

Multi-block Bregman proximal alternating linearized minimization and its application to orthogonal nonnegative matrix factorization

Masoud Ahookhosh¹ · Le Thi Khanh Hien² · Nicolas Gillis² · Panagiotis Patrinos³

Received: 16 February 2020 / Accepted: 28 May 2021 / Published online: 9 June 2021 © The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2021

Abstract

We introduce and analyze BPALM and A-BPALM, two multi-block proximal alternating linearized minimization algorithms using Bregman distances for solving structured nonconvex problems. The objective function is the sum of a multi-block relatively smooth function (i.e., relatively smooth by fixing all the blocks except one) and block separable (nonsmooth) nonconvex functions. The sequences generated by our algorithms are subsequentially convergent to critical points of the objective function, while they are globally convergent under the KL inequality assumption. Moreover, the rate of convergence is further analyzed for functions satisfying the Łojasiewicz's gradient inequality. We apply this framework to orthogonal nonnegative matrix factorization (ONMF) that satisfies all of our assumptions and the related subproblems are solved in closed forms, where some preliminary numerical results are reported.

Keywords Nonsmooth nonconvex optimization · Proximal alternating linearized minimization · Bregman distance · Multi-block relative smoothness · KL inequality · Orthogonal nonnegative matrix factorization

Mathematics Subject Classification 90C06 · 90C25 · 90C26 · 49J52 · 49J53

1 Introduction

For a nonempty, convex, and open set $C \subseteq \mathbb{R}^n$, we consider the structured nonsmooth nonconvex minimization problem

Masoud Ahookhosh masoud.ahookhosh@uantwerp.be

Extended author information available on the last page of the article

$$\underset{\mathbf{x}=(x_1,\dots,x_N)\in\overline{C}}{\text{minimize}} \quad \varphi(\mathbf{x}) \equiv f(\mathbf{x}) + \sum_{i=1}^N g_i(x_i), \tag{1.1}$$

where we will assume the following hypotheses (see Sect. 3 for details):

Assumption 1 (requirements for composite minimization (1.1))

A1 $g_i : \mathbb{R}^{n_i} \to \overline{\mathbb{R}} := \mathbb{R} \cup \{\infty\}$ is proper and lower semicontinuous (lsc); A2 $h : \mathbb{R}^n \to \overline{\mathbb{R}}$ is proper, closed, and *multi-block* Legendre; A3 $f : \mathbb{R}^n \to \overline{\mathbb{R}}$ is $C^1(\text{int dom } h)$ and (L_1, \dots, L_N) -smooth relative to h; here $n = \sum_{i=1}^N n_i$; A4 $\arg \min \varphi \neq \emptyset, \overline{C} := \overline{\operatorname{dom } h}, C := \operatorname{int dom } h, g := \sum_{i=1}^N g_i$, and $C \subseteq \operatorname{dom } g$.

There is a huge number of algorithmic studies around solving the optimization problems of the form (1.1). Among all of such methodologies, we are interested in the class of *alternating minimization* algorithms such as block coordinate descent [15, 18, 42, 48, 53, 54, 61, 62], block coordinate [30, 32, 41], and Gauss-Seidel methods [11, 19, 36], which assume that all blocks are fixed except one and solve the corresponding auxiliary problem with respect to this block, update the latter block, and continue with the others. In particular, the proximal alternating minimization has received much attention in the last few years; see for example [6–10, 16]. Recently, the proximal alternating linearized minimization and its variants has been developed to handle (1.1); see for example [23, 51, 57].

Traditionally, the Lipschitz (Hölder) continuity of partial gradients of f in (1.1) is a necessary tool for providing the convergence analysis of optimization algorithms; see, e.g., [23, 51]. It is, however, well-known that it is not the Lipschitz (Hölder) continuity of gradients playing a key role in such analysis, but one of its consequences: an upper estimation of f including a Bregman distance called *descent lemma*; cf. [13, 46]. This idea is central to convergence analysis of many optimization schemes requiring such an upper estimation; see, e.g., [1–3, 12, 13, 24, 37, 38, 46, 59, 60]. In this paper, we propose a multi-block extension of the descent lemma given in [13, 46] and propose a Bregman proximal alternating linearized minimization (BPALM) algorithm and its adaptive version (A-BPALM) for (1.1).

1.1 Contribution

Our contribution is summarized as follows:

(i) (Relative smoothness with possibly nonconvex kernel) An extension of the relative smoothness [13, 46] for problems with block separable structure entails an application of multi-block convex kernel functions (see Definition 3.1) which is not necessarily a jointly convex function. This paves the way toward the development of non-Euclidean alternating minimization methodologies.

- (ii) (*Bregman proximal alternating linearized minimization*) We introduce BPALM, a multi-block generalization of the proximal alternating linearized minimization (PALM) [23] using Bregman distances, and its adaptive version (A-BPALM). To do so, we extend the notion of relative smoothness [13, 46] to its multi-block counterpart to support a structured problem of the form (1.1). Owing to multi-block relative smoothness of *f*, our algorithm does not need to know the local Lipschitz moduli of partial gradients $\nabla_i f$ (i = 1, ..., N) and their lower and upper bounds, which are hard to provide in practice.
- (iii) (Efficient framework for ONMF) Exploiting a suitable kernel for Bregman distance, it turns out that the objective of orthogonal nonnegative matrix factorization (ONMF) is multi-block relatively smooth, and the subproblems of our algorithms are solved in closed forms making them suitable for large-scale data analysis problems. To the best of our knowledge, BPALM and A-BPALM are the first algorithms with rigorous convergence theory for ONMF.

1.2 Related works

Closely related to our framework, there are two papers [43, 63]. However, we notice that [63] uses a sum separable kernel function which is a special case of our multiblock kernel functions (see Example 3.2), and the paper provides a subsequential convergence theory. Regarding [43], an algorithm (named B-PALM) was proposed that is a special case of our BPALM when N = 2, $g_1 = g_2 = 0$, and $f \in C^2$. Hence, the proposed algorithms in [43] are not applicable to ONMF, which requires non-negativity constraints. We stress that involving the block separable nonsmooth non-convex functions g_i and considering N > 2 make our analysis different from those of [43].

1.3 Organization

This paper has five sections, besides this introductory section. In Sect. 2, we describe some preliminaries, and in Sect. 3, we introduce the notion of multi-block relative smoothness, and verify the fundamental properties of Bregman proximal alternating linearized mapping. In Sect. 4, we introduce BPALM and A-BPALM and investigate their convergence analysis. In Sect. 5, we show that ONMF satisfies our assumptions, the related subproblems are solved in closed forms and report our numerical results. Finally, Sect. 6 delivers some conclusions.

2 Notation and preliminaries

We denote by $\overline{\mathbb{R}} := \mathbb{R} \cup \{\infty\}$ the extended-real line. For the identity matrix I_n , we set $U_i \in \mathbb{R}^{n \times n_i}$ such that $I_n = (U_1, \dots, U_N) \in \mathbb{R}^{n \times n}$. For notation clarity, we will use bold lower-case letters (e.g., x, y, z) for vectors in $\mathbb{R}^{\sum n_i}$ and use normal lower-case letters (e.g., z, x_i, y_i) for vectors in \mathbb{R}^{n_i} . The open ball of radius r > 0 centered in $x \in \mathbb{R}^p$ is denoted as $\mathbf{B}(x;r)$. For a set $C \subseteq \mathbb{R}^p$, \overline{C} denotes its closure. The set

of cluster points of $(x^k)_{k \in \mathbb{N}}$ is denoted as $\omega(x^0)$. A function $f : \mathbb{R}^p \to \overline{\mathbb{R}}$ is proper $f \not\equiv \infty$, in which case its *domain* is defined as the set **dom** $f := \{x \in \mathbb{R}^p\}[f(x) < \infty]$. For $\alpha \in \mathbb{R}$, $[f \le \alpha] := \{x \in \mathbb{R}^p\}[f(x) \le \alpha]$ is the α -(sub)level set of f; $[f \ge \alpha]$ and $[f = \alpha]$ are defined similarly. We say that f is level bounded if $[f \le \alpha]$ is bounded for all $\alpha \in \mathbb{R}$. A vector $v \in \partial f(x)$ is a subgradient of f at x, and the set of all such vectors is called the subdifferential $\partial f(x)$ [55, Definition 8.3], i.e.

$$\partial f(x) = \{ v \in \mathbb{R}^p \} [\exists (x^k, v^k)_{k \in \mathbb{N}} \text{ s.t. } x^k \to x, f(x^k) \to f(x), \ \widehat{\partial} f(x^k) \ni v^k \to v],$$

and $\partial f(x)$ is the set of *regular subgradients* of *f* at *x*, namely

$$\widehat{\partial}f(x) = \{v \in \mathbb{R}^p\}[f(z) \ge f(x) + \langle v, z - x \rangle + o(||z - x||), \forall z \in \mathbb{R}^p].$$

The Fenchel conjugate function $f^* : \mathbb{R}^p \to \overline{\mathbb{R}}$ of $f : \mathbb{R}^p \to \overline{\mathbb{R}}$ is given by

$$f^*(z^*) = \sup\{\langle z^*, x \rangle - f(x)\} [x \in \mathbb{R}^p].$$

For a given nonempty closed convex set $S \subseteq \mathbb{R}^p$, the function $\delta_S : \mathbb{R}^p \to \mathbb{R}$ denotes the *indicator function*, namely $\delta_S(x) = 0$ if $x \in S$ and $\delta_S(x) = \infty$ otherwise. Moreover, **Proj**_S : $\mathbb{R}^p \to \mathbb{R}^p$ denotes the *projection function* given by **Proj**_S(x) = **argmin** ||z - x||.

2.1 Bregman proximal mapping

For a kernel function $h : \mathbb{R}^n \to \overline{\mathbb{R}}$ and a proper lower semicontinuous function $g : \mathbb{R}^n \to \overline{\mathbb{R}}$, the *Bregman proximal mapping* is given by

$$\mathbf{prox}_{\gamma g}^{h}(x) := \underset{z \in \mathbb{R}^{n}}{\operatorname{arg\,min}} \{ g(z) + \frac{1}{\gamma} \mathbf{D}_{h}(z, x) \}.$$
(2.1)

which is a generalization of the classical one by using the Bregman distance (3.2) in place of the Euclidean distance; see, e.g., [28] and references therein. We note that

$$\mathbf{prox}_{\gamma g}^{h}(x) = \{ y \in \mathbf{dom} \ g \cap \mathbf{dom} \ h \mid g(y) + \frac{1}{\gamma} \mathbf{D}_{h}(y, x) = \mathbf{inf}_{z} \{ g(z) + \frac{1}{\gamma} \mathbf{D}_{h}(z, x) \} < +\infty \},\$$

which implies **dom prox**^{*h*}_{γg} \subset **intdom** *h*, **rangeprox**^{*h*}_{γg} \subset **dom** *g* \cap **dom** *h*. The function *g* is *h*-*prox*-*bounded* if there exists $\gamma > 0$ such that $\inf_{z} \{g(z) + \frac{1}{\gamma} \mathbf{D}_{h}(z, x)\} > -\infty$ for some $x \in \mathbb{R}^{n}$; see for example [3].

2.2 Kurdyka-Łojasiewicz (KL) function

Let us consider the class of functions that are satisfying the celebrated Kurdyka-Łojasiewicz inequality, which we present next.

Definition 2.1 (KL property) A proper and lsc function $\varphi : \mathbb{R}^{n_1} \times ... \times \mathbb{R}^{n_N} \to \mathbb{R}$ has the *Kurdyka-Łojasiewicz* property (KL property) at $x^* \in \operatorname{dom} \varphi$ if there exist a

concave *desingularizing function* ψ : $[0, \eta] \rightarrow [0, +\infty[$ (for some $\eta > 0$) and neighborhood **B**($x^*; \epsilon$) with $\epsilon > 0$, such that

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

$$\psi'(\varphi(\mathbf{x}) - \varphi(\mathbf{x}^*))\mathbf{dist}(0, \partial\varphi(\mathbf{x})) \ge 1.$$
(2.2)

If φ satisfies the KL property at each point of **dom** $\partial \varphi$, then φ is called a KL function.

The first inequality of this type is given in the seminal work of Łojasiewicz [44, 45] for analytic functions, which we nowadays call Łojasiewicz's gradient inequality. Later, Kurdyka [40] showed that this inequality is valid for C^1 functions whose graph belong to an *o-minimal structure* (see its definition in [31]). The first extensions of the KL property to nonsmooth functions were given by Bolte et al. [20–22].

