x_j^{k,i}) + \left( \tfrac{|\alpha_i^k|}{\sigma_i \gamma_i^k} - \tfrac{1-\gamma_i^k L_i}{\gamma_i^k} \right) \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1})$$
$$+ \tfrac{|\alpha_i^k|}{\sigma_i \gamma_i^k} \mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1}),$$

which yields

$$\Phi(\mathbf{x}^{k,i}) \le \Phi(\mathbf{x}^{k,i-1}) + \left( \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} - \frac{1-\gamma_i^k L_i}{\gamma_i^k} \right) \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) + \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} \mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1}). \quad (3.7)$$

Now, let us sum up both sides of (3.7) for $i = 1, \ldots, N$, i.e.,

$$\Phi(\mathbf{x}^{k+1}) - \Phi(\mathbf{x}^k) = \sum_{i=1}^{N} \left( \Phi(\mathbf{x}^{k,i}) - \Phi(\mathbf{x}^{k,i-1}) \right)$$
$$\le \sum_{i=1}^{N} \left( \left( \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} - \frac{1-\gamma_i^k L_i}{\gamma_i^k} \right) \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) + \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} \mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1}) \right).$$

$\square$

We notice that Lemma 3.4 does not guarantee the monotonicity of the sequence $(\Phi(\mathbf{x}^k))_{k\in\mathbb{N}}$. For $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and $\delta_i \ge 0$, we define the *Lyapunov function* $\mathcal{L} : \mathbb{R}^n \times \mathbb{R}^n \to \overline{\mathbb{R}}$,

$$\mathcal{L}(\mathbf{x}, \mathbf{y}) := \Phi(\mathbf{x}) + \sum_{i=1}^{N} \delta_i \mathbf{D}_{h_i}((x_1, \ldots, x_i, y_{i+1}, \ldots, y_N), (x_1, \ldots, x_{i-1}, y_i, \ldots, y_N)), \quad (3.8)$$

Note that $\mathcal{L}(\mathbf{x}^{k+1}, \mathbf{x}^k) := \Phi(\mathbf{x}^{k+1}) + \sum_{i=1}^{N} \delta_i \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1})$. We denote by $\mathcal{L}^{k+1}$ and $\mathcal{L}^k$ the terms $\mathcal{L}(\mathbf{x}^{k+1}, \mathbf{x}^k)$ and $\mathcal{L}(\mathbf{x}^k, \mathbf{x}^{k-1})$, respectively. We next indicate the monotonicity of $(\mathcal{L}^k)_{k\in\mathbb{N}}$.

**Proposition 3.5** (descent property of the Lyapunov function). *Let Assumption I and Assumption II hold, let $(\mathbf{x}^k)_{k\in\mathbb{N}}$ be generated by BIBPA, and let $h_i$ ($i = 1, \ldots, N$) be $\sigma_i$-block strongly convex. If $\lim_{k\to\infty} \alpha_i^k = \alpha_i$ and $0 < \gamma_i \le \frac{\sigma_i - 2|\alpha_i|}{\sigma_i L_i}$ and*

$$|\alpha_i^k| < \frac{\sigma_i}{2}, \quad 0 < \gamma_i \le \gamma_i^k \le \frac{\sigma_i - 2|\alpha_i^k|}{\sigma_i L_i}, \quad \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} \le \delta_i \le \frac{1-\gamma_i^k L_i}{\gamma_i^k} - \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} \quad i = 1, \ldots, N, \quad (3.9)$$

*then, setting $a_i := \frac{1-\gamma_i^k L_i}{\gamma_i^k} - \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} - \delta_i$ and $b_i := \delta_i - \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k}$ for $i = 1, \ldots, N$, we get*

$$\mathcal{L}^{k+1} - \mathcal{L}^k \le -\sum_{i=1}^{N} \left( a_i \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) + b_i \mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1}) \right), \quad (3.10)$$

*i.e., the sequence* $(\mathcal{L}^k)_{k\in\mathbb{N}}$ *is non-increasing and consequently* $\lim_{k\to\infty} \mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1}) = 0$, *i.e.,* $\lim_{k\to\infty} \|x^{k,i} - x^{k,i-1}\| = 0$, *for all* $i = 1, \ldots, N$.

*Proof.* Using (3.6) and applying the Lyapunov function (3.8), we have

$$\mathcal{L}^{k+1} - \mathcal{L}^k = \Phi(x^{k+1}) - \Phi(x^k) + \sum_{i=1}^N \delta_i \mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1}) - \sum_{i=1}^N \delta_i \mathbf{D}_{h_i}(x^{k-1,i}, x^{k-1,i-1})$$

$$\leq \sum_{i=1}^N \left( \left( \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} - \frac{1-\gamma_i^k L_i}{\gamma_i^k} + \delta_i \right) \mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1}) + \left( \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} - \delta_i \right) \mathbf{D}_{h_i}(x^{k-1,i}, x^{k-1,i-1}) \right),$$

as claimed in (3.10). In order to guarantee the non-increasing property of the sequence $(\mathcal{L}^k)_{k\in\mathbb{N}}$, the inequalities $a_i = \frac{1-\gamma_i^k L_i}{\gamma_i^k} - \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} - \delta_i \geq 0$, $b_i = \delta_i - \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} \geq 0$ should be satisfied, for $i = 1, \ldots, N$, i.e., $\frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} \leq \delta_i \leq \frac{1-\gamma_i^k L_i}{\gamma_i^k} - \frac{|\alpha_i^k|}{\sigma_i \gamma_i^k} \leq \frac{1-\gamma_i L_i}{\gamma_i}$ $i = 1, \ldots, N$, which is guaranteed by (3.9), i.e., $\mathcal{L}^{k+1} \leq \mathcal{L}^k$. Together with (3.10), this yields that

$$\sum_{k=0}^p \sum_{i=1}^N a_i \mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1}) + b_i \mathbf{D}_{h_i}(x^{k-1,i}, x^{k-1,i-1}) \leq \sum_{k=0}^p \left( \mathcal{L}^k - \mathcal{L}^{k+1} \right)$$

$$= \mathcal{L}^0 - \mathcal{L}^{p+1} \leq \mathcal{L}^0 - \inf \mathcal{L} < +\infty.$$

Let $p \to +\infty$, the result follows from $\mathbf{D}_{h_i}(\cdot, \cdot) \geq 0$ and block strong convexity of $h_i$. $\qquad\square$

In convergence analysis of proximal algorithms, one usual assumption is the boundedness of $(x^k)_{k\in\mathbb{N}}$; cf., [5, 20]. A sufficient condition for this is given next.

**Corollary 3.6** (boundedness of iterations). *Suppose that all assumptions of Proposition 3.5 hold. Further, if* $\varphi$ *has bounded level sets, then the sequence* $(x^k)_{k\in\mathbb{N}}$ *is bounded.*

*Proof.* It follows from Proposition 3.5 that $\mathcal{L}(x^{k+1}, x^k)$ is non-increasing, hence

$$\Phi(x^{k+1}) \leq \mathcal{L}(x^{k+1}, x^k) = \Phi(x^{k+1}) + \sum_{i=1}^N \delta_i \mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1})$$

$$\leq \mathcal{L}(x^1, x^0) = \Phi(x^1) + \sum_{i=1}^N \delta_i \mathbf{D}_{h_i}(x^{0,i}, x^{0,i-1}) < \infty.$$

Hence, $\mathcal{N}(x^1, x^0) := \left\{ x \in \mathbb{R}^n \mid \Phi(x) \leq \Phi(x^1) + \sum_{i=1}^N \delta_i \mathbf{D}_{h_i}(x^{0,i}, x^{0,i-1}) \right\}$ encompasses $(x^k)_{k\in\mathbb{N}}$, i.e., $(x^k)_{k\in\mathbb{N}} \subseteq \mathcal{N}(x^1, x^0)$. Since $\varphi$ has bounded level sets, we have $(x^k)_{k\in\mathbb{N}}$ are bounded. $\quad\square$

The next proposition provides a lower bound for $\sum_{i=1}^N \sqrt{\mathbf{D}_h(x^{k,i}, x^{k,i-1})} + \sqrt{\mathbf{D}_h(x^{k-1,i}, x^{k-1,i-1})}$.