The following two facts constitute the crucial steps for establishing the global convergence of the algorithms given in Sect. 4.

Fact 2.2 (uniformized KL property) [23, Lemma 6] Let Ω be a compact set and $\zeta : \mathbb{R}^d \to \overline{\mathbb{R}}$ be a proper and lower semicontinuous function. Assume that ζ is constant on Ω and satisfies the KL property at each point of Ω . Then, there exist $\varepsilon > 0$, $\eta > 0$, and a desingularizing function ψ such that for $\overline{u} \in \Omega$ and all u in the intersection

$$\{u \in \mathbb{R}^d \mid \operatorname{dist}(u, \Omega) < \varepsilon\} \cap [\zeta(\overline{u}) < \zeta(u) < \zeta(\overline{u}) + \eta]$$

we have

$$\psi'(\zeta(u) - \zeta(\overline{u}))$$
dist $(0, \partial \zeta(u)) \ge 1$.

Fact 2.3 [25, Lemma 2.3] Let $(a_k)_{k \in \mathbb{N}}$ and $(b_k)_{k \in \mathbb{N}}$ be sequences in $[0, +\infty)$ such that $\sum_{k=1}^{\infty} b_k < \infty$ and $a_{k+1} \le \alpha a_k + b_k$ for all $k \in \mathbb{N}$ and some $\alpha < 1$. Then, $\sum_{k=1}^{\infty} a_k < \infty$.

3 Multi-block Bregman proximal mapping

In this section, we first establish the notion of multi-block relative smoothness, which is an extension of the relative smoothness [13, 46] for problems with block structure. Afterwards, we introduce Bregman alternating linearized mapping and study some of its basic properties.

3.1 Multi-block relative smoothness

In order to extend the definition of Bregman distances for the multi-block problem (1.1), we first need to introduce the notion of *multi-block kernel* functions, which will coincide with the standard one (cf. [3, Definition 2.1]) if N = 1.

Definition 3.1 (multi-block convexity and kernel function) Let $h : \mathbb{R}^n \to \overline{\mathbb{R}}$ be a proper and lsc function with **int dom** $h \neq \emptyset$ and such that $h \in C^1(\text{int dom } h)$. For a fixed vector $\mathbf{x} \in \mathbb{R}^n$, we define the function $h_x^i : \mathbb{R}^{n_i} \to \overline{\mathbb{R}}$ given by

$$h_{\mathbf{x}}^{i}(z) := h(\mathbf{x} + U_{i}(z - x_{i})).$$
 (3.1)

Then, we say that *h* is

- (i) multi-block (strongly/strictly) convex if the function $h_x^i(\cdot)$ is (strongly/strictly) convex for all $x \in \text{dom } h$ and i = 1, ..., N;
- (ii) multi-block locally strongly convex around $\mathbf{x}^{\star} = (x_1^{\star}, \dots, x_N^{\star})$ if, for $i = 1, \dots, N$, there exists $\delta > 0$ and $\sigma_h^i > 0$ such that

$$h_{\boldsymbol{x}}^{i}(x_{i}) \geq h_{\boldsymbol{x}}^{i}(y_{i}) + \langle \nabla_{i}h(\boldsymbol{x} + U_{i}(y_{i} - x_{i})), x_{i} - y_{i} \rangle + \frac{\sigma_{h}^{i}}{2} \|x_{i} - y_{i}\|^{2} \quad \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbf{B}(\boldsymbol{x}^{\star}; \delta);$$

- (iii) a multi-block kernel function if h is multi-block convex and $h_x^i(\cdot)$ is 1-coercive for all $x \in \operatorname{dom} h$ and $i = 1, \ldots, N$, i.e., $\lim_{\|z\|\to\infty} \frac{h_x^i(z)}{\|z\|} = \infty$;
- (iv) essentially smooth, if for every sequence $(\mathbf{x}^k)_{k \in \mathbb{N}} \subseteq$ int dom h converging to a boundary point of dom h, we have $\|\nabla h(\mathbf{x}^k)\| \to \infty$;
- (v) of *multi-block Legendre* type if it is essentially smooth and multi-block strictly convex.

Example 3.2 (popular kernel functions) There are many kernel functions satisfying the conditions given in Definition 3.1(iii). For example, for N = 1, energy, Boltzmann-Shannon entropy, Fermi-Dirac entropy (cf. [14, Example 2.3]) and several examples in [46, Section 2]; and for N = 2 see two examples in [43, Section 2]. Two important classes of multi-block kernels are *sum separable kernels*, i.e., $h(x_1, \ldots, x_N) = h_1(x_1) + \ldots + h_N(x_N)$, and *product separable kernels*, i.e., $h(x_1, \ldots, x_N) = h_1(x_1) \times \ldots \times h_N(x_N)$, see such a kernel for ONMF in Proposition 5.1.

We now give the definition of *Bregman distances* (cf. [27]) for multi-block kernels.

Definition 3.3 (Bregman distance) For a multi-block kernel function *h*, the *Bregman distance* $\mathbf{D}_h : \mathbb{R}^n \times \mathbb{R}^n \to \overline{\mathbb{R}}$ is given by

$$\mathbf{D}_{h}(\mathbf{y}, \mathbf{x}) := \begin{cases} h(\mathbf{y}) - h(\mathbf{x}) - \langle \nabla h(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle & \mathbf{y} \in \operatorname{dom} h, \mathbf{x} \in \operatorname{int} \operatorname{dom} h \\ \infty & \operatorname{otherwise.} \end{cases}$$
(3.2)

Fixing all blocks except the *i*-th one, the Bregman distance with respect to this block is given by

$$\mathbf{D}_{h}(\mathbf{x} + U_{i}(y_{i} - x_{i}), \mathbf{x}) = h(\mathbf{x} + U_{i}(y_{i} - x_{i})) - h(\mathbf{x}) - \langle \nabla h(\mathbf{x}), U_{i}(y_{i} - x_{i}) \rangle$$
$$= h_{\mathbf{x}}^{i}(y_{i}) - h_{\mathbf{x}}^{i}(x_{i}) - \langle \nabla_{i}h(\mathbf{x}), y_{i} - x_{i} \rangle,$$

Deringer

which measures the proximity between y and x with respect to the *i*-th block of variables. Moreover, the kernel h is multi-block convex if and only if $\mathbf{D}_h(\mathbf{x} + U_i(y_i - x_i), \mathbf{x}) \ge 0$ for all $\mathbf{y} \in \mathbf{dom} h$ and $\mathbf{x} \in \mathbf{int} \mathbf{dom} h$ and i = 1, ..., N. Note that if h is multi-block strictly convex, then $\mathbf{D}_h(\mathbf{x} + U_i(y_i - x_i), \mathbf{x}) = 0$ (i = 1, ..., N) if and only if $x_i = y_i$.

We are now in a position to present the notion of *multi-block relative smoothness*, which is the central tool for our analysis in Sect. 4.

Definition 3.4 (multi-block relative smoothness) Let $h : \mathbb{R}^n \to \overline{\mathbb{R}}$ be a multi-block kernel and let $f : \mathbb{R}^n \to \overline{\mathbb{R}}$ be a proper and lower semicontinuous function. If there exists $L_i > 0$ (i = 1, ..., N) such that the functions $\phi_x^i : \mathbb{R}^{n_i} \to \overline{\mathbb{R}}$ given by

$$\phi_{\mathbf{x}}^{i}(z) := L_{i}h(\mathbf{x} + U_{i}(z - x_{i})) - f(\mathbf{x} + U_{i}(z - x_{i}))$$

are convex for all $x, x + U_i(z - x_i) \in \text{int dom } h$ and i = 1, ..., N, then f is called $(L_1, ..., L_N)$ -smooth relative to h.

Note that if N = 1, the multi-block relative smoothness is reduced to standard relative smoothness, which was introduced only recently in [13, 46]. In the Euclidean case, relative smoothness is equivalent to the one-sided descent lemma, which is implied by the Lipschitz continuity of gradients ∇f (i.e., $\frac{L}{2} || \cdot ||^2 - f$ is convex). Therefore, the relative smoothness of f generalizes the notions of Lipschitz smoothness using Bregman distances. If N = 2, this definition will be reduced to the relative bi-smoothness given in [43] for $h, f \in C^2$.

We next characterize the notion of multi-block relative smoothness.

Proposition 3.5 (characterization of multi-block relative smoothness) Let $h : \mathbb{R}^n \to \overline{\mathbb{R}}$ be a multi-block kernel and let $f : \mathbb{R}^n \to \overline{\mathbb{R}}$ be a proper lower semicontinuous function and $f \in C^1$. Then, the following statements are equivalent:

- (a) f is (L_1, \ldots, L_N) -smooth relative to h;
- (b) for all $(\mathbf{x}, \mathbf{y}) \in \operatorname{int} \operatorname{dom} h \times \operatorname{int} \operatorname{dom} h$ and $i = 1, \dots, N$,

$$f(\mathbf{x} + U_i(y_i - x_i)) \le f(\mathbf{x}) + \langle \nabla_i f(\mathbf{x}), y_i - x_i \rangle + L_i \mathbf{D}_h(\mathbf{x} + U_i(y_i - x_i), \mathbf{x}); \quad (3.3)$$

(c) for all $(\mathbf{x}, \mathbf{y}) \in \mathbf{int} \operatorname{dom} h \times \mathbf{int} \operatorname{dom} h$ and $i = 1, \dots, N$,

$$\langle \nabla_i f(\mathbf{x}) - \nabla_i f(\mathbf{y}), x_i - y_i \rangle \le L_i \langle \nabla_i h(\mathbf{x}) - \nabla_i h(\mathbf{y}), x_i - y_i \rangle;$$
(3.4)

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

$$L_i \nabla_{x_i x_i}^2 h(\mathbf{x}) - \nabla_{x_i x_i}^2 f(\mathbf{x}) \ge 0, \qquad (3.5)$$

for i = 1, ..., N and for all $x \in int \text{ dom } h$.

Proof Fixing all the blocks except one of them, the results can be concluded in the same way as [46, Proposition 1.1]. \Box

3.2 Bregman proximal alternating linearized mapping

The main objective of the remainder of this section is to provide a guarantee for existence of the Bregman proximal alternating linearized minimization mapping (see (3.9)) and to investigate its fundamental properties.

Let us begin with an extension of *h*-prox-boundedness to our multi-block setting, where both of them are extensions of the classical prox-boundedness condition; see, e.g., [55, Definition 1.23].

Definition 3.6 (multi-block *h*-prox-boundedness) A function $g : \mathbb{R}^n \to \overline{\mathbb{R}}$ is *multiblock h-prox-bounded* if for each $i \in \{1, ..., N\}$ there exists $\gamma_i > 0$ and $x \in \mathbb{R}^n$ such that

$$g^{h/\gamma_i}(\boldsymbol{x}) := \inf_{z \in \mathbb{R}^{n_i}} \{ g(\boldsymbol{x} + U_i(z - x_i)) + \frac{1}{\gamma_i} \mathbf{D}_h(\boldsymbol{x} + U_i(z - x_i), \boldsymbol{x}) \} > -\infty$$

The supremum of the set of all such γ_i is the threshold $\gamma_{i,p}^h$ of the multi-block *h*-proxboundedness, i.e.,

$$\gamma_{i,g}^h := \sup\{\gamma_i > 0 \mid \exists x \in \mathbb{R}^n \text{ s.t. } g^{h/\gamma_i}(x) > -\infty\}.$$
(3.6)

For the problem (1.1), we have $g = \sum_{i=1}^{N} g_i$ leading to

$$g^{h/\gamma_i}(\boldsymbol{x}) = \sum_{j \neq i} g_j(x_j) + \inf_{z \in \mathbb{R}^{n_i}} \{ g_i(z) + \frac{1}{\gamma_i} \mathbf{D}_h(\boldsymbol{x} + U_i(z - x_i), \boldsymbol{x}) \}.$$
(3.7)

If g is multi-block h-prox-bounded for $\overline{\gamma}_i > 0$, so is for all $\gamma_i \in (0, \overline{\gamma}_i)$. We next present equivalent conditions to this notion.

Proposition 3.7 (characterization of multi-block h-prox-boundedness) For a multi-block kernel function $h : \mathbb{R}^n \to \overline{\mathbb{R}}$ and proper and lsc functions $g_i : \mathbb{R}^{n_i} \to \overline{\mathbb{R}}$ with i = 1, ..., N, the following statements are equivalent:

- (a) g = ∑_{i=1}^N g_i is multi-block h-prox-bounded;
 (b) for all i = 1,..., N and hⁱ_x given in (3.1), g_i + r_ihⁱ_x is bounded below on ℝ^{n_i} for some $r_i \in \mathbb{R}$;
- (c) for all $i = 1, \dots, N$, $\liminf_{\|z\| \to \infty} \frac{g_i(z)}{h_i(z)} > -\infty$.

Proof We first show (a) \Leftrightarrow (b). Let us assume $g^{h/\gamma_i}(\mathbf{x}) > -\infty$, and let $r_i > \frac{1}{\gamma_i}$. Then, for all i = 1, ..., N, it holds that

$$g_i(z) + r_i h_{\mathbf{x}}^i(z) = g_i(z) + \frac{1}{\gamma_i} \mathbf{D}_h(\mathbf{x} + U_i(z - x_i), \mathbf{x}) + r_i h_{\mathbf{x}}^i(z) - \frac{1}{\gamma_i} \mathbf{D}_h(\mathbf{x} + U_i(z - x_i), \mathbf{x})$$

$$\geq g^{h/\gamma_i}(\mathbf{x}) - \sum_{j \neq i} g_j(x_j) + \frac{r_i \gamma_i - 1}{\gamma_i} h_{\mathbf{x}}^i(z) + \frac{1}{\gamma_i} (h(\mathbf{x}) + \langle \nabla_i h(\mathbf{x}), z - x_i \rangle) =: \tilde{g}_i(z),$$

for the function $\tilde{g}_i : \mathbb{R}^{n_i} \to \overline{\mathbb{R}}$. We notice that \tilde{g}_i is convex and coercive, and as such is lower bounded. Conversely, suppose that $\alpha_i := \inf g_i + r_i h_x^i > -\infty$. Then, from (3.7), we obtain

$$\begin{split} g^{h/\gamma_i}(\mathbf{x}) &= \sum_{j \neq i} g_j(x_j) + \min_{z \in \mathbb{R}^{n_i}} \{ g_i(z) + \frac{1}{\gamma_i} \mathbf{D}_h(\mathbf{x} + U_i(z - x_i), \mathbf{x}) \}, \\ &\geq \sum_{j \neq i} g_j(x_j) + \alpha_i + \inf_z \{ -r_i h^i_{\mathbf{x}}(z) + \frac{1}{\gamma_i} \mathbf{D}_h(\mathbf{x} + U_i(z - x_i), \mathbf{x}) \} \\ &\geq \sum_{j \neq i} g_j(x_j) + \alpha_i - \frac{1}{\gamma_i} h(\mathbf{x}) + \frac{1}{\gamma_i} \langle \nabla_i h(\mathbf{x}), x_i \rangle + \inf_z \{ \frac{1 - \gamma_i r_i}{\gamma_i} h^i_{\mathbf{x}}(z) - \frac{1}{\gamma_i} \langle \nabla_i h(\mathbf{x}), z \rangle \}, \end{split}$$

which is finite, owing to 1-coercivity of $z \mapsto \frac{1-\gamma_i r_i}{\gamma_i} h_x^i(z) - \frac{1}{\gamma_i} \langle \nabla_i h(x), z \rangle$.

We now show $(b) \Leftrightarrow (c)$. Suppose that $\alpha_i := \inf g_i + r_i h_x^i > -\infty$. Since $h_x^i(\cdot)$ is 1-coercive, we have

$$\liminf_{\|z\|\to\infty}\frac{g_i(z)}{h_x^i(z)}\geq -r_i+\liminf_{\|z\|\to\infty}\frac{\alpha_i}{h_x^i(z)}=-r_i>-\infty.$$

Conversely, suppose $\liminf_{\|z\|\to\infty} \frac{g_i(z)}{h_x^i(z)} > -\infty$. Then, there exist $\ell_i, M_i \in \mathbb{R}$ such that $\frac{g_i(z)}{h_x^i(z)} \ge \ell_i$ whenever $\|z\| \ge M_i$. Thus, it holds that

$$\inf_{\|z\| \ge M_i} g_i(z) + r_i h_x^i(z) \ge \inf_{\|x\| \ge M_i} h_x^i(z) (\ell_i + r_i) > -\infty,$$

where the last inequality follows from coercivity of h_x^i , i,e, h_x^i is level-bounded and consequently lower bounded (cf. [55, Corollary 1.10]). Since the function $g_i(z) + r_i h_x^i(z)$ is lower semicontinuous, it follows from [55, Corollary 1.10] that $\inf_{\|z\| \le M_i} g_i(z) + r_i h_x^i(z) > -\infty$. Therefore, we conclude that $g_i + r_i h_x^i$ is lower bounded on \mathbb{R}^n .

Let us now define the function $\mathcal{M}_{h/\gamma}$: $\mathbb{R}^n \times \mathbb{R}^n \to \overline{\mathbb{R}}$ as

$$\mathcal{M}_{h/\gamma}(z, \boldsymbol{x}) := f(\boldsymbol{x}) + \langle \nabla f(\boldsymbol{x}), \boldsymbol{z} - \boldsymbol{x} \rangle + \frac{1}{\gamma} \mathbf{D}_{h}(z, \boldsymbol{x}) + \sum_{i=1}^{N} g_{i}(z_{i})$$
(3.8)

and the set-valued *Bregman* proximal alternating linearized mapping $\mathbf{T}_{h/ri}^{i}$: $\mathbb{R}^{n} \Rightarrow \mathbb{R}^{n_{i}}$ as

$$\mathbf{T}_{h/\gamma i}^{i}(\boldsymbol{x}) := \operatorname*{arg\,min}_{z \in \mathbb{R}^{n_{i}}} \mathcal{M}_{h/\gamma_{i}}(\boldsymbol{x} + U_{i}(z - x_{i}), \boldsymbol{x}), \tag{3.9}$$

which reduces to the Bregman forward-backward splitting mapping if N = 1; cf. [3, 24].

Remark 3.8 (majorization model) Note that, for $x \in \text{int dom } h$, invoking Proposition 3.5(b) and Assumption 1, the multi-block (L_1, \ldots, L_N) -relative smoothness assumption of f entails a block majorization model

$$\begin{split} \varphi(\mathbf{x} + U_i(y_i - x_i)) &\leq f(\mathbf{x}) + \langle \nabla_i f(\mathbf{x}), y_i - x_i \rangle + L_i \mathbf{D}_h(\mathbf{x} + U_i(y_i - x_i), \mathbf{x}) + g_i(y_i) + \sum_{j \neq i} g_j(x_j) \\ &\leq f(\mathbf{x}) + \langle \nabla_i f(\mathbf{x}), y_i - x_i \rangle + \frac{1}{\gamma_i} \mathbf{D}_h(\mathbf{x} + U_i(y_i - x_i), \mathbf{x}) + g_i(y_i) + \sum_{j \neq i} g_j(x_j), \end{split}$$

for $\gamma_i \in (0, 1/L_i)$.

In the next lemma, we show that the cost function φ is monotonically decreasing by minimizing the model (3.8) with respect to each block of variables.

Lemma 3.9 (Bregman proximal alternating inequality) Let the conditions in Assumption 1 hold, and let $\overline{z} \in \mathbf{T}_{h/vi}^{i}(\mathbf{x})$ with $\gamma_{i} \in (0, 1/L_{i})$ and $\mathbf{x} \in \mathbf{int} \operatorname{dom} h$. Then,

$$\varphi(\mathbf{x} + U_i(\overline{z} - x_i)) \le \varphi(\mathbf{x}) - \frac{1 - \gamma_i L_i}{\gamma_i} \mathbf{D}_h(\mathbf{x} + U_i(\overline{z} - x_i), \mathbf{x}),$$
(3.10)

for all i = 1, ..., N.

Proof For $i \in \{1, ..., N\}$, (3.9) is simplified in the form

$$\mathbf{T}_{h/\gamma i}^{i}(\mathbf{x}) = \underset{z \in \mathbb{R}^{n_{i}}}{\arg\min} \{ \langle \nabla f(\mathbf{x}), U_{i}(z-x_{i}) \rangle + \frac{1}{\gamma_{i}} \mathbf{D}_{h}(\mathbf{x}+U_{i}(z-x_{i}),\mathbf{x}) + g_{i}(z) + \sum_{j \neq i}^{N} g_{j}(x_{j}) \}$$
$$= \underset{z \in \mathbb{R}^{n_{i}}}{\arg\min} \{ \langle \nabla_{i} f(\mathbf{x}), z-x_{i} \rangle + \frac{1}{\gamma_{i}} \mathbf{D}_{h}(\mathbf{x}+U_{i}(z-x_{i}),\mathbf{x}) + g_{i}(z) \}.$$
(3.11)

Considering $\overline{z} \in \mathbf{T}_{h/\gamma i}^{i}(\mathbf{x})$, we have

$$\langle \nabla_i f(\mathbf{x}), \overline{z} - x_i \rangle + \frac{1}{\gamma_i} \mathbf{D}_h(\mathbf{x} + U_i(\overline{z} - x_i), \mathbf{x}) + g_i(\overline{z}) \le g_i(x_i).$$

Since *f* is $(L_1, ..., L_N)$ -smooth relative to *h*, it follows from Proposition 3.5(b) for *x* and $y_i = \overline{z}$ that

$$\begin{split} f(\mathbf{x} + U_i(\overline{z} - x_i)) &\leq f(\mathbf{x}) + \langle \nabla_i f(\mathbf{x}), \overline{z} - x_i \rangle + L_i \mathbf{D}_h(\mathbf{x} + U_i(\overline{z} - x_i), \mathbf{x}) \\ &\leq f(\mathbf{x}) + L_i \mathbf{D}_h(\mathbf{x} + U_i(\overline{z} - x_i), \mathbf{x}) + g_i(x_i) - g_i(\overline{z}) - \frac{1}{\gamma_i} \mathbf{D}_h(\mathbf{x} + U_i(\overline{z} - x_i), \mathbf{x}) \\ &= f(\mathbf{x}) + g_i(x_i) - g_i(\overline{z}) - \frac{1 - \gamma_i L_i}{\gamma_i} \mathbf{D}_h(\mathbf{x} + U_i(\overline{z} - x_i), \mathbf{x}), \end{split}$$

giving (3.10).

Recall that a function $\vartheta : \mathbb{R}^n \times \mathbb{R}^m \to \overline{\mathbb{R}}$ with values $\vartheta(x, u)$ is *level-bounded* in x locally uniformly in u if for each $\overline{u} \in \mathbb{R}^m$ and $\alpha \in \mathbb{R}$ there is a neighborhood \mathcal{U} of \overline{u} along with a bounded set $B \subset \mathbb{R}^n$ such that $\{x \in \mathbb{R}^n \mid \vartheta(x, u) \le \alpha\} \subset B$ for all $u \in \mathcal{U}$, cf. [55, Definition 1.16]. Using this definition, we next investigate the fundamental properties of the mapping $\mathbf{T}_{h/\gamma i}^i$, which imply that if $x \in \text{int dom } h$, then the set $\mathbf{T}_{h/\gamma i}^i(\mathbf{x})$ is nonempty. This is an essential assertion to show the well-definedness of the algorithm (BPALM) in the next section. Let us remind that such statement is a common hypothesis for methods assuming relative smoothness in nonconvex setting; see, e.g., [24].