**Proposition 3.7** (subgradient lower bound for iterations gap). *Let Assumption I and Assumption II hold, let* $(x^k)_{k\in\mathbb{N}}$ *be generated by* BIBPA, *and let* $h_i$ $(i \in [N])$ *be* $\sigma_i$-*block strongly convex. Furthermore, suppose that* $\nabla_i f$, $\nabla_i h$, $(i = 1, \ldots, N)$ *are locally Lipschitz on bounded sets with Lipschitz moduli* $\widehat{L}$ *and* $\widetilde{L}_i > 0$, $\nabla_{ii}^2 h_i$ *is bounded on bounded set with constants* $\overline{L}_i$ $(i \in [N])$ *and that the sequence* $(x^k)_{k\in\mathbb{N}}$ *is bounded. For a fixed* $k \in \mathbb{N}$ *and* $j \in [N]$, *we define*

$$\mathcal{G}_j^{k+1} := (\mathcal{V}_j^{k+1}, \mathcal{W}_j^{k+1}), \tag{3.11}$$

*where*

$$\mathcal{V}_j^{k+1} := \sum_{i=j}^N \delta_i(\nabla_j h_i(x^{k,i}) - \nabla_j h_i(x^{k,i-1})) + \frac{1}{\gamma_j^k}(\nabla_j h_j(x^{k,j-1}) - \nabla_j h_j(x^{k,j}))$$

$$+ \frac{\alpha_j^k}{\gamma_j^k}(x_j^k - x_j^{k-1}) + \nabla_j f(x^{k+1}) - \nabla_j f(x^{k,j-1})$$

$$\mathcal{W}_j^{k+1} := \sum_{i=1}^{j-1} \delta_i(\nabla_j h_i(x^{k,i}) - \nabla_j h_i(x^{k,i-1})) - \nabla_{jj}^2 h_j(x^{k,j-1})(x_j^{k+1} - x_j^k).$$

*If* $h_i$, $i \in [N]$, *is block strongly convex, then* $\mathcal{G}^{k+1} := \left( \mathcal{G}_1^{k+1}, \ldots, \mathcal{G}_N^{k+1} \right) \in \partial \mathcal{L}(x^{k+1}, x^k)$ *and*

$$\|\mathcal{G}^{k+1}\| \leq \overline{c} \sum_{i=1}^N \sqrt{\mathbf{D}_h(x^{k,i}, x^{k,i-1})} + \widehat{c} \sum_{i=1}^N \sqrt{\mathbf{D}_h(x^{k-1,i}, x^{k-1,i-1})}, \tag{3.12}$$

*with*

$$\overline{c} := \mathbf{max}\left\{ \sqrt{2/\sigma_1}, \ldots, \sqrt{2/\sigma_N} \right\}\left( N\left(\widehat{L} + \mathbf{max}\left\{\delta_1\widetilde{L}_1, \ldots, \delta_N\widetilde{L}_N\right\}\right) + \mathbf{max}\left\{\frac{\widetilde{L}_1}{\gamma_1} + \overline{L}_1, \ldots, \frac{\widetilde{L}_N}{\gamma_N} + \overline{L}_N\right\}\right),$$

$$\widehat{c} := \mathbf{max}\left\{ \sqrt{2/\sigma_1}, \ldots, \sqrt{2/\sigma_N} \right\}\mathbf{max}\left\{\frac{\sigma_1(1-\gamma_1 L_1)}{\gamma_1}, \ldots, \frac{\sigma_N(1-\gamma_N L_N)}{\gamma_N}\right\}.$$

*Proof.* Following [42, Chapter 10], the subdifferential of $\mathcal{L}$ at $(x^{k+1}, x^k)$ is given by

$$\partial\mathcal{L}(x^{k+1}, x^k) = \left(\partial_{x^{k+1}}\mathcal{L}(x^{k+1}, x^k), \partial_{x^k}\mathcal{L}(x^{k+1}, x^k)\right), \tag{3.13}$$

where, for $j = 1, \ldots, N$, by applying [42, Exercise 8.8] we have

$$\partial_{x_j^{k+1}}\mathcal{L}(x^{k+1}, x^k) = \nabla_j f(x^{k+1}) + \partial g_j(x_j^{k+1}) + \sum_{i=j}^{N} \delta_i(\nabla_j h_i(x^{k,i}) - \nabla_j h_i(x^{k,i-1})); \tag{3.14}$$

$$\partial_{x_j^k}\mathcal{L}(x^{k+1}, x^k) = \sum_{i=1}^{j-1} \delta_i(\nabla_j h_i(x^{k,i}) - \nabla_j h_i(x^{k,i-1})) - \nabla_{jj}^2 h_j(x^{k,j-1})(x_j^{k+1} - x_j^k). \tag{3.15}$$

Writing the first-order optimality conditions for the subproblem (3.3) implies that there exists a subgradient $\eta_j^{k+1} \in \partial g_j(x_j^{k+1})$ such that

$$\nabla_j f(x^{k,j-1}) - \frac{\alpha_j^k}{\gamma_j^k}(x_j^k - x_j^{k-1}) + \frac{1}{\gamma_j^k}\left(\nabla_j h_j(x^{k,j}) - \nabla_j h_j(x^{k,j-1})\right) + \eta_j^{k+1} = 0 \quad j \in [N],$$

which implies $\eta_j^{k+1} = \frac{1}{\gamma_j^k}\left(\nabla_j h_j(x^{k,j-1}) - \nabla_j h_j(x^{k,j})\right) + \frac{\alpha_j^k}{\gamma_j^k}(x_j^k - x_j^{k-1}) - \nabla_j f(x^{k,j-1})$, $j \in [N]$. Therefore, we have $\mathcal{V}_j^{k+1} = \nabla_j f(x^{k+1}) + \eta_j^{k+1} + \sum_{i=j}^{N} \delta_i(\nabla_j h_i(x^{k,i}) - \nabla_j h_i(x^{k,i-1})) \in \partial_{x_j^{k+1}}\mathcal{L}(x^{k+1}, x^k)$, which implies $\mathcal{G}^{k+1} \in \partial\mathcal{L}(x^{k+1}, x^k)$. Together with the Lipschitz continuity of $\nabla_i f$, $\nabla_i h_i$ and the boundedness of $\nabla_{ii}^2 h_i$ on bounded sets, the boundedness of $(x^k)_{k\in\mathbb{N}}$, and the triangle inequality, this implies that there exist constants $\widehat{L}$, $\widehat{L}_i$, $\overline{L}_i > 0$ (for $i \in [N]$) such that

$$\|\mathcal{G}_j^{k+1}\| = \|\mathcal{V}_j^{k+1}\| + \|\mathcal{W}_j^{k+1}\| \leq \frac{\alpha_j^k}{\gamma_j^k}\|x_j^k - x_j^{k-1}\| + \|\nabla_j f(x^{k+1}) - \nabla_j f(x^{k,j-1})\|$$

$$+ \sum_{i=1}^{N} \delta_i\|\nabla_j h_i(x^{k,i}) - \nabla_j h_i(x^{k,i-1})\| + \frac{1}{\gamma_j^k}\|\nabla_j h_j(x^k) - \nabla_j h_j(x^{k,1})\| + \|\nabla_{jj}^2 h_j(x^{k,j-1})\|\,\|x_j^{k+1} - x_j^k\|$$

$$\leq \frac{\alpha_j^k}{\gamma_j^k}\|x_j^k - x_j^{k-1}\| + \widehat{L}\sum_{i=1}^{N}\|x_i^{k+1} - x_i^k\| + \sum_{i=1}^{N}\delta_i\widetilde{L}_i\|x_i^{k+1} - x_i^k\| + \left(\frac{\widetilde{L}_j}{\gamma_j^k} + \overline{L}_j\right)\|x_j^{k+1} - x_j^k\|.$$

Combining the last two inequalities with (3.9), it can be deduced that

$$\|\mathcal{G}^{k+1}\| \leq \left(N\left(\widehat{L} + \mathbf{max}\left\{\delta_1\widetilde{L}_1, \ldots, \delta_N\widetilde{L}_N\right\}\right) + \mathbf{max}\left\{\frac{\widetilde{L}_1}{\gamma_1^k} + \overline{L}_1, \ldots, \frac{\widetilde{L}_N}{\gamma_N^k} + \overline{L}_N\right\}\right)\sum_{i=1}^{N}\|x_i^{k+1} - x_i^k\|$$

$$+ \mathbf{max}\left\{\frac{\alpha_1^k}{\gamma_1^k}, \ldots, \frac{\alpha_N^k}{\gamma_N^k}\right\}\sum_{i=1}^{N}\|x_i^k - x_i^{k-1}\|$$

$$\leq \left(N\left(\widehat{L} + \mathbf{max}\left\{\delta_1\widetilde{L}_1, \ldots, \delta_N\widetilde{L}_N\right\}\right) + \mathbf{max}\left\{\frac{\widetilde{L}_1}{\gamma_1} + \overline{L}_1, \ldots, \frac{\widetilde{L}_N}{\gamma_N} + \overline{L}_N\right\}\right)\sum_{i=1}^{N}\|x_i^{k+1} - x_i^k\|$$

$$+ \mathbf{max}\left\{\frac{\sigma_1(1-\gamma_1 L_1)}{\gamma_1}, \ldots, \frac{\sigma_N(1-\gamma_N L_N)}{\gamma_N}\right\}\sum_{i=1}^{N}\|x_i^k - x_i^{k-1}\|$$

$$\leq \overline{c}\sum_{i=1}^{N}\|x_i^{k+1} - x_i^k\| + \widehat{c}\sum_{i=1}^{N}\|x_i^k - x_i^{k-1}\|.$$