Proposition 3.10 (properties of Bregman proximal alternating linearized mapping) Under conditions given in Assumption 1 and setting $\gamma_i \in (0, \gamma_{i,g}^h)$ for i = 1, ..., N, the following statements are true:

- (i) $\mathbf{T}_{h/\gamma i}^{i}$ is nonempty, compact, and outer semicontinuous (osc) for all $\mathbf{x} \in \mathbf{int} \operatorname{dom} h$;
- (ii) dom $\mathbf{T}_{h/\nu i}^{i}$ = int dom h;

Proof Let us define the function Φ_i : $\mathbb{R}^{n_i} \times \operatorname{int} \operatorname{dom} h \times (0, \gamma_{i,g}^h) \to \overline{\mathbb{R}}$ given by

$$\Phi_i(z, \mathbf{x}, \gamma_i) := g_i(z) + \langle \nabla_i f(\mathbf{x}), z - x_i \rangle + \begin{cases} \frac{1}{\gamma_i} \mathbf{D}_h(\mathbf{x} + U_i(z - x_i), \mathbf{x}) & \text{if } \gamma_i \in (0, \gamma_i^0], \\ 0 & \text{if } \gamma_i = 0 \text{ and } z = x_i, \\ +\infty & \text{otherwise.} \end{cases}$$

Since *f* is C^1 , g_i is proper and lsc, and *h* is C^1 proper and lsc, the function Φ_i is proper and lsc on the sets of the form $\{(z, \boldsymbol{x}, \gamma_i)\}[||z - x_i|| \le \mu\gamma_i, 0 \le \gamma_i \le \gamma_i^0]$, for any $\mu > 0$ and a fixed $\gamma_i^0 \in (0, \gamma_{i,g}^h)$. We now show that Φ_i is level-bounded in *z* locally uniformly at every $(\boldsymbol{x}, \gamma_i) \in \operatorname{int} \operatorname{dom} h \times (0, \gamma_i^0)$. If it is not true, then there exist $(\overline{\boldsymbol{x}}, \overline{\gamma}_i) \in \operatorname{int} \operatorname{dom} h \times (0, \gamma_{i,g}^h)$, $(\boldsymbol{x}^k)_{k \in \mathbb{N}} \subset \operatorname{int} \operatorname{dom} h$, $(z^k)_{k \in \mathbb{N}}$ with $\boldsymbol{x}^k + U_i(z^k - x_i^k) \subset \operatorname{int} \operatorname{dom} h$, and $(\gamma_i^k)_{k \in \mathbb{N}} \subset (0, \gamma_i^0]$ such that $\Phi_i(z^k, \boldsymbol{x}^k, \gamma_i^k) \le \beta < \infty$ with $(\boldsymbol{x}^k, \gamma_i^k) \to (\overline{\boldsymbol{x}}, \overline{\gamma}_i)$ and $||z^k|| \to \infty$. This guarantees that, for sufficiently large *k*, $z^k \neq x_i^k$, i.e., $\gamma_i^k \in (0, \gamma_i^0]$ and

$$g_i(z^k) + \langle \nabla_i f(\boldsymbol{x}^k), z^k - x_i^k \rangle + \frac{1}{\gamma_i^k} \mathbf{D}_h(\boldsymbol{x}^k + U_i(z^k - x_i^k), \boldsymbol{x}^k) \le \beta$$

Setting $\tilde{\gamma}_i \in (\gamma_i^0, \gamma_i^h)$, Definition 3.6 and Proposition 3.7(b) ensure that there exists $\tilde{\beta} \in \mathbb{R}$ such that

$$g_i(z^k) + \frac{1}{\tilde{\gamma}_i}h(\mathbf{x}^k + U_i(z^k - x_i^k)) \ge \tilde{\beta}.$$

Subtracting the last two inequalities, it holds that

$$\langle \nabla_i f(\mathbf{x}^k), z^k - x_i^k \rangle + \frac{1}{\gamma_i^k} \mathbf{D}_h(\mathbf{x}^k + U_i(z^k - x_i^k), \mathbf{x}^k) - \frac{1}{\tilde{\gamma}_i} h(\mathbf{x}^k + U_i(z^k - x_i^k)) \le \beta - \tilde{\beta}.$$

By expanding $\mathbf{D}_h(\mathbf{x}^k + U_i(z^k - x_i^k), \mathbf{x}^k)$, dividing both sides by $||z^k||$, and taking limit from both sides of this inequality as $k \to \infty$, it can be deduced that

$$\lim_{k\to\infty} \left(\left\langle \nabla_i f(\mathbf{x}^k) - \frac{1}{\gamma_i^k} \nabla_i h(\mathbf{x}^k), \frac{z^k - x_i^k}{\|z^k\|} \right\rangle - \frac{1}{\gamma_i^k} \frac{h(\mathbf{x}^k)}{\|z^k\|} \right) + \left(\frac{1}{\gamma_i^k} - \frac{1}{\tilde{\gamma}_i} \right) \lim_{k\to\infty} \frac{h_{\mathbf{x}^k}^i(z^k)}{\|z^k\|} \le \lim_{k\to\infty} \frac{\beta - \tilde{\beta}}{\|z^k\|}.$$

This leads to the contradiction $+\infty \leq 0$, which implies that Φ_i is level-bounded in *z* locally uniformly at every $(\mathbf{x}, \gamma_i) \in \mathbf{int} \operatorname{dom} h \times (0, \gamma_i^0)$.

Considering the above result, all assumptions of the parametric minimization theorem given in [55, Theorem 1.17] are satisfied, i.e., Proposition 3.10(i) holds true. If $x \in \operatorname{int} \operatorname{dom} h$, it follows from Proposition 3.10(i) that $\mathbf{T}_{h/\gamma i}^{i}(x) \neq \emptyset$, i.e., int dom $h \subseteq \operatorname{dom} \mathbf{T}_{h/\gamma i}^{i}$. For $x \in \operatorname{dom} \mathbf{T}_{h/\gamma i}^{i}$, the essential smoothness of h yields $x \in \operatorname{int} \operatorname{dom} h$, i.e., dom $\mathbf{T}_{h/\gamma i}^{i} \subseteq \operatorname{int} \operatorname{dom} h$, giving Proposition 3.10(ii).

Remark 3.11 (sum or product separable kernel) Let Assumption 1 hold, and let $h_i : \mathbb{R}^{n_i} \to \mathbb{R}$ (i = 1, ..., N) be strictly convex, 1-coercive, and essentially smooth. Then, the following hold:

(i) If *h* is an *additive separable* function, i.e., $h(x_1, ..., x_N) = h_1(x_1) + ... + h_N(x_N)$, and $\nabla h_i(x_i) - \gamma_i \nabla_i f(\mathbf{x}) \in \mathbf{dom} \nabla h_i^*$, then (3.9) can be written in the form

$$\begin{aligned} \mathbf{T}_{h/\gamma i}^{i}(\mathbf{x}) &= \operatorname*{arg\,min}_{z \in \mathbb{R}^{n_{i}}} \{g_{i}(z) + \langle \nabla_{i}f(\mathbf{x}), z - x_{i} \rangle + \frac{1}{\gamma_{i}} \mathbf{D}_{h}(\mathbf{x} + U_{i}(z - x_{i}), \mathbf{x}) \} \\ &= \operatorname*{arg\,min}_{z \in \mathbb{R}^{n_{i}}} \{g_{i}(z) + \frac{1}{\gamma_{i}}(h_{i}(z) - h_{i}(x_{i}) - \langle \nabla h_{i}(x_{i}) - \gamma_{i} \nabla_{i}f(\mathbf{x}), z - x_{i} \rangle) \} \\ &= \operatorname*{arg\,min}_{z \in \mathbb{R}^{n_{i}}} \{g_{i}(z) + \frac{1}{\gamma_{i}} \mathbf{D}_{h_{i}}(z, \nabla h_{i}^{*}(\nabla h_{i}(x_{i}) - \gamma_{i} \nabla_{i}f(\mathbf{x})) \} \\ &= \operatorname{prox}_{\gamma_{i}g_{i}}^{h_{i}}(\nabla h_{i}^{*}(\nabla h_{i}(x_{i}) - \gamma_{i} \nabla_{i}f(\mathbf{x}))). \end{aligned}$$

(ii) If *h* is product separable, i.e., $h(x_1, ..., x_N) = h_1(x_1) \times ... \times h_N(x_N)$, each h_i (i = 1, ..., N) is positive, and $\nabla h_i(x_i) - \gamma_i \nabla_i f(\mathbf{x}) \in \mathbf{dom} \nabla h_i^*$ holds, then

$$\begin{split} \mathbf{T}_{h/\gamma i}^{i}(\mathbf{x}) &= \operatorname*{arg\,min}_{z \in \mathbb{R}^{n_{i}}} \{g_{i}(z) + \langle \nabla_{i}f(\mathbf{x}), z - x_{i} \rangle + \frac{1}{\gamma_{i}} \mathbf{D}_{h}(\mathbf{x} + U_{i}(z - x_{i}), \mathbf{x}) \} \\ &= \operatorname*{arg\,min}_{z \in \mathbb{R}^{n_{i}}} \{g_{i}(z) + \frac{\eta_{x}^{i}}{\gamma_{i}}(h_{i}(z) - h_{i}(x_{i}) - \langle \nabla h_{i}(x_{i}) - \frac{\gamma_{i}}{\eta_{x}^{i}} \nabla_{i}f(\mathbf{x}), z - x_{i} \rangle) \} \\ &= \operatorname*{arg\,min}_{z \in \mathbb{R}^{n_{i}}} \{g_{i}(z) + \frac{1}{\mu_{i}} \mathbf{D}_{h_{i}}(z, \nabla h_{i}^{*}(\nabla h_{i}(x_{i}) - \mu_{i} \nabla_{i}f(\mathbf{x})) \} \\ &= \operatorname{prox}_{\mu_{i}g_{i}}^{h_{i}}(\nabla h_{i}^{*}(\nabla h_{i}(x_{i}) - \mu_{i} \nabla_{i}f(\mathbf{x}))), \\ &\text{where } \mu_{i} := \frac{\gamma_{i}}{\eta_{i}^{i}} \text{ and } \eta_{\mathbf{x}}^{i} := \prod_{j \neq i} h_{j}(x_{j}) > 0. \end{split}$$

Note that in above-mentioned multi-block kernel function, h is not necessarily convex, but it is only multi-block convex that is a much weaker notion than the convexity.

4 Multi-block Bregman proximal alternating linearized minimization

We here introduce a multi-block proximal alternating linearized minimization algorithm and investigate its subsequential and global convergence, along with its convergence rate.

For given points $\mathbf{x}^{k} = (x_{1}^{k}, ..., x_{N}^{k})$ and $\mathbf{x}^{k+1} = (x_{1}^{k+1}, ..., x_{N}^{k+1})$, we set

$$\mathbf{x}^{k,i} := (x_1^{k+1}, \dots, x_i^{k+1}, x_{i+1}^k, \dots, x_N^k) = \mathbf{x}^k + \sum_{j=1}^i U_j (x_j^{k+1} - x_j^k),$$

🖄 Springer

so that $\mathbf{x}^{k,0} = \mathbf{x}^k$ and $\mathbf{x}^{k,N} = \mathbf{x}^{k+1}$, where $U_j \in \mathbb{R}^{n \times n_j}$, $I_n = (U_1, \dots, U_N) \in \mathbb{R}^{n \times n}$ for the identity matrix I_n . Using this notation and (3.9), we next introduce the multi-block Bregman proximal alternating linearized minimization (BPALM) algorithm.

Algorithm 1 (BPALM) Bregman Proximal Alternating Linearized Minimization

REQUIRE $\gamma_i \in (0, \frac{1}{L_i}), i = 1, ..., N, x^0 \in \operatorname{int} \operatorname{dom} h, I_n = (U_1, ..., U_N) \in \mathbb{R}^{n \times n}$ with $U_i \in \mathbb{R}^{n \times n_i}$ and the identity matrix I_n . INITIALIZE k = 0. 1: while some stopping criterion is not met do 2: $x^{k,0} = x^k$; 3: for i = 1, ..., N do $x_i^{k,i} \in \mathbf{T}^i_{h/\gamma_i}(x^{k,i-1}), \quad x^{k,i} = x^{k,i-1} + U_i(x_i^{k,i} - x_i^{k,i-1});$ (4.1) 4: $x^{k+1} = x^{k,N}, k = k + 1;$ ENSURE A vector x^k .

We note that each iteration of BPALM requires one call of the first-order oracle for the information needed in (4.1). In addition, notice that if N = 1, this algorithm reduces to the common (Bregman) proximal gradient (forward-backward) method [13, 17, 24]; if N = 2, $h, f \in C^2$, and $g_1 = g_2 = 0$, then it reduces to B-PALM [43]; if N = 2 and $h(\mathbf{x}) = \frac{1}{2}(||x_1||^2 + ||x_2||^2)$, it reduces to PALM [23]; if $h(\mathbf{x}) = \frac{1}{2}\sum_{i=1}^{N} ||x_i||^2$, then this algorithm is reduced to C-PALM [57].

From Proposition 3.10, we know that the operator $\mathbf{T}_{h/\gamma i}^{i}(\mathbf{x}^{k,i-1})$ is nonempty and compact for all $\mathbf{x}^{k,i-1} \in \mathbf{int} \operatorname{dom} h$; however, after the substitution $\mathbf{x}^{k,i} = \mathbf{x}^{k,i-1} + U_i(\mathbf{x}_i^{k,i} - \mathbf{x}_i^{k,i-1})$, we need to be sure $\mathbf{x}^{k,i} \in \mathbf{int} \operatorname{dom} h$ to guarantee the nonemptymess of $\mathbf{T}_{h/\gamma i}^{i}(\mathbf{x}^{k,i})$ in the next step of the algorithm. To do so, we require to make the following extra assumption:

Assumption 2 For i = 1, ..., N and all $z \in \mathbf{rangeT}_{h/\gamma i}^{i}(\mathbf{x})$, we have $\mathbf{x} + U_{i}(z - x_{i}) \in \mathbf{int} \operatorname{dom} h$.