Hence, it follows from the block strong convexity of $h_i$ $(i = 1, \ldots, N)$ that

$$\|\mathcal{G}^{k+1}\| \leq \overline{c}\sum_{i=1}^{N}\|x_i^{k+1} - x_i^k\| + \widehat{c}\sum_{i=1}^{N}\|x_i^k - x_i^{k-1}\|$$

$$\leq \overline{c}\sum_{i=1}^{N}\sqrt{\mathbf{D}_h(x^{k,i}, x^{k,i-1})} + \widehat{c}\sum_{i=1}^{N}\sqrt{\mathbf{D}_h(x^{k-1,i}, x^{k-1,i-1})},$$

giving our desired result. $\qquad\square$

**Remark 3.8.** Note that a uniformly continuous function maps bounded sets to bounded sets. Therefore, in Proposition 3.7, if the function $\nabla_{ii}^2 h_i$ ($i = 1, \ldots, N$) is uniformly continuous, it is bounded on bounded sets. □

Applying Proposition 3.7, the *subsequential convergence* of $(x^k)_{k \in \mathbb{N}}$ generated by BIBPA is presented next. On top of that we explain some basic properties of $\omega(x^0)$.

**Assumption III.** $\overline{C} \subseteq \mathbf{int} \, \mathbf{dom} \, h_1$.

**Theorem 3.9** (subsequential convergence and properties of $\omega(x^0)$). *Let all assumptions of Proposition 3.7 and Assumption II hold. Then, the following assertions are satisfied:*

(i) *every cluster point of $(x^k)_{k \in \mathbb{N}}$ is a critical point of $\Phi$, i.e., $\omega(x^0) \subset \mathbf{crit} \, \Phi$;*

(ii) $\lim_{k \to \infty} \mathbf{dist} \left( x^k, \omega(x^0) \right) = 0$;

(iii) $\omega(x^0)$ *is a nonempty, compact, and connected set;*

(iv) *the Lyapunov function $\mathcal{L}$ is finite and constant on $\omega(x^0)$.*

*Proof.* Let us assume $x^\star = (x_1^\star, \ldots, x_N^\star) \in \omega(x^0)$. The boundedness of $(x^k)_{k \in \mathbb{N}}$ implies that there exists an infinite index set $\mathcal{J} \subset \mathbb{N}$ such that the subsequence $(x^k)_{k \in \mathcal{J}} \to x^\star$ as $k \to \infty$. It follows from (3.5) that

$$
\begin{aligned}
&\langle \nabla_i f(x^{k,i-1}) - \tfrac{\alpha_i^k}{\gamma_i^k}(x_i^k - x_i^{k-1}), x_i^{k+1} - x_i^k \rangle + \tfrac{1}{\gamma_i^k} \mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1}) + g_i(x_i^{k+1}) \\
&\leq \langle \nabla_i f^k(x^{k,i-1}) - \tfrac{\alpha_i^k}{\gamma_i^k}(x_i^k - x_i^{k-1}), x_i^\star - x_i^k \rangle + \tfrac{1}{\gamma_i^k} \mathbf{D}_{h_i}(x^\star, x^{k,i-1}) + g_i(x_i^\star).
\end{aligned}
\tag{3.16}
$$

Invoking Proposition 3.5 and using block strong convexity of $h_i$, there exist $\varepsilon_i^\star > 0$, $k_i^0 \in \mathbb{N}$, and a neighborhood $\mathbf{B}(x_i^\star, \varepsilon_i^\star)$ such that $\lim_{k \to \infty} \tfrac{\sigma_i}{2}\|x_i^{k+1} - x_i^k\|^2 \leq \lim_{k \to \infty} \mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1}) = 0$, $x_i^k \in \mathbf{B}(x_i^\star, \varepsilon_i^\star)$, $i \in [N]$, for $k \geq k_i^0$ and $k \in \mathcal{J}$, i.e., $\lim_{k \to \infty}(x_i^{k+1} - x_i^k) = 0$. Hence, substituting $k = k_j - 1$ for $k_j \in \mathcal{J}$ into (3.16) and taking the limit from both sides of (3.16), we derive $\limsup_{j \to \infty} g_i(x_i^{k_j}) \leq g_i(x_i^\star)$ $i = 1 \in [N]$. Furthermore, since $g_i$ is lsc, this yields that $\lim_{j \to \infty} g_i(x_i^{k_j}) = g_i(x_i^\star)$, then

$$
\lim_{j \to \infty} \mathcal{L}(x^{k_j+1}, x^{k_j}) = \lim_{j \to \infty} \left( f(x_1^{k_j}, \ldots, x_N^{k_j}) + \sum_{i=1}^N g_i(x_i^{k_j}) + \sum_{i=1}^N \delta_i \mathbf{D}_{h_i}(x^{k_j,i}, x^{k_j,i-1}) \right) = \mathcal{L}(x^\star, x^\star).
$$

Hence, from (3.12) and Proposition 3.5, we obtain

$$
\lim_{k \to +\infty} \|\mathcal{G}^{k+1}\| \leq \lim_{k \to +\infty} \left( \overline{c} \sum_{i=1}^N \sqrt{\mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1})} + \widehat{c} \sum_{i=1}^N \sqrt{\mathbf{D}_{h_i}(x^{k-1,i}, x^{k-1,i-1})} \right) = 0,
$$

which consequently yields $\lim_{k \to \infty} \mathcal{G}^{k+1} = 0$. As a result, we have $0 \in \partial \mathcal{L}(x^\star, x^\star)$, owing to the closedness of the subdifferential mapping $\partial \mathcal{L}$. The result of Theorem 3.9(i) follows from the fact $\partial \mathcal{L}(x^\star, x^\star) = (\partial \Phi(x^\star), 0)$. Moreover, Theorem 3.9(ii) is a straightforward consequence of Theorem 3.9(i), and Theorem 3.9(iii) and Theorem 3.9(iv) can be proved in the same way as [19, Lemma 5(iii)-(iv)]. □

3.1. **Global convergence for KŁ functions.** In this section, we consider the class of Kurdyka-Łojasiewicz (KŁ) functions (see [32, 34]) and show that for such functions the sequence $(x^k)_{k \in \mathbb{N}}$ converges to a critical point $x^\star$.

**Definition 3.10** (KŁ property). *A proper and lsc function $\varphi : \mathbb{R}^n \to \overline{\mathbb{R}}$ has the KŁ property at $x^\star \in \mathbf{dom} \, \varphi$ if there exist a concave function $\psi : [0, \eta] \to [0, +\infty[$ (with $\eta > 0$) and*

neighborhood $\mathbf{B}(x^\star; \varepsilon)$ with $\varepsilon > 0$, such that (i) $\psi(0) = 0$; (ii) $\psi$ is of class $C^1$ with $\psi > 0$ on $(0, \eta)$; (ii) for all $x \in \mathbf{B}(x^\star; \varepsilon)$ such that $\varphi(x^\star) < \varphi(x) < \varphi(x^\star) + \eta$ it holds that

$$\psi'(\varphi(x) - \varphi(x^\star)) \, \mathbf{dist}(0, \partial\varphi(x)) \geq 1. \tag{3.17}$$

*If this property holds for each point of* $\mathbf{dom}\, \partial\varphi$, *the* $\varphi$ *is a KŁ function.*

In [33, 34], Stanisław Łojasiewicz showed for the first time that every real analytic function[1] satisfies (3.17) with $\psi(s) := \frac{\kappa}{1-\theta} s^{1-\theta}$ with $\theta \in [0, 1)$. In 1998, Kurdyka [32] proved that this inequality is valid for $C^1$ functions whose graph belong to an *o-minimal structure* (see its definition in [25]). Later, (3.17) was extended for nonsmooth functions in [17, 16, 18].

The KŁ property (3.17) of the underlying objective function plays a key role in establishing the global convergence of a generic algorithm for nonconvex problems; however, this is not sufficient and one also needs some additional conditions to be guaranteed by the algorithm (see below). In particular, for several algorithms the cost functions satisfy the sufficient decrease condition (cf. [2, 6, 19]), while for some others the sufficient decrease condition is satisfied for some Lyapunov functions (cf. [26, 40, 39, 41, 51]).

As shown in Proposition 3.5, Proposition 3.7, and Theorem 3.9 (see its proof), the sequence $(x^k)_{k\in\mathbb{N}}$ generated by BIBPA satisfies the following conditions that are non-Euclidean extension of those given in [6, 19] for the structured problem (2.1):

1) (*sufficient descent condition*) For each $k \in \mathbb{N}$ and $a_i, b_i \geq 0$ $(i = 1, \ldots, N)$,

$$\sum_{i=1}^{N} \left( a_i \, \mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1}) + b_i \, \mathbf{D}_{h_i}(x^{k-1,i}, x^{k-1,i-1}) \right) \leq \mathcal{L}(x^k, x^{k-1}) - \mathcal{L}(x^{k+1}, x^k);$$

2) (*subgradient lower bound of iteration gap*) For each $k \in \mathbb{N}$, there exists a subgradient $\mathcal{G}^{k+1} \in \partial\mathcal{L}(x^{k+1}, x^k)$ and $\overline{c}, \widehat{d} \geq 0$ such that

$$\|\mathcal{G}^{k+1}\| \leq \overline{c} \sum_{i=1}^{N} \sqrt{\mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1})} + \widehat{c} \sum_{i=1}^{N} \sqrt{\mathbf{D}_{h_i}(x^{k-1,i}, x^{k-1,i-1})};$$

3) (*continuity condition*) The function $\mathcal{L}$ is a KŁ function, and each cluster point $x^\star$ of $(x^k)_{k\in\mathbb{N}}$ $(x^\star \in \omega(x^0))$ satisfies $(x^\star, x^\star) \in \mathrm{crit}\mathcal{L}$

We now use the above three conditions to prove that the whole sequence $(x^k)_{k\in\mathbb{N}}$ converges.

**Theorem 3.11** (global convergence). *Let all assumptions of* Proposition 3.7 *and* Assumption II *hold. If* $\mathcal{L}$ *is a KŁ function, then the following statements are true:*

(i) *The sequence* $(x^k)_{k\in\mathbb{N}}$ *has finite length, i.e.,*

$$\sum_{k=1}^{\infty} \|x_i^{k+1} - x_i^k\| < \infty \quad i = 1, \ldots, N; \tag{3.18}$$

(ii) *The sequence* $(x^k)_{k\in\mathbb{N}}$ *converges to a stationary point* $x^\star$ *of* $\Phi$.

*Proof.* Define the sequence $(d_k)_{k\in\mathbb{N}}$ as $d_k := \sum_{i=1}^{N} \sqrt{\mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1})} + \sqrt{\mathbf{D}_{h_i}(x^{k-1,i}, x^{k-1,i-1})}$. From Proposition 3.7 for $\widetilde{c} := \mathbf{max}\{\overline{c}, \widehat{c}\}$, we obtain

$$\begin{aligned} \|\mathcal{G}^{k+1}\| &\leq \overline{c} \sum_{i=1}^{N} \sqrt{\mathbf{D}_h(x^{k,i}, x^{k,i-1})} + \widehat{c} \sum_{i=1}^{N} \sqrt{\mathbf{D}_h(x^{k-1,i}, x^{k-1,i-1})} \\ &\leq \widetilde{c} \sum_{i=1}^{N} \left( \sqrt{\mathbf{D}_{h_i}(x^{k,i}, x^{k,i-1})} + \sqrt{\mathbf{D}_{h_i}(x^{k-1,i}, x^{k-1,i-1})} \right) = \widetilde{c} d_k. \end{aligned} \tag{3.19}$$

---

[1] A function $\varphi : \mathbb{R}^n \to \overline{\mathbb{R}}$ said to be real analytic if it can be represented by a convergent power series.

Applying twice the root-mean square and arithmetic mean inequalitywe come to

$$
\begin{aligned}
d_k &\leq \sqrt{N \sum_{i=1}^{N} \mathbf{D}_{h_i}(\boldsymbol{x}^{k,i}, \boldsymbol{x}^{k,i-1})} + \sqrt{N \sum_{i=1}^{N} \mathbf{D}_{h_i}(\boldsymbol{x}^{k-1,i}, \boldsymbol{x}^{k-1,i-1})} \\
&\leq \sqrt{2N \sum_{i=1}^{N} \left( \mathbf{D}_{h_i}(\boldsymbol{x}^{k,i}, \boldsymbol{x}^{k,i-1}) + \mathbf{D}_{h_i}(\boldsymbol{x}^{k-1,i}, \boldsymbol{x}^{k-1,i-1}) \right)}.
\end{aligned}
\tag{3.20}
$$

Then, it can be concluded from Proposition 3.5 and (3.20) that

$$
\begin{aligned}
\mathcal{L}^k - \mathcal{L}^{k+1} &\geq \sum_{i=1}^{N} \left( a_i \, \mathbf{D}_{h_i}(\boldsymbol{x}^{k,i}, \boldsymbol{x}^{k,i-1}) + b_i \, \mathbf{D}_{h_i}(\boldsymbol{x}^{k-1,i}, \boldsymbol{x}^{k-1,i-1}) \right) \\
&\geq \varrho \sum_{i=1}^{N} \left( \mathbf{D}_{h_i}(\boldsymbol{x}^{k,i}, \boldsymbol{x}^{k,i-1}) + \mathbf{D}_{h_i}(\boldsymbol{x}^{k-1,i}, \boldsymbol{x}^{k-1,i-1}) \right) \geq \tfrac{\varrho}{2N} d_k^2,
\end{aligned}
$$

where $\varrho := \min\{a_1, b_1, \ldots, a_N, b_N\}$. Together with (3.19) and Theorem 3.9*(i)*, this implies that [39, Assumption H] holds true with $a_k = \frac{\varrho}{2N}, b_k = 1, b = \widetilde{c}, I = \{1\}, \varepsilon_k = 0$. Therefore, since $\mathcal{L}$ is a proper lower semicontinuous KŁ function, [39, Theorem 10] yields that Theorem 3.11*(i)* holds true and the sequence $(\boldsymbol{x}^k)_{k \in \mathbb{N}}$ converges to $\boldsymbol{x}^\star$ in which $(\boldsymbol{x}^\star, \boldsymbol{x}^\star)$ is a stationary point of the Lyapunov function $\mathcal{L}$ (3.8), i.e., $0 \in \partial \mathcal{L}(\boldsymbol{x}^\star, \boldsymbol{x}^\star)$. Finally, the result follows from the fact $\partial \mathcal{L}(\boldsymbol{x}^\star, \boldsymbol{x}^\star) = (\partial \Phi(\boldsymbol{x}^\star), 0)$. $\qquad\square$

## 3.2. Rate of convergence for Łojasiewicz-type KŁ functions.

We now investigate the convergence rate of the generated sequence under KŁ inequality of Łojasiewicz-type at $x^\star$ ($\psi(s) := \frac{\kappa}{1-\theta} s^{1-\theta}$ with $\theta \in [0, 1)$), i.e., there exists $\varepsilon > 0$ such that

$$
|\varphi(\boldsymbol{x}) - \varphi^\star|^\theta \leq \kappa \, \mathbf{dist}(0, \partial \varphi(\boldsymbol{x})) \quad \forall \boldsymbol{x} \in \mathbf{B}(\boldsymbol{x}^\star; \varepsilon).
\tag{3.21}
$$

**Fact 3.12** (convergence rate of a sequence with positive elements). [22, Lemma 15] Let $(s_k)_{k \in \mathbb{N}}$ be a monotonically decreasing sequence in $\mathbb{R}_+$ and let $\theta \in [0, 1)$ and $\beta > 0$. Suppose that $s_k^{2\theta} \leq \beta(s_k - s_{k+1})$ holds for all $k \in \mathbb{N}$. Then, the following assertions hold:

  *(i)* If $\theta = 0$, the sequences $(s_k)_{k \in \mathbb{N}}$ converges in a finite time;

 *(ii)* If $\theta \in (0, 1/2]$, there exist $\lambda > 0$ and $\tau \in [0, 1)$ such that $0 \leq s_k \leq \lambda \tau^k$ for every $k \in \mathbb{N}$.

*(iii)* If $\theta \in (1/2, 1)$, there exists $\mu > 0$ such that $0 \leq s_k \leq \mu k^{-\frac{1}{2\theta-1}}$ for every $k \in \mathbb{N}$

Let $(\mathcal{S}_k)_{k \in \mathbb{N}}$ given by $\mathcal{S}_k := \mathcal{L}(\boldsymbol{x}^k, \boldsymbol{x}^{k-1}) - \mathcal{L}(\boldsymbol{x}^\star, \boldsymbol{x}^\star)$. We next derive the *convergence rates* of $(\boldsymbol{x}^k)_{k \in \mathbb{N}}$ and $(\mathcal{S}_k)_{k \in \mathbb{N}}$ when $\mathcal{L}$ satisfies the KŁ inequality of Łojasiewicz type.