Let us emphasize that together with Proposition 3.10, the latter assumption implies that BPALM is well-defined. Now, we can begin with showing some basic properties of the sequence generated by BPALM, involving a *sufficient decrease condition*.

Proposition 4.1 (sufficient decrease condition) Let Assumption 1 and 2 hold, and let $(\mathbf{x}^k)_{k \in \mathbb{N}}$ be generated by BPALM. Then, the following statements are true:

(i) the sequence $(\varphi(\mathbf{x}^k))_{k \in \mathbb{N}}$ is nonincreasing and

$$\rho \sum_{i=1}^{N} \mathbf{D}_{h}(\boldsymbol{x}^{k,i}, \boldsymbol{x}^{k,i-1}) \leq \varphi(\boldsymbol{x}^{k}) - \varphi(\boldsymbol{x}^{k+1}),$$
(4.2)

where
$$\rho := \min\left\{\frac{1-\gamma_1 L_1}{\gamma_1}, \dots, \frac{1-\gamma_N L_N}{\gamma_N}\right\};$$

(ii) we have

$$\sum_{k=1}^{\infty} \sum_{i=1}^{N} \mathbf{D}_h(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) < \infty,$$

$$(4.3)$$

i.e., $\lim_{k\to\infty} \mathbf{D}_h(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) = 0$ for i = 1, ..., N. **Proof** Plugging $\overline{z} = x_i^{k,i}$ and $\mathbf{x} = \mathbf{x}^{k,i-1}$ into Lemma 3.9, it holds that

$$\varphi(\mathbf{x}^{k,i}) \le \varphi(\mathbf{x}^{k,i-1}) - \frac{1 - \gamma_i L_i}{\gamma_i} \mathbf{D}_h(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}).$$
(4.4)

Summing up both sides of (4.4) from i = 1 to N, it follows that

$$\sum_{i=1}^{N} \frac{1-\gamma_i L_i}{\gamma_i} \mathbf{D}_h(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) \leq \sum_{i=1}^{N} [\varphi(\mathbf{x}^{k,i-1}) - \varphi(\mathbf{x}^{k,i})] = \varphi(\mathbf{x}^k) - \varphi(\mathbf{x}^{k+1}),$$

giving (4.2). Let us sum up both sides of (4.2) from k = 0 to q:

$$\rho \sum_{k=0}^{q} \sum_{i=1}^{N} \mathbf{D}_{h}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) \leq \sum_{k=0}^{q} \varphi(\mathbf{x}^{k}) - \varphi(\mathbf{x}^{k+1}) = \varphi(\mathbf{x}^{0}) - \varphi(\mathbf{x}^{q+1}) \leq \varphi(\mathbf{x}^{0}) - \inf \varphi < \infty.$$

Taking the limit as $q \to +\infty$, (4.3) holds true. Together with $\mathbf{D}_h(\cdot, \cdot) \ge 0$, this proves the claim.

Let us consider the condition

$$\sum_{i=1}^{N} \mathbf{D}_{h}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) \le \varepsilon$$
(4.5)

as a *stopping criterion*, for the accuracy parameter $\epsilon > 0$. Then, the first main consequence of Proposition 4.1 will provide us the *iteration complexity* of BPALM, which is the number of iterations needed for the stopping criterion (4.5) to be satisfied.

The boundedness of the sequence $(\mathbf{x}^k)_{k\in\mathbb{N}}$ is a typical assumption in convergence analysis of proximal-type algorithms for solving general non-convex non-smooth composite optimization problem, see e.g., [7, 24]. We next provide the iteration complexity of BPALM and a sufficient condition guaranteeing the boundedness of $(\mathbf{x}^k)_{k\in\mathbb{N}}$ as a simple consequence of Proposition 4.1.

Corollary 4.2 (iteration complexity and boundedness) Let Assumption 1 and 2 hold, and let the sequence $(\mathbf{x}^k)_{k \in \mathbb{N}}$ be generated by BPALM with the stopping criterion (4.5). Then, for the constant $\rho = \min\{\frac{1-\gamma_1L_1}{\gamma_1}, \dots, \frac{1-\gamma_NL_N}{\gamma_N}\}$,

- (i) **BPALM** will be terminated within $k \le 1 + \frac{\varphi(x^0) \inf \varphi}{\rho \varepsilon}$ iterations;
- (ii) If φ has bounded level sets, then the sequence $(\mathbf{x}^k)_{k \in \mathbb{N}}$ is bounded.

Proof Summing both sides of (4.2) over the first $\mathcal{K} \in \mathbb{N}$ iterations and telescoping the right hand side, it holds that

$$\rho \sum_{k=0}^{\mathcal{K}-1} \sum_{i=1}^{N} \mathbf{D}_{h}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) \leq \sum_{k=0}^{\mathcal{K}-1} \left(\varphi(\mathbf{x}^{k}) - \varphi(\mathbf{x}^{k+1}) \right) = \varphi(\mathbf{x}^{0}) - \varphi(\mathbf{x}^{\mathcal{K}}) \leq \varphi(\mathbf{x}^{0}) - \inf \varphi.$$

Assuming that for all the first $(\mathcal{K} - 1)$ iterations the stopping criterion (4.5) is not satisfied, i.e., $\sum_{i=1}^{N} \mathbf{D}_{h}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) > \varepsilon$, which leads to $\mathcal{K} \le 1 + \frac{\varphi(x^{0}) - \inf \varphi}{\rho \varepsilon}$, giving the desired result.

Proposition 4.1(i) shows that $\varphi(\mathbf{x}^k)$ is non-increasing; hence, the sequence $(\mathbf{x}^k)_{k \in \mathbb{N}}$ is encompassed within the lower level set $[\varphi \leq \varphi(\mathbf{x}^0)]$; i.e., $(\mathbf{x}^k)_{k \in \mathbb{N}} \subseteq [\varphi \leq \varphi(\mathbf{x}^0)]$. Since φ has bounded level sets, the sequence $(\mathbf{x}^k)_{k \in \mathbb{N}}$ would be bounded.

In order to show subsequential convergence of the sequence $(\mathbf{x}^k)_{k \in \mathbb{N}}$ generated by BPALM, the next proposition will provide a lower bound for iterations gap $\|\mathbf{x}^{k+1} - \mathbf{x}^k\|$ using the subdifferential of $\partial \varphi(\mathbf{x}^{k+1})$. In the remainder of this section, we assume that the kernel *h* has a full domain, i.e., $\overline{C} = \mathbb{R}^n$.

Assumption 3 The multi-block Legendre kernel h has full domain, i.e., dom $h = \mathbb{R}^n$.

Proposition 4.3 (subgradient lower bound for iterations gap) *Let Assumption* 1, 2 *and* 3 *hold, and let* $(\mathbf{x}^k)_{k \in \mathbb{N}}$ *be generated by* BPALM *that we assume to be bounded. For a fixed* $k \in \mathbb{N}$ *, we define*

$$\mathcal{G}_{i}^{k+1} := \frac{1}{\gamma_{i}} (\nabla_{i} h(\boldsymbol{x}^{k,i-1}) - \nabla_{i} h(\boldsymbol{x}^{k,i})) + \nabla_{i} f(\boldsymbol{x}^{k+1}) - \nabla_{i} f(\boldsymbol{x}^{k,i-1}) \quad i = 1, \dots, N.$$
(4.6)

If $\nabla_i f$ and $\nabla_i h$ (i = 1, ..., N) are locally Lipschitz on bounded sets involving all iterations, then $(\mathcal{G}_1^{k+1}, ..., \mathcal{G}_N^{k+1}) \in \partial \varphi(\mathbf{x}^{k+1})$ and

$$\|(\mathcal{G}_{1}^{k+1},\ldots,\mathcal{G}_{N}^{k+1})\| \leq \overline{c} \sum_{i=1}^{N} \|x_{i}^{k+1} - x_{i}^{k}\|,$$
(4.7)

with $\overline{c} := \max\{\frac{\widetilde{L}+\gamma_1\widehat{L}}{\gamma_1}, \dots, \frac{\widetilde{L}+\gamma_N\widehat{L}}{\gamma_N}\}$ in which \widehat{L} and $\widetilde{L} > 0$ are Lipschitz moduli of $\nabla_i f$, $\nabla_i h \ (i = 1, \dots, N)$ on bounded sets involving all iterations.

Proof The optimality conditions for (4.1) ensure that there exist $q_i^{k,i} \in \partial g_i(x_i^{k,i})$ such that

$$\nabla_i f(\boldsymbol{x}^{k,i-1}) + \frac{1}{\gamma_i} \left(\nabla_i h(\boldsymbol{x}^{k,i}) - \nabla_i h(\boldsymbol{x}^{k,i-1}) \right) + q^{k,i} = 0 \quad i = 1, \dots, N,$$

leading to

$$q^{k,i} = \frac{1}{\gamma_i} \left(\nabla_i h(\boldsymbol{x}^{k,i-1}) - \nabla_i h(\boldsymbol{x}^{k,i}) \right) - \nabla_i f(\boldsymbol{x}^{k,i-1}) \quad i = 1, \dots, N.$$

$$(4.8)$$

On the other hand, owing to [7, Proposition 2.1], the subdifferential of φ is given by

🖉 Springer

$$\partial \varphi(\mathbf{x}) = (\partial_1 \varphi(\mathbf{x}), \dots, \partial_N \varphi(\mathbf{x})) = (\nabla_1 f(\mathbf{x}) + \partial g_1(x_1), \dots, \nabla_N f(\mathbf{x}) + \partial g_N(x_N)),$$

i.e., for $x = x^{k+1}$,

$$\nabla_i f(\boldsymbol{x}^{k+1}) + \partial g_i(\boldsymbol{x}_i^{k+1}) = \partial_i \varphi(\boldsymbol{x}^{k+1}) \quad i = 1, \dots, N,$$

which means $(\mathcal{G}_1^{k+1}, \dots, \mathcal{G}_N^{k+1}) \in \partial \varphi(\mathbf{x}^{k+1})$. It follows from Assumption 3 and the Lipschitz continuity of $\nabla_i f$, $\nabla_i h$ on bounded sets involving all iterations and the assumption of $(\mathbf{x}^k)_{k\in\mathbb{N}}$ being bounded that there exist $\hat{L}, \tilde{L} > 0$ such that

$$\|\mathcal{G}_{i}^{k+1}\| \leq \frac{1}{\gamma_{i}} \|\nabla_{i}h(\mathbf{x}^{k,i-1}) - \nabla_{i}h(\mathbf{x}^{k,i})\| + \|\nabla_{i}f(\mathbf{x}^{k+1}) - \nabla_{i}f(\mathbf{x}^{k,i-1})\| \leq \frac{\widetilde{L} + \gamma_{i}\widehat{L}}{\gamma_{i}} \sum_{j=1}^{N} \|x_{j}^{k+1} - x_{j}^{k}\|,$$

for i = 1, ..., N. Invoking the last two inequalities, it can be concluded that

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

as claimed.

Next, we proceed to derive the *subsequential convergence* of the sequence $(\mathbf{x}^k)_{k\in\mathbb{N}}$ generated by BPALM: every cluster point of $(\mathbf{x}^k)_{k\in\mathbb{N}}$ is a critical point of φ . Further, we explain some fundamental properties of the set of all cluster points $\omega(\mathbf{x}^0)$ of the sequence $(\mathbf{x}^k)_{k\in\mathbb{N}}$.

Theorem 4.4 (subsequential convergence and properties of $\omega(\mathbf{x}^0)$) Let Assumption 1, 2 and 3 hold, let the kernel h be locally multi-block strongly convex, and let $(\mathbf{x}^k)_{k\in\mathbb{N}}$ be generated by BPALM. If the sequence $(\mathbf{x}^k)_{k\in\mathbb{N}}$ is bounded and $\nabla_i f$ and $\nabla_i h$ (i = 1, ..., N) are locally Lipschitz around its cluster points, then all the cluster points of $(\mathbf{x}^k)_{k\in\mathbb{N}}$ are critical points φ , i.e., $\omega(\mathbf{x}^0) \subset \operatorname{crit} \varphi$.