**Theorem 3.13** (convergence rate). *Let all assumptions of Proposition 3.7 and Assumption II hold, and $(\boldsymbol{x}^k)_{k \in \mathbb{N}}$ converges to $\boldsymbol{x}^\star$. If $\mathcal{L}$ satisfies the KŁ inequality of Łojasiewicz type, then the following assertions hold:*

  *(i) if $\theta = 0$, then the sequences $(\boldsymbol{x}^k)_{k \in \mathbb{N}}$ and $(\Phi(\boldsymbol{x}^k))_{k \in \mathbb{N}}$ converge in a finite number of steps to $\boldsymbol{x}^\star$ and $\Phi(\boldsymbol{x}^\star)$, respectively;*

 *(ii) if $\theta \in (0, 1/2]$, then there exist $\lambda_1 > 0$, $\mu_1 > 0$, $\tau, \overline{\tau} \in [0, 1)$, and $\overline{k} \in \mathbb{N}$ such that*

$$
0 \leq \|\boldsymbol{x}^k - \boldsymbol{x}^\star\| \leq \lambda_1 \tau^k, \quad 0 \leq \mathcal{S}_k \leq \mu_1 \overline{\tau}^k \quad \forall k \geq \overline{k};
$$

*(iii) if $\theta \in (1/2, 1)$, then there exist $\lambda_2 > 0$, $\mu_2 > 0$, and $\overline{k} \in \mathbb{N}$ such that*

$$
0 \leq \|\boldsymbol{x}^k - \boldsymbol{x}^\star\| \leq \lambda_2 k^{-\frac{1-\theta}{2\theta-1}}, \quad 0 \leq \mathcal{S}_k \leq \mu_2 k^{-\frac{1-\theta}{2\theta-1}} \quad \forall k \geq \overline{k} + 1.
$$

*Proof.* We first set $\varepsilon > 0$ to be that a constant described in (3.21) and $x^k \in \mathbf{B}(x^\star; \varepsilon)$ for all $k \geq \tilde{k}$ and $\tilde{k} \in \mathbb{N}$. Let us define $\Delta_k := \psi(\mathcal{L}(\boldsymbol{x}^k, \boldsymbol{x}^{k-1}) - \mathcal{L}(\boldsymbol{x}^\star, \boldsymbol{x}^\star)) = \psi(\mathcal{S}_k)$. Then, it follows

from the concavity of $\psi$ and 2) that

$$
\begin{aligned}
\Delta_k - \Delta_{k+1} &= \psi(\mathcal{S}_k) - \psi(\mathcal{S}_{k+1}) \geq \psi'(\mathcal{S}_k)(\mathcal{S}_k - \mathcal{S}_{k+1}) \\
&= \psi'(\mathcal{S}_k)(\mathcal{L}(\mathbf{x}^k, \mathbf{x}^{k-1}) - \mathcal{L}(\mathbf{x}^{k+1}, \mathbf{x}^k)) \geq \frac{\mathcal{L}(\mathbf{x}^k, \mathbf{x}^{k-1}) - \mathcal{L}(\mathbf{x}^{k+1}, \mathbf{x}^k)}{\text{dist}(0, \partial \mathcal{L}(\mathbf{x}^k, \mathbf{x}^{k-1}))} \\
&\geq \frac{\sum_{i=1}^{N} \left( a_i \, \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) + b_i \, \mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1}) \right)}{\bar{c} \sum_{i=1}^{N} \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1})} + \widehat{c} \sum_{i=1}^{N} \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-2,i}, \mathbf{x}^{k-2,i-1})}} \\
&\geq \frac{1}{c} \frac{\sum_{i=1}^{N} \left( \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) + \mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1}) \right)}{\sum_{i=1}^{N} \left( \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1})} + \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-2,i}, \mathbf{x}^{k-2,i-1})} \right)},
\end{aligned}
$$

with $c := \max\{\bar{c}, \widehat{c}\} / \min\{a_1, b_1, \ldots, a_N, b_N\}$. Using (3.20) and applying the arithmetic mean and geometric mean inequality, it can be concluded that

$$
\begin{aligned}
d_k &\leq \sqrt{2N \sum_{i=1}^{N} \left( \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) + \mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1}) \right)} \\
&\leq \sqrt{2cN(\Delta_k - \Delta_{k+1}) \sum_{i=1}^{N} \left( \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1})} + \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-2,i}, \mathbf{x}^{k-2,i-1})} \right)} \qquad (3.22) \\
&\leq cN(\Delta_k - \Delta_{k+1}) + \tfrac{1}{2} \sum_{i=1}^{N} \left( \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1})} + \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-2,i}, \mathbf{x}^{k-2,i-1})} \right)
\end{aligned}
$$

We now define the sequences $(a_k)_{k \in \mathbb{N}}$ and $(b_k)_{k \in \mathbb{N}}$ given by

$$
p_{k+1} := \sum_{i=1}^{N} \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1})} + \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1})}, \quad q_k = cN(\Delta_k - \Delta_{k+1}), \quad \alpha := \tfrac{1}{2}, \quad (3.23)
$$

where $\sum_{i=1}^{\infty} q_k = 2cN \sum_{i=1}^{\infty} (\Delta_i - \Delta_{i+1}) = \Delta_1 - \Delta_\infty = \Delta_1 < \infty$. This and (3.22) yield $p_{k+1} \leq \tfrac{1}{2} p_k + q_k$ for all $k \geq \tilde{k}$. Since $(\Phi)_{k \in \mathbb{N}}$ is non-increasing,

$$
\sum_{j=k}^{\infty} p_{j+1} \leq \tfrac{1}{2} \sum_{j=k}^{\infty} (p_j - p_{j+1} + p_{j+1}) + 2cN \sum_{j=k}^{\infty} \left( \Delta_j - \Delta_{j+1} \right) = \tfrac{1}{2} \sum_{j=k}^{\infty} p_{j+1} + \tfrac{1}{2} p_k + 2cN\Delta_k.
$$

From the root-mean square, the arithmetic mean inequality, $\psi(\mathcal{S}_k) \leq \psi(\mathcal{S}_{k-1})$, and Proposition 3.5, this lead to

$$
\begin{aligned}
\sum_{j=k}^{\infty} p_{j+1} &\leq p_k + 4cN\Delta_k = \sum_{i=1}^{N} \left( \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1})} + \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-2,i}, \mathbf{x}^{k-2,i-1})} \right) + 4cN\psi(\mathcal{S}_k) \\
&\leq \sqrt{N \sum_{i=1}^{N} \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1})} + \sqrt{N \sum_{i=1}^{N} \mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1})} + 4cN\psi(\mathcal{S}_k) \\
&\leq \sqrt{2N \sum_{i=1}^{N} \left( \mathbf{D}_{h_i}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) + \mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1}) \right)} + 4cN\psi(\mathcal{S}_k) \\
&\leq \sqrt{2N/\varrho} \sqrt{\mathcal{S}_{k-1} - \mathcal{S}_k} + 4cN\psi(\mathcal{S}_{k-1}),
\end{aligned}
$$
$$(3.24)$$

with $\varrho := \min\{a_1, b_1, \ldots, a_N, b_N\}$. Since $\mathbf{D}_{h_i}(\cdot, \cdot) \geq 0$, for $i = 1, \ldots, N$, it holds that

$$
\begin{aligned}
\|x_i^k - x_i^\star\| &\leq \|x_i^{k+1} - x_i^k\| + \|x_i^{k+1} - x_i^\star\| \leq \ldots \leq \sum_{j=k}^{\infty} \|x_i^{j+1} - x_i^j\| \\
&\leq \sum_{j=k}^{\infty} \sqrt{\tfrac{2}{\sigma_i} \mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1})} \leq \sqrt{\tfrac{2}{\sigma_i}} \sum_{j=k}^{\infty} \left( \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1})} + \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-2,i}, \mathbf{x}^{k-2,i-1})} \right).
\end{aligned}
$$