Proof For a limit point $\mathbf{x}^* = (x_1^*, \dots, x_N^*)$ of the sequence $(\mathbf{x}^k)_{k \in \mathbb{N}}$, it follows from the boundedness of this sequence that there exists an infinite index set $\mathcal{J} \subset \mathbb{N}$ such that the subsequence $(\mathbf{x}^k)_{k \in \mathbb{N}}[k \in \mathcal{J}]$ converges to \mathbf{x}^* as $k \to \infty$. From the lower semicontinuity of g_i $(i = 1, \dots, N)$ and for $j \in \mathcal{J}$, it can be deduced that

$$\liminf_{j \to \infty} g_i(x_i^{k_j}) \ge g(x_i^{\star}) \quad i = 1, \dots, N.$$

$$(4.9)$$

By (4.1), we get

$$\langle \nabla_{i}f(\boldsymbol{x}^{k,i-1}), \boldsymbol{x}_{i}^{k+1} - \boldsymbol{x}_{i}^{k} \rangle + \frac{1}{\gamma_{i}} \mathbf{D}_{h}(\boldsymbol{x}^{k,i}, \boldsymbol{x}^{k,i-1}) + g_{i}(\boldsymbol{x}_{i}^{k+1}) \qquad \leq \langle \nabla_{i}f(\boldsymbol{x}^{k,i-1}), \boldsymbol{x}_{i}^{\star} - \boldsymbol{x}_{i}^{k} \rangle$$
$$+ \frac{1}{\gamma_{i}} \mathbf{D}_{h}(\boldsymbol{x}^{\star}, \boldsymbol{x}^{k,i-1}) + g_{i}(\boldsymbol{x}_{i}^{\star}).$$
(4.10)

Using multi-block local strong convexity of *h* around \mathbf{x}^* and invoking Proposition 4.1(ii), there exist a neighborhood $\mathbf{B}(x_i^*, \varepsilon_i^*)$ for $\varepsilon_i^* > 0$, $\sigma_i^* > 0$, and $k_i^0 \in \mathbb{N}$ such that for $k \ge k_i^0$ and $k \in \mathcal{J}$

$$\lim_{k \to \infty} \frac{\sigma_i^{\star}}{2} \|x_i^{k+1} - x_i^k\|^2 \le \lim_{k \to \infty} \mathbf{D}_h(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) = 0, \quad x_i^k \in \mathbf{B}(x_i^{\star}, \varepsilon_i^{\star}), \ i = 1, \dots, N.$$
(4.11)

This indicates that the distance between two successive iterations goes to zero for large enough k, i.e., for $\lim_{k\to\infty} x_i^{k+1} = \lim_{k\to\infty} x_i^k$. Together with the boundedness of the sequence $(\mathbf{x}^k)_{k\in\mathbb{N}}$, the continuity of $\nabla_i f$ and $\nabla_i h$ around cluster points of $(\mathbf{x}^k)_{k\in\mathbb{N}}$, the substitution of $k = k_j - 1$ in (4.10) for $j \in \mathcal{J}$, taking the limit in both sides of this inequality as $k \to \infty$, this implies that

$$\limsup_{j\to\infty} g_i(x_i^{k_j}) \le g_i(x_i^{\star}) \quad i=1,\ldots,N,$$

and consequently by lsc of g_i ,

$$\lim_{j \to \infty} \varphi(\mathbf{x}^{k_j}) = \lim_{j \to \infty} \left(f(x_1^{k_j}, \dots, x_N^{k_j}) + \sum_{i=1}^N g_i(x_i^{k_j}) \right) = f(x_1^{\star}, \dots, x_N^{\star}) + \sum_{i=1}^N g_i(x_i^{\star}).$$
(4.12)

Further, Proposition 4.1(ii) and Proposition 4.3 ensure $(\mathcal{G}_1^{k+1}, \dots, \mathcal{G}_N^{k+1}) \in \partial \varphi(\mathbf{x}^{k+1})$ and

$$\lim_{k \to +\infty} \|\mathcal{G}_1^{k+1}, \dots, \mathcal{G}_N^{k+1}\| \le \overline{c} \lim_{k \to +\infty} \sum_{i=1}^N \|x_i^{k+1} - x_i^k\| \le \overline{c} \lim_{k \to +\infty} \left(\sum_{i=1}^N \sqrt{\frac{2}{\sigma_i^*} \mathbf{D}_h(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1})} \right) = 0,$$

i.e., $\lim_{k\to\infty} (\mathcal{G}_1^{k+1}, \dots, \mathcal{G}_N^{k+1}) = (0_{n_1}, \dots, 0_{n_N})$. Together with the closedness of the subdifferential mapping $\partial \varphi$ and (4.12), this implies that $(0_{n_1}, \dots, 0_{n_N}) \in \partial \varphi(x_1^*, \dots, x_N^*)$, giving our desired result.

4.1 Global convergence under Kurdyka-Łojasiewicz inequality

Our following main result indicates that the sequence $(x^k)_{k\in\mathbb{N}}$ generated by BPALM converges to a critical point x^* of φ if it satisfies the KL property; cf. Definition 2.1.

Theorem 4.5 (global convergence) Let Assumption 1, 2 and 3 hold, let the kernels *h* be multi-block globally strongly convex with modulus σ_i (i = 1, ..., N), and let $(\mathbf{x}^k)_{k \in \mathbb{N}}$ be generated by BPALM that we assume to be bounded. If φ is a KL function, then the following statements are true:

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

$$\sum_{k=1}^{\infty} \|x_i^{k+1} - x_i^k\| < \infty \quad i = 1, \dots, N;$$
(4.13)

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

Proof Let us set $\varphi^* := \lim_{k \to \infty} \varphi(\mathbf{x}^k)$ and define the sequence $(\mathcal{S}_k)_{k \in \mathbb{N}}$ given by $\mathcal{S}_k := \varphi(\mathbf{x}^k) - \varphi^*$, which is decreasing by Proposition 4.1(i), i.e., $(\mathcal{S}_k)_{k \in \mathbb{N}} \to 0$. We now consider two cases: (i) there exists $k \in \mathbb{N}$ such that $\mathcal{S}_k = 0$; (ii) $\mathcal{S}_k > 0$ for all $k \ge 1$.

In Case (i), Proposition 4.1(i) implies that $\varphi(\mathbf{x}^k) = \varphi^*$ for all $k \ge \overline{k}$. It follows from Proposition 4.1(ii) and multi-block strong convexity of *h* that

$$\frac{\sigma_i}{2} \| x_i^{k+1} - x_i^k \|^2 \le \mathbf{D}_h(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) = 0 \quad i = 1, \dots, N,$$

implying $\mathbf{x}^{k+1} = \mathbf{x}^k$ for all $k \ge \overline{k}$, which leads to Theorem 4.5(i).

In Case (ii), it holds that $\varphi(\mathbf{x}^k) > \varphi^*$ for all $k \ge 1$. Thus, Lemma 8.1 yields that there exists $\varepsilon, \eta > 0$ and the desingularizing function ψ such that

$$\psi'(\varphi(\mathbf{x}^k) - \varphi^{\star})$$
dist $(0, \partial \varphi(\mathbf{x}^k)) \ge 1$ for $k \ge k_0$.

Let us define $\Delta_k := \psi(\varphi(\mathbf{x}^k) - \varphi^*) = \psi(\mathcal{S}_k)$. Then, it follows from the concavity of ψ and Proposition 4.3 that

$$\begin{split} \Delta_{k} - \Delta_{k+1} &= \psi(\mathcal{S}_{k}) - \psi(\mathcal{S}_{k+1}) \geq \psi'(\mathcal{S}_{k})(\mathcal{S}_{k} - \mathcal{S}_{k+1}) = \psi'(\mathcal{S}_{k})(\varphi(\mathbf{x}^{k}) - \varphi(\mathbf{x}^{k+1})) \\ &\geq \frac{\varphi(\mathbf{x}^{k}) - \varphi(\mathbf{x}^{k+1})}{\operatorname{dist}(0, \partial \varphi(\mathbf{x}^{k}))} \geq \frac{\rho \sum_{i=1}^{N} \mathbf{D}_{h}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1})}{\overline{c} \sum_{i=1}^{N} \|x_{i}^{k} - x_{i}^{k-1}\|} \geq \frac{1}{\widehat{c}} \frac{\sum_{i=1}^{N} \|x_{i}^{k+1} - x_{i}^{k}\|^{2}}{\sum_{i=1}^{N} \|x_{i}^{k} - x_{i}^{k-1}\|}, \end{split}$$

with $\hat{c} := \frac{\bar{c}}{(\rho \min\{\sigma_1, \dots, \sigma_N\})}$. Using the arithmetic and quadratic mean inequalities, and applying the arithmetic and geometric mean inequalities, it can be concluded that

$$\sum_{i=1}^{N} \|x_i^{k+1} - x_i^k\| \le \sqrt{\widehat{c}N(\Delta_k - \Delta_{k+1})} \sum_{i=1}^{N} \|x_i^k - x_i^{k-1}\| \le \frac{1}{2} \sum_{i=1}^{N} \|x_i^k - x_i^{k-1}\| + \frac{\widehat{c}N}{2} (\Delta_k - \Delta_{k+1}).$$
(4.14)

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

$$a_{k+1} := \sum_{i=1}^{N} \|x_i^{k+1} - x_i^k\|, \quad b_k = \frac{\hat{c}N}{2} (\Delta_k - \Delta_{k+1}), \quad \alpha := \frac{1}{2},$$
(4.15)

where $\sum_{i=1}^{\infty} b_k = \frac{\partial N}{2} \sum_{i=1}^{\infty} (\Delta_i - \Delta_{i+1}) = \frac{\partial N}{2} (\Delta_1 - \Delta_{\infty}) = \frac{\partial N}{2} \Delta_1 < \infty$. According to Fact 2.3, we infer $\sum_{k=1}^{\infty} a_k < \infty$, which proves Theorem 4.5(i).

By (4.13), the sequence $(\mathbf{x}^k)_{k \in \mathbb{N}}$ is a Cauchy sequence, i.e., it converges to a stationary point \mathbf{x}^* , giving the desired result.

4.2 Convergence rate under Łojasiewicz-type inequality

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

$$|\varphi(\mathbf{x}) - \varphi^{\star}|^{\theta} \le \kappa \operatorname{dist}(0, \partial \varphi(\mathbf{x})) \quad \forall \mathbf{x} \in \mathbf{B}(\mathbf{x}^{\star}; \varepsilon).$$
(4.16)

The following fact plays a key role in studying the convergence rate of the sequence generated by BPALM.

Fact 4.6 (convergence rate of a sequence with positive elements) [26, 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 $s_k^{2\theta} \le \beta(s_k - s_{k+1})$ 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 for every $k \in \mathbb{N}$

$$0 \le s_k \le \lambda \tau^k;$$

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

$$0 \le s_k \le \mu k^{-\frac{1}{2\theta - 1}}.$$

We next derive the *convergence rates* of the sequences $(\mathbf{x}^k)_{k \in \mathbb{N}}$ and $(\varphi(\mathbf{x}^k))_{k \in \mathbb{N}}$ under the additional assumption that the function φ satisfies the KL inequality of Łojasiewicz type.

Theorem 4.7 (convergence rate) Let Assumption 1, 2 and 3 hold, let the kernel h be multi-block globally strongly convex with modulus $(\sigma_1, ..., \sigma_N)$, and let the sequence $(\mathbf{x}^k)_{k\in\mathbb{N}}$ generated by BPALM converging to \mathbf{x}^* . If φ satisfies the KL inequality of Lojasiewicz type (4.16), then the following assertions hold:

- (i) If $\theta = 0$, then the sequences $(\mathbf{x}^k)_{k \in \mathbb{N}}$ and $(\varphi(\mathbf{x}^k))_{k \in \mathbb{N}}$ converge in a finite number of steps to \mathbf{x}^* and $\varphi(\mathbf{x}^*)$, respectively;
- (ii) If $\theta \in (0, 1/2]$, then there exist $\lambda_1 > 0$, $\mu_1 > 0$, $\tau \in [0, 1)$, and $\overline{k} \in \mathbb{N}$ such that

$$0 \le \|\boldsymbol{x}^{k} - \boldsymbol{x}^{\star}\| \le \lambda_{1} \tau^{\frac{k}{2}}, \quad 0 \le S_{k} \le \mu_{1} \tau^{k} \quad \forall k \ge \overline{k};$$

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

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

Proof See the proof in the appendix.