Combining this with (3.24) and setting $\rho := \max\left\{ \sqrt{2/\sigma_1}, \ldots, \sqrt{2/\sigma_N} \right\}$, we come to

$$
\begin{aligned}
\sum_{i=1}^{N} \|x_i^k - x_i^\star\| &\leq \rho \sum_{j=k}^{\infty} \sum_{i=1}^{N} \left( \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1})} + \sqrt{\mathbf{D}_{h_i}(\mathbf{x}^{k-2,i}, \mathbf{x}^{k-2,i-1})} \right) \\
&\leq \rho \sqrt{2N/\varrho} \sqrt{\mathcal{S}_{k-1} - \mathcal{S}_k} + 4c\rho N\psi(\mathcal{S}_{k-1}),
\end{aligned}
$$

which consequently yields

$$
\|x_i^k - x_i^\star\| \leq \nu \max\left\{ \sqrt{\mathcal{S}_{k-1}}, \psi(\mathcal{S}_{k-1}) \right\} \quad i = 1, \ldots, N, \qquad (3.25)
$$

with $\nu := \rho \sqrt{2N/\varrho} + 4c\rho N$ and $\psi(s) := \frac{\kappa}{1-\theta} s^{1-\theta}$. Furthermore, the nonlinear equation $\sqrt{\mathcal{S}_{k-1}} - \frac{\kappa}{1-\theta} \mathcal{S}_{k-1}^{1-\theta} = 0$ has a solution at $\mathcal{S}_{k-1} = ((1-\theta)/\kappa)^{\frac{2}{1-2\theta}}$. For $\hat{k} \in \mathbb{N}$ and $k \geq \hat{k}$, we

assume (3.25) holds and $\mathcal{S}_{k-1} \leq \left(\frac{\kappa}{1-\theta}\right)^{\frac{2}{1-2\theta}}$. Two cases are recognized: (a) $\theta \in (0, 1/2]$; (b) $\theta \in (1/2, 1)$. In Case (a), if $\theta \in (0, 1/2)$, then $\psi(\mathcal{S}_{k-1}) \leq \sqrt{\mathcal{S}_{k-1}}$. For $\theta = 1/2$, we get $\psi(\mathcal{S}_{k-1}) = \frac{\kappa}{1-\theta} \sqrt{\mathcal{S}_{k-1}}$, which implies $\mathbf{max}\left\{\sqrt{\mathcal{S}_{k-1}}, \psi(\mathcal{S}_{k-1})\right\} = \mathbf{max}\left\{1, \frac{\kappa}{1-\theta}\right\} \sqrt{\mathcal{S}_{k-1}}$. Then, $\mathbf{max}\left\{\sqrt{\mathcal{S}_{k-1}}, \psi(\mathcal{S}_{k-1})\right\} \leq \mathbf{max}\left\{1, \frac{\kappa}{1-\theta}\right\} \sqrt{\mathcal{S}_{k-1}}$. In Case (b), it holds that $\psi(\mathcal{S}_{k-1}) \geq \sqrt{\mathcal{S}_{k-1}}$, i.e., $\mathbf{max}\left\{\sqrt{\mathcal{S}_{k-1}}, \psi(\mathcal{S}_{k-1})\right\} = \frac{\kappa}{1-\theta} \mathcal{S}_{k-1}^{1-\theta}$. Combining both cases, for all $k \geq \bar{k} := \mathbf{max}\{\tilde{k}, \hat{k}\}$, we end up with

$$\|x_i^k - x_i^\star\| \leq \begin{cases} \nu \, \mathbf{max}\left\{1, \frac{\kappa}{1-\theta}\right\} \sqrt{\mathcal{S}_{k-1}} & \text{if } \theta \in (0, 1/2], \\ \nu \frac{\kappa}{1-\theta} \mathcal{S}_{k-1}^{1-\theta} & \text{if } \theta \in (1/2, 1). \end{cases} \tag{3.26}$$

On the other hand, it follows from Proposition 3.5 that

$$\begin{aligned} &\mathcal{S}_{k-1} - \mathcal{S}_k \\ &= \mathcal{L}(x^{k-1}, x^{k-2}) - \mathcal{L}(x^k, x^{k-1}) \geq \varrho \sum_{i=1}^N \left(\mathbf{D}_h(x^{k-1,i}, x^{k-1,i-1}) + \mathbf{D}_h(x^{k-2,i}, x^{k-2,i-1})\right) \\ &\geq \frac{\varrho}{2N}\left(\sqrt{\mathbf{D}_{h_i}(x^{k-1,i}, x^{k-1,i-1})} + \sqrt{\mathbf{D}_{h_i}(x^{k-2,i}, x^{k-2,i-1})}\right)^2 \\ &\geq \frac{\varrho}{2Nc^2}\|(\mathcal{G}_1^k, \ldots, \mathcal{G}_N^k)\|^2 \geq \frac{\varrho}{2Nc^2}\,\mathbf{dist}(0, \partial\mathcal{L}(x^k, x^{k-1}))^2 \geq \frac{\varrho}{2Nc^2\kappa^2}\mathcal{S}_{k-1}^\theta = c_2 \mathcal{S}_{k-1}^\theta, \end{aligned}$$

where $c_2 := \frac{\varrho}{2Nc^2\kappa^2}$. The results then follow from $\mathcal{S}_k \to 0$, (3.26) and Fact 3.12. $\qquad\square$

## 4. APPLICATION TO SYMMETRIC NONNEGATIVE MATRIX TRI-FACTORIZATION

A natural way of analyzing large data sets is finding an effective way to represent them using dimensionality reduction methodologies. *Nonnegative matrix factorization* (NMF) is one such technique that has received much attention in the last few years; see, e.g., [24, 27] and the references therein. In order to extract hidden and important features from data, NMF decomposes the data matrix into two factor matrices (usually much smaller than the original data matrix) by imposing componentwise nonnegativity and (possibly) other constraints such as sparsity to take prior information into account. More precisely, let the data matrix be $X = [x_1, x_2, \ldots, x_n] \in \mathbb{R}_+^{m \times n}$ where each $x_i$ represents some data point. NMF seeks a decomposition of $X$ into a nonnegative $n \times r$ basis matrix $U = [u_1, u_2, \ldots, u_r] \in \mathbb{R}_+^{m \times r}$ and a nonnegative $r \times n$ coefficient matrix $V = [v_1, v_2, \ldots, v_r]^T \in \mathbb{R}_+^{r \times n}$ such that

$$X \approx UV, \tag{4.1}$$

where $\mathbb{R}_+^{m \times n}$ is the set of $m \times n$ nonnegative matrices. Extensive research has been carried out on variants of NMF, and most studies have focused on algorithmic developments, but with very limited convergence theory. This motivates us to study the application of BIBPA to a variant of NMF, namely SymTriNMF; see (2.6) for the formulation of SymTriNMF as an optimization problem.

One popular application of SymTriNMF is community detection. Let $X$ be the adjacency matrix of graph so that $X_{ij} = 1$ if item $i$ is connected to item $j$, and $X_{ij} = 0$ otherwise. Let also $X \approx UVU^T$ be a SymTriNMF decomposition of $X$. Each column of $U$ corresponds to a community, that is, to a subset of items highly connected. In other words, the entry $U_{jk}$ of $U$ indicates the membership of item $j$ within community $k$, and $U_{jk} > 0$ if $j$ belongs to community $k$. The $r$-by-$r$ matrix $V$ indicates the relationship between communities, that is, whether the items within two communities are likely to interact: $V_{kp}$ is the "strength" of the interaction between the $k$th and $p$th communities. We have $X \approx \sum_{k=1}^r \sum_{p=1}^r U_{:k}V_{k,p}U_{:p}^T$, so that $X$ is decomposed via the sum of $r^2$ rank-one factors corresponding to the $r$ communities and their interactions; see [49, 52] for more details. Note that SymTriNMF is closely related to the mixed membership stochastic blockmodel [4].

Given $U^k$ and $V^k$, we next derive the closed-form solutions for $U^{k+1}$ and $V^{k+1}$.

**Theorem 4.1** (closed-form solutions of the subproblem (3.5) for SymTriNMF). *Let $h_1$ and $h_2$ be the kernel functions given in* (2.7) *and* (2.8) *and $U^k$ and $V^k$ are given. Then,*

(i) *the iteration $U^{k+1}$ of the subproblem* (3.5) *is given by*

$$U^{k+1} = \frac{1}{t_k} \max\left\{ \frac{1}{\gamma_1^k}(\nabla_U h_1(U^k, V^k) - \gamma_1^k \nabla_U f(U^k, V^k) + \alpha_1^k(U^k - U^{k-1})), 0 \right\} \quad (4.2)$$

*with*

$$\nabla_U f(U^k, V^k) = -X U^k(V^k)^T - X^T U^k V^k + U^k V^k(U^k)^T U^k(V^k)^T + U^k(V^k)^T(U^k)^T U^k V^k,$$

$$\nabla_U h_1(U^k, V^k) = \left(a_1 \|U^k\|_F^2 \|V^k\|_F^2 + b_1(\|X\|_F \|V^k\|_F + \varepsilon_1)\right) U^k,$$

*and*

$$t_k = \frac{\tau_1}{3} + \sqrt[3]{\frac{\tau_2 + \sqrt{\Delta_1}}{2} + \frac{\tau_1^3}{27}} + \sqrt[3]{\frac{\tau_2 - \sqrt{\Delta_1}}{2} + \frac{\tau_1^3}{27}}, \quad (4.3)$$

*where $\tau_1 = b_1(\|X\|_F\|V^k\|_F + \varepsilon_1)$,   $\tau_2 = a_1\|V^k\|_F^2\| \max\{G^k, 0\}\|_F^2$,   $\Delta_1 = \tau_2^2 + \frac{4}{27}\tau_2^2\tau_1^3$ with $G^k := \frac{1}{\gamma_1^k}\left(\nabla_U h_1(U^k, V^k) - \gamma_1^k \nabla_U f(U^k, V^k) + \alpha_1^k(U^k - U^{k-1})\right)$.*

(ii) *for $\eta_k := a_2\|U^{k+1}\|^4 + \varepsilon_2$, the iteration $V^{k+1}$ of the subproblem* (3.5) *is given by*

$$V^{k+1} = \max\left\{ V^k - \frac{1}{\eta_k}\left(\alpha_2^k(V^k - V^{k-1}) - \gamma_2^k \nabla_V f(U^{k+1}, V^k)\right), 0 \right\}, \quad (4.4)$$

*with $\nabla_V f(U^{k+1}, V^k) = (U^{k+1})^T X U^{k+1} + (U^{k+1})^T U^{k+1} V^k(U^{k+1})^T U^{k+1}$.*

*Proof.* Setting $g_1 := \delta_{U \geq 0}$ and $f(U, V) = \frac{1}{2}\|X - UVU^T\|_F^2$, it follows from (3.5) that

$$U^{k+1} = \arg\min_{U \in \mathbb{R}^{m \times r}} \left\{ \langle \nabla_U f(U^k, V^k) - \frac{\alpha_1^k}{\gamma_1^k}(U^k - U^{k-1}), U - U^k \rangle \right.$$
$$\left. + \frac{1}{\gamma_1^k} \mathbf{D}_{h_1}((U, V^k), (U^k, V^k)) + g_1(U) \right\}$$
$$= \arg\min_{U \geq 0} \left\{ \frac{1}{\gamma_1^k}\langle \gamma_1^k \nabla_U f(U^k, V^k) - \nabla_U h_1(U^k, V^k) - \alpha_1^k(U^k - U^{k-1}), U \rangle + \frac{1}{\gamma_1^k} h_1(U, V^k) \right\}.$$
$$(4.5)$$

By [44, Corollary 3.5], the normal cone of the nonnegativity constraint $U \geq 0$ is $\mathcal{N}_{U \geq 0}(U^k) = \left\{ P \in \mathbb{R}^{m \times r} \mid U^k \odot P = 0, \ P \leq 0 \right\}$ where $U^k \odot P$ denotes the *Hadamard products* given pointwise by $(U^k \odot P)_{ij} := U_{ij}^k P_{ij}$ for $i \in 1, \ldots, m$ and $j \in 1, \ldots, r$. The first-order optimality conditions for the subproblem (4.5) yields that $G^k - (a_1\|U^{k+1}\|_F^2\|V^k\|_F^2 + b_1(\|X\|_F\|V^k\|_F + \varepsilon_1))U^{k+1} \in \mathcal{N}_{U \geq 0}(U^{k+1})$.

We now consider two cases: (i) $G_{ij} \leq 0$; (ii) $G_{ij} > 0$. In Case (i), we have

$$P_{ij} = G_{ij}^k - (a_1\|U^{k+1}\|_F^2\|V^k\|_F^2 + b_1(\|X\|_F\|V^k\|_F + \varepsilon_1))U_{ij}^{k+1} \leq 0,$$

hence $U_{ij}^{k+1} = 0$. In Case (ii), if $U_{ij}^{k+1} = 0$, then $P_{ij} = G_{ij}^k > 0$, which contradicts $P \leq 0$; hence $G_{ij}^k - (a_1\|U^{k+1}\|_F^2\|V^k\|_F^2 + b_1(\|X\|_F\|V^k\|_F + \varepsilon_1))U_{ij}^{k+1} = 0$. Combining both cases, we get $(a_1\|U^{k+1}\|_F^2\|V^k\|_F^2 + b_1(\|X\|_F\|V^k\|_F + \varepsilon_1))U^{k+1} = \mathbf{Proj}_{G \geq 0}(G^k)$. Denote $t_k = a_1\|U^{k+1}\|_F^2\|V^k\|_F^2 + b_1\|X\|_F\|V^k\|_F$, then $\|U^{k+1}\|_F^2 = (t_k - b_1\|X\|_F\|V^k\|_F)/(a_1\|V^k\|_F^2)$. We have $t_k^3 - b_1\|X\|_F\|V^k\|_F t_k^2 - a_1\|V^k\|_F^2\|\mathbf{Proj}_{G \geq 0}(G^k)\|_F^2 = 0$. Note that the third order polynomial equation $y^2(y - a) = c$ has the unique real solution $y = \frac{a}{3} + \sqrt[3]{\frac{c + \sqrt{\Delta}}{2} + \frac{a^3}{27}} + \sqrt[3]{\frac{c - \sqrt{\Delta}}{2} + \frac{a^3}{27}}$, where $\Delta = c^2 + \frac{4}{27}ca^3$. Then we get (4.3). Finally, the result follows from $U^{k+1} = \frac{\mathbf{Proj}_{G \geq 0}(G^k)}{t_k}$.

By setting $g_2 := \delta_{V \geq 0}$ and invoking (3.5), we get

$$
\begin{aligned}
V^{k+1} &= \arg\min_{V \in \mathbb{R}^{r \times r}} \Big\{ \langle \nabla_V f(U^{k+1}, V^k) - \tfrac{\alpha_2^k}{\gamma_2^k}(V^k - V^{k-1}), V - V^k \rangle \\
&\qquad + \tfrac{1}{\gamma_2^k} \mathbf{D}_{h_2}((U^{k+1}, V), (U^{k+1}, V^k)) + g_2(V) \Big\} \\
&= \arg\min_{V \geq 0} \tfrac{1}{\gamma_2^k} \langle \gamma_2^k \nabla_V f(U^{k+1}, V^k) - \alpha_2^k(V^k - V^{k-1}) - \nabla h_2(U^{k+1}, V^k), V \rangle + \tfrac{1}{\gamma_2^k} h_2(U^{k+1}, V) \\
&= \arg\min_{V \geq 0} \Big\{ \Big\| V - \tfrac{1}{a_2 \|U^{k+1}\|^4 + \varepsilon_2} (\alpha_2^k(V^k - V^{k-1}) + \nabla h_2(U^{k+1}, V^k) - \gamma_2^k \nabla_V f(U^{k+1}, V^k)) \Big\|_F^2 \Big\} \\
&= \mathbf{Proj}_{V \geq 0} \Big( V^k - \tfrac{1}{\eta_k}(\alpha_2^k(V^k - V^{k-1}) - \gamma_2^k \nabla_V f(U^{k+1}, V^k)) \Big),
\end{aligned}
$$

which proves (4.4).                                                                                                          $\square$

## 5. Final remarks

The descent lemma is a key factor for analyzing the first-order methods in both Euclidean and non-Euclidean settings. Owing to the notion of block relative smoothness, it was shown that the descent lemma is still valid for each block of variables for structured nonsmooth nonconvex problems with non-Lipschitz gradients. Based on this development, `BIBPA` was introduced to deal with such problems, and it was shown to be globally convergent for KŁ functions and its convergence rate was also studied. Besides, it was shown that the objective of the symmetric nonnegative matrix tri-factorization (SymTriNMF) problem is block relatively smooth, and the corresponding subproblems can be solved in closed forms. To our knowledge, `BIBPA` is the first algorithm with rigorous theoretical guarantee of convergence for this problem. We emphasize that the main objective of this paper is to provide a theoretical and algorithmic framework that can handle block structured nonsmooth nonconvex problems under the block relative smoothness assumption. Hence, a comprehensive numerical experiments for such structured problems are postponed to a future work.