4.3 Adaptive BPALM

The tightness of the *i*-th block upper bound of the function f given in Proposition 3.5(b) is dependent on the parameter $L_i > 0$; however, in general, this

parameter is a global information and it might not be tight locally, i.e., one may find a $L_i > \overline{L}_i(\mathbf{x}) \ge 0$ such that

$$f(\mathbf{x} + U_i(y_i - x_i)) \le f(\mathbf{x}) + \langle \nabla_i f(\mathbf{x}), y_i - x_i \rangle + L_i(\mathbf{x}) \mathbf{D}_h(\mathbf{x} + U_i(y_i - x_i), \mathbf{x})$$

for all $\mathbf{y} \in \mathbf{B}(\mathbf{x}; \varepsilon_1)$ with a small enough $\varepsilon_1 > 0$. Consequently, the block majorization model described by \mathcal{M}_{h/γ_i} may not be tight enough, which will consequently lead to smaller step-sizes $\gamma_i \in (0, 1/L_i)$. In this case and in the case where (L_1, \dots, L_N) is not available, one can retrieve them adaptively by applying a line search starting from a lower estimate; see, e.g., [3, 5, 47, 49].

Putting together the above discussions, we propose an adaptive version of BPALM using a line search; see Algorithm 2.

Algorithm 2 (A-BPALM) adaptive BPALM

REQUIRE $x^0 \in \operatorname{int} \operatorname{dom} h, v > 1, \overline{L}_i^0 > 0$ for $i = 1, \dots, N, I_n = (U_1, \dots, U_N) \in \mathbb{R}^{n \times n}$ with $U_i \in \mathbb{R}^{n \times n_i}$ and the identity matrix I_n . INITIALIZE $k = 0, p = 0, \gamma_i^0 \in (0, 1/\overline{L}_i^0)$ for i = 1, ..., N. 1: while some stopping criterion is not met do $\boldsymbol{x}^{k,0} = \boldsymbol{x}^k;$ 2: for i = 1, ..., N do 3: 4: repeat $\begin{array}{ll} \operatorname{set} \overline{L}_i^{k+1} = \nu^p \overline{L}_i^k, \quad \gamma_i^{k+1} = \gamma_i^k / \nu^p, \quad p = p+1; \\ \operatorname{compute} x_i^{k,i} \in \mathbf{T}_{y_i j_i^{k+1}}(\boldsymbol{x}^{k,i-1}), \quad \boldsymbol{x}^{k,i} = \boldsymbol{x}^{k,i-1} + U_i(x_i^{k,i} - x_i^{k,i-1}); \end{array}$ 5: 6: until $f(x^{k,i}) \leq f(x^{k,i-1}) + \langle \nabla_i f(x^{k,i-1}), x_i^{k,i} - x_i^{k,i-1} \rangle + \overline{L}_i^{k+1} \mathbf{D}_h(x^{k,i}, x^{k,i-1})$ 7: $p_i^k = p - 1; \ p = 0;$ 8: $x^{k+1} = x^{k,N}, k = k+1;$ 9: ENSURE A vector x^k .

Let us assume that the solution of the subproblem $x_i^{k,i} \in \mathbf{T}_{h/\gamma_i^{k+1}}^i(\mathbf{x}^{k,i-1})$ can be computed exactly. Beside of the computational cost of solving the subproblem, we note that in each iteration of A-BPALM, we only needs two calls of the firstorder oracle (one call for *f* and the other for *h* in Line 7). We next provide an upper bound on the *total number of calls of oracle* after *k* iterations of A-BPALM and those needed to satisfy (4.5).

Proposition 4.8 (worst-case oracle calls) Let $(\mathbf{x}^k)_{k\in\mathbb{N}}$ be generated by A-BPALM. *Then*,

- (i) after at most $\max\left\{\frac{1}{\ln v}\ln\frac{vL_i}{\overline{L}_i^0}, 0\right\}$ iterations the line search (Lines 4 to 7 of A-BPALM) will be satisfied;
- (ii) the number of oracle calls N_k after k full cycle is bounded by

$$\mathcal{N}_k \leq \left(2N + 2\sum_{i=1}^N \max\left\{\frac{1}{\ln\nu}\ln\frac{\nu L_i}{\overline{L}_i^0}, 0\right\}\right)k;$$

(iii) the worst-case number of oracle calls to satisfy (4.5) is given by

$$\left(2N+2\sum_{i=1}^{N}\max\left\{\frac{1}{\ln\nu}\ln\frac{\nu L_{i}}{\overline{L}_{i}^{0}},0\right\}\right)\left(1+\frac{(\varphi(\mathbf{x}^{0})-\inf\varphi)}{\overline{\rho}\varepsilon}\right)$$

with $\overline{\rho}$:= min $\left\{\frac{(1-\gamma_{1}^{0}\overline{L}_{1}^{0})}{\gamma_{1}^{0}},\ldots,\frac{(1-\gamma_{N}^{0}\overline{L}_{N}^{0})}{\gamma_{N}^{0}}\right\}$.

Proof According to step 5 and step 8 of A-BPALM, if $\ln \overline{L}_i^0 \ge \ln(\nu L_i)$ then from step 5 we have $\overline{L}_i^{k+1} \ge \overline{L}_i^k \ge \overline{L}_i^0 > L_i$ for all $k \ge 0$, hence $p_i^k = 0$ since the condition in step 7 is satisfied; otherwise, we have $\overline{L}_i^{k+1} = \nu^{p_i^k} \overline{L}_i^k$, i.e.,

$$p_i^k = \frac{1}{\ln\nu} \left(\mathbf{ln} \overline{L}_i^{k+1} - \mathbf{ln} \overline{L}_i^k \right) \le \frac{1}{\ln\nu} \left(\mathbf{ln}(\nu L_i) - \mathbf{ln} \overline{L}_i^0 \right) \quad i = 1, \dots, N,$$

giving Proposition 4.8(i). Hence, the total number of calls of oracle after k iterations is given by

$$\mathcal{N}_{k} = \sum_{j=0}^{k-1} \sum_{i=1}^{N} 2(p_{i}^{j} + 1) \leq 2 \sum_{i=1}^{N} \left[k + \sum_{j=0}^{k-1} \max\left\{ \frac{1}{\ln \nu} \ln \frac{\nu L_{i}}{L_{i}^{0}}, 0 \right\} \right],$$

giving Proposition 4.8(ii).

The step 5 implies that the sequences $\left((1 - \gamma_i^k \overline{L}_i^k) / \gamma_i^k\right)_{k \in \mathbb{N}}$ (i = 1, ..., N) are increasing with respect to k, i.e., $(1 - \gamma_i^{k+1} \overline{L}_i^{k+1}) / \gamma_i^{k+1} \ge (1 - \gamma_i^0 \overline{L}_i^0) / \gamma_i^0$, i = 1, ..., N. Now, following the proof of Proposition 4.1, it is easy to see that

$$\overline{\rho} \sum_{i=1}^{N} \mathbf{D}_{h}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) = \min\{(1 - \gamma_{1}^{0} \overline{L}_{1}^{0}) / \gamma_{1}^{0}, \dots, (1 - \gamma_{N}^{0} \overline{L}_{N}^{0}) / \gamma_{N}^{0}\} \sum_{i=1}^{N} \mathbf{D}_{h}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) \\ \leq \min\{\frac{1 - \gamma_{1}^{k+1} \overline{L}_{1}^{k+1}}{\gamma_{1}^{k+1}}, \dots, \frac{1 - \gamma_{N}^{k+1} \overline{L}_{N}^{k+1}}{\gamma_{N}^{k+1}}\} \sum_{i=1}^{N} \mathbf{D}_{h}(\mathbf{x}^{k,i}, \mathbf{x}^{k,i-1}) \\ \leq \varphi(\mathbf{x}^{k}) - \varphi(\mathbf{x}^{k+1}).$$
(4.17)

From the proof of Corollary 4.2, **BPALM** will be terminated within $k \le 1 + \frac{\varphi(x^0) - \inf \varphi}{\overline{\rho} \varepsilon}$ iterations. Together with Proposition 4.8(ii), this implies that Proposition 4.8(iii) is true.

In light of (4.17), Proposition 4.1 holds true by replacing ρ with $\overline{\rho}$. Considering this replacement, all the results of Proposition 4.3, Theorem 4.4, Theorem 4.5, and Theorem 4.7 remain valid for A-BPALM.

Remark 4.9 (A-BPALM variant) We here notice that one may change Line 5 of A-BPALM as "set $\overline{L}_i^{k+1} = v^p \overline{L}_i^0$, $\gamma_i^{k+1} = \frac{\gamma_i^0}{v^p}$, p = p + 1;", which always start the procedure from \overline{L}_i^0 and γ_i^0 . It is easy to see that the results of Proposition 4.8 are still valid for this variant of A-BPALM.

5 Application to orthogonal nonnegative matrix 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., [29, 33, 34] 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) sparsity constraints on these factor matrices. More precisely, let the data matrix be $X = (x_1, x_2, \dots, x_q) \in \mathbb{R}^{p \times q}_+$ where each x_i represents some data point. NMF seeks a decomposition of X into a nonnegative $p \times r$ basis matrix $U = (u_1, u_2, \dots, u_r) \in \mathbb{R}^{p \times r}_+$ and a nonnegative $r \times q$ coefficient matrix $V = (v_1, v_2, \dots, v_r)^T \in \mathbb{R}^{r \times q}_+$ such that

$$X \approx UV,$$
 (5.1)

where $\mathbb{R}^{p\times q}_+$ is the set of $p \times q$ element-wise nonnegative matrices. Extensive research has been carried out on variants of NMF, and most studies in this area have focused on algorithmic developments, but with very limited convergence theory. This motivates us to study the application of BPALM and A-BPALM to a variant of NMF, namely orthogonal NMF (ONMF).

5.1 Orthogonal nonnegative matrix factorization

Besides the decomposition (5.1), the *orthogonal nonnegative matrix factorization* (ONMF) involves an additional orthogonality constraint $VV^T = I_r$ leading to the constrained optimization problem

minimize
$$\frac{1}{2} ||X - UV||_F^2$$

subject to $U \ge 0, V \ge 0, VV^T = I_r$ (5.2)

where $I_r \in \mathbb{R}^{r \times r}$ is the identity matrix. By imposing the matrix *V* to be orthogonal (as well as nonnegative), ONMF imposes that each data point is only associated with one basis vector, hence ONMF is closely related to clustering problems; see [52] and the references therein. Since the projection onto the set $\mathcal{C} := \{(U, V) \in \mathbb{R}^{p \times r} \times \mathbb{R}^{r \times q} \mid U \ge 0, V \ge 0, VV^T = I_r\}$ is costly, we here consider the penalized formulation

minimize
$$\frac{1}{2} \|X - UV\|_F^2 + \frac{\lambda}{2} \|I_r - VV^T\|_F^2$$

subject to $U \ge 0, V \ge 0,$ (5.3)

for the penalty parameter $\lambda > 0$. By fixing *U*, the objective function of the problem (5.3) is not Lipschitz smooth with respect to *V*, and hence standard NMF block proximal gradient descent algorithms cannot be applied to solve (5.3). The unconstrained version of (5.3) is given by

$$\underset{(U,V)}{\text{minimize}} \quad \frac{1}{2} \| X - UV \|_F^2 + \frac{\lambda}{2} \| I_r - VV^T \|_F^2 + \delta_{U \ge 0} + \delta_{V \ge 0}, \tag{5.4}$$

where $\delta_{U \ge 0}$ and $\delta_{V \ge 0}$ are the indicator functions of the sets $C_1 := \{U \in \mathbb{R}^{p \times r} \mid U \ge 0\}$ and $C_2 := \{V \in \mathbb{R}^{r \times q} \mid V \ge 0\}$, respectively. This problem can be put in the form of (1.1) using

$$f(U,V) := \frac{1}{2} \|X - UV\|_F^2 + \frac{\lambda}{2} \|I_r - VV^T\|_F^2, \quad g_1(U) := \delta_{U \ge 0}, \quad g_2(V) := \delta_{V \ge 0},$$
(5.5)

in which both $g_1(U)$ and $g_2(V)$ are nonsmooth convex functions.

We now apply BPALM and A-BPALM to solve (5.4). Specifically, in Proposition 5.1, we provide a full-domain multi-block strongly convex kernel function hsuch that f(U, V) is (L_1, L_2) -smooth relative to h, and then in Theorem 5.2 we give closed-form solutions of the subproblem (4.1). Let us emphasize that Assumption 1, 2 and 3 hold for the problem (5.4) and the kernel (5.6). Now, applying Theorem 4.5, we have that BPALM and A-BPALM converges globally to a stationary point of the objective function of (5.4). For the ONMF problem (5.4), since the function $f(U, V) = \frac{1}{2} ||X - UV||_F^2 + \frac{\lambda}{2} ||I_r - VV^T||_F^2$ is a polynomial and the functions $g_1(U) = \delta_{U>0}$ and $g_2(V) = \delta_{V>0}^2$ are indicator functions of semialgebraic sets, they are semialgebraic functions. Therefore, it follows from [20, Theorem 3.1] that $\varphi(U, V) = f(U, V) + g_1(U) + g_2(V)$ is a KL function with exponent $\theta \in [0, 1)$. We however could not provide a KL exponent (see (4.16)) for this problem, and so we are not able to derive the rate of convergence of our methods for this problem, which remains an open question.