## References

[1] Ahookhosh, M.: Accelerated first-order methods for large-scale convex optimization: nearly optimal complexity under strong convexity. Math. Methods of Operations Research **89**(3), 319–353 (2019)

[2] Ahookhosh, M., Hien, L.T.K., Gillis, N., Patrinos, P.: Multi-block Bregman proximal alternating linearized minimization and its application to sparse orthogonal nonnegative matrix factorization. arXiv:1908.01402 (2019)

[3] Ahookhosh, M., Themelis, A., Patrinos, P.: A bregman forward-backward linesearch algorithm for nonconvex composite optimization: superlinear convergence to nonisolated local minima. arXiv:1905.11904 (2019)

[4] Airoldi, E.M., Blei, D.M., Fienberg, S.E., Xing, E.P.: Mixed membership stochastic blockmodels. Journal of Machine Learning Research **9**(Sep), 1981–2014 (2008)

[5] Attouch, H., Bolte, J., Redont, P., Soubeyran, A.: 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)

[6] Attouch, H., Bolte, J., Svaiter, B.F.: Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized Gauss-Seidel methods. Mathematical Programming **137**(1), 91–129 (2013)

[7] Attouch, H., Redont, P., Soubeyran, A.: A new class of alternating proximal minimization algorithms with costs-to-move. SIAM Journal on Optimization **18**(3), 1061–1081 (2007)

[8] Auslender, A.: Optimisation méthodes numériques. 1976. Mason, Paris (1976)

[9] Bauschke, H.H., Bolte, J., Chen, J., Teboulle, M., Wang, X.: On linear convergence of non-euclidean gradient methods without strong convexity and lipschitz gradient continuity. Journal of Optimization Theory and Applications **182**(3), 1068–1087 (2019)

[10] Bauschke, H.H., Bolte, J., Teboulle, M.: A descent lemma beyond Lipschitz gradient continuity: first-order methods revisited and applications. Mathematics of Operations Research **42**(2), 330–348 (2016)

[11] Bauschke, H.H., Combettes, P.L.: Convex analysis and monotone operator theory in Hilbert spaces. CMS Books in Mathematics. Springer (2017). DOI 10.1007/978-3-319-48311-5

[12] Beck, A.: First-Order Methods in Optimization. Society for Industrial and Applied Mathematics, Philadelphia, PA (2017). DOI 10.1137/1.9781611974997

[13] Beck, A., Pauwels, E., Sabach, S.: The cyclic block conditional gradient method for convex optimization problems. SIAM Journal on Optimization **25**(4), 2024–2049 (2015)

[14] Beck, A., Tetruashvili, L.: On the convergence of block coordinate descent type methods. SIAM journal on Optimization **23**(4), 2037–2060 (2013)

[15] Bertsekas, D.P., Tsitsiklis, J.N.: Parallel and Distributed Computation: Numerical Methods. Prentice-Hall, Inc. (1989)

[16] Bolte, J., Daniilidis, A., Lewis, A.: The Łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. SIAM Journal on Optimization **17**(4), 1205–1223 (2007)

[17] Bolte, J., Daniilidis, A., Lewis, A., Shiota, M.: Clarke subgradients of stratifiable functions. SIAM Journal on Optimization **18**(2), 556–572 (2007)

[18] Bolte, J., Daniilidis, A., Ley, O., Mazet, L.: Characterizations of Łojasiewicz inequalities: subgradient flows, talweg, convexity. Transactions of the American Mathematical Society **362**(6), 3319–3363 (2010)

[19] Bolte, J., Sabach, S., Teboulle, M.: Proximal alternating linearized minimization for nonconvex and nonsmooth problems. Mathematical Programming **146**(1–2), 459–494 (2014)

[20] Bolte, J., Sabach, S., Teboulle, M., Vaisbourd, Y.: First order methods beyond convexity and Lipschitz gradient continuity with applications to quadratic inverse problems. SIAM Journal on Optimization **28**(3), 2131–2151 (2018)

[21] Boţ, R.I., Csetnek, E.R., László, S.C.: 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)

[22] Bot, R.I., Nguyen, D.K.: The proximal alternating direction method of multipliers in the nonconvex setting: convergence analysis and rates. arXiv:1801.01994 (2018)

[23] Bregman, L.M.: The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. USSR Computational Mathematics and Mathematical Physics **7**(3), 200–217 (1967)

[24] Cichocki, A., Zdunek, R., Phan, A.H., Amari, S.i.: Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation. John Wiley & Sons (2009)

[25] Van den Dries, L.: Tame Topology and o-Minimal Structures, vol. 248. Cambridge university press (1998)

[26] Frankel, P., Garrigos, G., Peypouquet, J.: Splitting methods with variable metric for Kurdyka–Łojasiewicz functions and general convergence rates. Journal of Optimization Theory and Applications **165**(3), 874–900 (2015)

[27] Gillis, N.: The why and how of nonnegative matrix factorization. Regularization, optimization, kernels, and support vector machines **12**(257), 257–291 (2014)

[28] Grippo, L., Sciandrone, M.: On the convergence of the block nonlinear Gauss–Seidel method under convex constraints. Operations Research Letters **26**(3), 127–136 (2000)

[29] Gutman, D.H., Peña, J.F.: Perturbed fenchel duality and first-order methods. aarXiv:1812.10198 (2018)

[30] Hanzely, F., Richtárik, P.: Fastest rates for stochastic mirror descent methods. arXiv preprint arXiv:1803.07374 (2018)

[31] Hanzely, F., Richtarik, P., Xiao, L.: Accelerated Bregman proximal gradient methods for relatively smooth convex optimization. arXiv:1808.03045 (2018)

[32] Kurdyka, K.: On gradients of functions definable in o-minimal structures. Annales de l'institut Fourier **48**(3), 769–783 (1998)

[33] Łojasiewicz, S.: Une propriété topologique des sous-ensembles analytiques réels. Les équations aux dérivées partielles pp. 87–89 (1963)

[34] Łojasiewicz, S.: Sur la géométrie semi- et sous- analytique. Annales de l'institut Fourier **43**(5), 1575–1595 (1993)

[35] Lu, H., Freund, R.M., Nesterov, Y.: Relatively smooth convex optimization by first-order methods, and applications. SIAM Journal on Optimization **28**(1), 333–354 (2018)

[36] Nesterov, Y.: Introductory lectures on convex optimization: A basic course, vol. 87. Springer (2003)

[37] Nesterov, Y.: Efficiency of coordinate descent methods on huge-scale optimization problems. SIAM Journal on Optimization **22**(2), 341–362 (2012)

[38] Nesterov, Y.: Universal gradient methods for convex optimization problems. Mathematical Programming **152**(1-2), 381–404 (2015)

[39] Ochs, P.: Unifying abstract inexact convergence theorems and block coordinate variable metric ipiano. SIAM Journal on Optimization **29**(1), 541–570 (2019)

[40] Ochs, P., Chen, Y., Brox, T., Pock, T.: iPiano: Inertial proximal algorithm for nonconvex optimization. SIAM Journal on Imaging Sciences **7**(2), 1388–1419 (2014)

[41] Pock, T., Sabach, S.: Inertial proximal alternating linearized minimization (iPALM) for nonconvex and nonsmooth problems. SIAM Journal on Imaging Sciences **9**(4), 1756–1787 (2016)

[42] Rockafellar, R.T., Wets, R.J.B.: Variational Analysis, vol. 317. Springer Science (2011)

[43] Shefi, R., Teboulle, M.: On the rate of convergence of the proximal alternating linearized minimization algorithm for convex problems. EURO Journal on Computational Optimization **4**(1), 27–46 (2016)

[44] Tam, M.K.: Regularity properties of non-negative sparsity sets. Journal of Mathematical Analysis and Applications **447**(2), 758–777 (2017)

[45] Teboulle, M.: A simplified view of first order methods for optimization. Math. Prog. pp. 1–30 (2018)

[46] Tseng, P.: Convergence of a block coordinate descent method for nondifferentiable minimization. Journal of OptimizationTheory and Applications **109**(3), 475–494 (2001)

[47] Tseng, P., Yun, S.: A coordinate gradient descent method for nonsmooth separable minimization. Mathematical Programming **117**(1-2), 387–423 (2009)

[48] Van Nguyen, Q.: Forward-backward splitting with bregman distances. Vietnam Journal of Mathematics **45**(3), 519–539 (2017)

[49] Wang, H., Huang, H., Ding, C.: Simultaneous clustering of multi-type relational data via symmetric non-negative matrix tri-factorization. In: Proceedings of the 20th ACM CIKM'11, pp. 279–284 (2011)

[50] Wang, X., Yuan, X., Zeng, S., Zhang, J., Zhou, J.: Block coordinate proximal gradient method for nonconvex optimization problems: convergence analysis. http://www.optimization-online.org/DB_HTML/2018/04/6573.html (2018)

[51] Zhang, X., Zhang, H., Peng, W.: Inertial bregman proximal gradient algorithm for nonconvex problem with smooth adaptable property. arXiv preprint arXiv:1904.04436 (2019)

[52] Zhang, Y., Yeung, D.Y.: Overlapping community detection via bounded nonnegative matrix tri-factorization. In: Proceedings of the 18th ACM SIGKDD, pp. 606–614 (2012)