Proposition 5.1 (multi-block relative smoothness of ONMF objective) Let $\alpha, \beta, \varepsilon_1, \varepsilon_2 > 0$ and let the function $h : \mathbb{R}^{p \times r} \times \mathbb{R}^{r \times q} \to \overline{\mathbb{R}}$ be a kernel given by

$$h(U,V) := \frac{\alpha}{2} \|U\|_F^2 \|V\|_F^2 + \frac{\beta}{4} \|V\|_F^4 + \frac{\varepsilon_1}{2} \|U\|_F^2 + \frac{\varepsilon_2}{2} \|V\|_F^2.$$
(5.6)

Then the function f given in (5.5) is (L_1, L_2) -smooth relative to h with

$$L_1 \ge \frac{1}{\alpha}, \quad L_2 \ge \max\{\frac{6\lambda}{\beta}, \frac{1}{\alpha}\}.$$
 (5.7)

Proof Using partial derivatives $\nabla_U f(U, V) = UVV^T - XV^T$, $\nabla^2_{UU} f(U, V)Z = ZVV^T$, and the Cauchy-Schwarz inequality, it can be concluded that $\langle Z, \nabla^2_{UU} f(U, V) Z \rangle \leq \|V\|_F^2 \|Z\|_F^2$. On the other hand, $\nabla_U h(U, V) = \varepsilon_1 U + \alpha \|V\|_F^2 U$ and

$$\langle Z, \nabla_{UU}^2 h(U, V) Z \rangle = (\varepsilon_1 + \alpha \|V\|_F^2) \langle Z, Z \rangle \ge \alpha \|V\|_F^2 \|Z\|_F^2.$$

Together with (5.7), this yields $L_1 \nabla^2_{UU} h(U, V) - \nabla^2_{UU} f(U, V) \ge 0$. From $\nabla_V f(U, \cdot)(V) = U^T UV - U^T X + 2\lambda (VV^T V - V)$ and the definition of directional derivative, we obtain

$$\nabla_{VV}^{2} f(U, V)Z = \lim_{t \to 0} \frac{U^{T} U(V + tZ) - U^{T} X + 2\lambda [(V + tZ)(V + tZ)^{T}(V + tZ) - (V + tZ)]}{t}$$
$$= \frac{U^{T} UV - U^{T} X + 2\lambda (VV^{T} V - V)}{t}$$
$$= U^{T} UZ + 2\lambda (ZV^{T} V + VZ^{T} V + VV^{T} Z - Z) \quad \forall Z \in \mathbb{R}^{r \times q}.$$

From this last equation, $\langle Y_1, Y_2 \rangle := \mathbf{tr}(Y_1^T Y_2)$, basic properties of the trace, the Cauchy-Schwarz inequality, and the submultiplicative property of the Frobenius norm, we obtain

$$\begin{aligned} \langle Z, \nabla^2_{VV} f(U, V) Z \rangle &= \langle Z, U^T U Z + 2\lambda (ZV^T V + VZ^T V + VV^T Z - Z) \rangle \\ &\leq 6\lambda \|Z\|_F^2 \|V\|_F^2 + \|U\|_F^2 \|Z\|_F^2. \end{aligned}$$

We have $\nabla_V h(U, V) = (\varepsilon_2 + \alpha ||U||_F^2)V + \beta ||V||^2 V$. Hence,

$$\nabla_{VV}^2 h(U, V) Z = \beta \left(\|V\|_F^2 Z + 2\langle V, Z \rangle V \right) + (\varepsilon_2 + \alpha \|U\|_F^2) Z,$$

which implies

$$\begin{aligned} \langle Z, \nabla^2_{VV} h(U, V) Z \rangle &= \beta \|Z\|_F^2 \|V\|_F^2 + 2\beta \langle V, Z \rangle^2 + (\varepsilon_2 + \alpha \|U\|_F^2) \|Z\|_F^2 \\ &\geq \beta \|Z\|_F^2 \|V\|_F^2 + \alpha \|U\|_F^2 \|Z\|_F^2. \end{aligned}$$

Hence, it follows from (5.7) that $L_2 \nabla^2_{VV} h(U, V) - \nabla^2_{VV} f(U, V) \ge 0$.

For given U^k and V^k , applying BPALM and A-BPALM to (5.4), U^{k+1} and V^{k+1} should be computed efficiently, which we study next.

Theorem 5.2 (closed-form solutions of the subproblem (4.1) **for ONMF)** *Let h be the kernel functions given in* (5.6). *For given* U^k *and* V^k , *the problem* (5.4), *and the subproblem* (4.1), *the following assertions hold*:

(i) Let
$$\eta_1 = \alpha \|V^k\|_F^2 + \varepsilon_1$$
. The iteration U^{k+1} is given by

$$U^{k+1} = \max\{U^k - \frac{\gamma_1}{n} (U^k V^k (V^k)^T - X (V^k)^T), 0\};$$
(5.8)

(ii) Let $\eta_2 = \varepsilon_2 + \alpha \|U^{k+1}\|_F^2$. The iteration V^{k+1} is given by

$$V^{k+1} = \frac{1}{t_k} \max\{G^k, 0\}$$
(5.9)

with $\nabla_V f(U^{k+1}, V^k) = (U^{k+1})^T U^{k+1} V^k - (U^{k+1})^T X + 2\lambda (V^k (V^k)^T V^k - V^k)$ and

$$G^{k} = \eta_{2}V_{k} + \beta \|V_{k}\|_{F}^{2}V_{k} - \gamma_{2}\nabla_{V}f(U^{k+1}, V^{k}),$$

$$t_{k} = \frac{\eta_{2}}{3} + \sqrt[3]{\frac{c+\sqrt{\Delta}}{2} + \frac{\eta_{2}^{3}}{27}} + \sqrt[3]{\frac{c-\sqrt{\Delta}}{2} + \frac{\eta_{2}^{3}}{27}},$$
(5.10)

🖄 Springer

where $c = \beta \| \max(0, G^k) \|_F^2$ and $\Delta = c^2 + \frac{4}{27} c \eta_2^3$. **Proof** It follows from (3.11) that

$$\begin{split} U^{k+1} &= \underset{U \in \mathbb{R}^{p \times r}}{\arg\min} \{ \langle \nabla_{U} f(U^{k}, V^{k}), U - U^{k} \rangle + \frac{1}{\gamma_{1}} \mathbf{D}_{h}((U, V^{k}), (U^{k}, V^{k})) + g_{1}(U) \} \\ &= \underset{U \ge 0}{\arg\min} \{ \langle \gamma_{1} \nabla_{U} f(U^{k}, V^{k}) - \nabla_{U} h(U^{k}, V^{k}), U \rangle + \frac{\eta_{1}}{2} \|U\|_{F}^{2} \} \\ &= \underset{U \ge 0}{\arg\min} \|U - (\frac{1}{\eta_{1}} \nabla_{U} h(U^{k}, V^{k}) - \frac{\gamma_{1}}{\eta_{1}} \nabla_{U} f(U^{k}, V^{k})) \|_{F}^{2} \\ &= \operatorname{Proj}_{U \ge 0} (U^{k} - \frac{\gamma_{1}}{\eta_{1}} \nabla_{U} f(U^{k}, V^{k})), \end{split}$$

with $\nabla_U f(U^k, V^k) = U^k V^k (V^k)^T - X(V^k)^T$, giving (5.8). We have

$$V^{k+1} = \underset{V \in \mathbb{R}^{n \times q}}{\arg\min} \{ \langle \nabla_{V} f(U^{k+1}, V^{k}), V - V^{k} \rangle + \frac{1}{\gamma_{2}} \mathbf{D}_{h}((U^{k+1}, V), (U^{k+1}, V^{k})) + g_{2}(V) \}$$

$$= \underset{V \ge 0}{\arg\min} \{ h(U^{k+1}, V) - \langle \nabla_{V} h(U^{k+1}V^{k}) - \gamma_{2} \nabla_{V} f(U^{k+1}, V^{k}), V \rangle \}.$$
(5.11)

Let us consider the normal cone $\mathcal{N}_{V \ge 0}(V^{k+1}) = \{P \in \mathbb{R}^{r \times q} \mid V^{k+1} \odot P = 0, P \le 0\}$ (see [58, Corollary 3.5]), where $V \odot P$ denotes the *Hadamard product* given pointwise by $(V \odot P)_{ij} := V_{ij}P_{ij}$ for $i \in 1, ..., r$ and $j \in 1, ..., q$. The first-order optimality conditions for (5.11) lead to $G^k - (\beta \| V^{k+1} \|_F^2 + \eta_2) V^{k+1} \in \mathcal{N}_{V \ge 0}(V^{k+1})$ with $G^k := \nabla_V h(U^{k+1}V^k) - \gamma_2 \nabla_V f(U^{k+1}, V^k)$.

Now we consider two cases: (i) $G_{ij}^k \le 0$; (ii) $G_{ij}^k > 0$. In Case (i), $P_{ij} = G_{ij}^k - (\beta || V^{k+1} ||_F^2 + \eta_2) V_{ij}^{k+1} \le 0$, hence, $V_{ij}^{k+1} = 0$ (otherwise, $V_{ij}^{k+1} < 0$ which implies $P_{ij} < 0$; this contradicts to the condition $V_{ij}^{k+1} P_{ij} = 0$). In Case (ii), if $V_{ij}^{k+1} = 0$, then $P_{ij} = G_{ij}^k > 0$, which contradicts $P \le 0$, hence, $G_{ii}^k - (\beta || V^{k+1} ||_F^2 + \eta_2) V_{ij}^{k+1} = 0$. Combining both cases, we get

$$(\beta \| V^{k+1} \|_F^2 + \eta_2) V^{k+1} = \mathbf{Proj}_{V \ge 0}(G^k).$$

Denote $t_k = (\beta ||V^{k+1}||_F^2 + \eta_2)$. Then we have $||V^{k+1}||_F^2 = (t_k - \eta_2)/\beta$. Therefore, t_k satisfies

$$t_k^3 - \eta_2 t_k^2 - \beta \|\mathbf{Proj}_{V \ge 0}(G^k)\|_F^2 = 0.$$

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 $\Delta = c^2 + \frac{4}{27}ca^3$, which gives (5.10).

5.2 Preliminary numerical experiments

In this section, we report preliminary numerical results of our experiments on BPALM and two variants of A-BPALM, namely,

- A-BPALM1: the algorithm A-BPALM;
- A-BPALM2: the variant of A-BPALM as described in Remark 4.9.

Since the unconstrained ONMF problem (5.4) involves the quadratic penalty term $\frac{\lambda}{2} ||I_r - VV^T||_F^2$, we also consider a "continuation" variant of these algorithm that starts from some $\lambda > 0$, run one of the above-mentioned algorithms until some stopping criterion holds and save its best point, and then it increases the penalty parameter and run the algorithm with the starting point as the best point of the last call, and it continues the procedure until we stop the algorithm. We refer to this heuristic procedure as *continuation*; see Algorithm 3.

Algorithm 3 Continuation procedure

REQUIRE x⁰ ∈ int dom h, λ > 0, c > 1.
1: repeat
2: starting from x⁰; run one of BPALM, A-BPALM1, or A-BPALM2 to attain an inexact solution x̄ of (5.4);
3: set x⁰ ← x̄, λ ← cλ;
4: until some stopping criterion holds
ENSURE x̄

In our implementation, all the codes were written in MATLAB¹ and runs were performed on a laptop with 1.8 GHz Intel Core i7 CPU and 16 GB RAM. On the basis of our preliminary experiments, we here set $\alpha = \beta = 1$ and $\varepsilon_1 = \varepsilon_2 = 10^{-9}$ to provide the relative smoothness constants as described in (5.7), and the related stepsizes are computed by $\gamma_i = \frac{1}{L_i} - \epsilon$, when ϵ is set as the machine precision. For A-BPALM1 and A-BPALM2, we set $\nu = 2$, and we also set $\overline{L_i^0} = 10^{-4}L_i$. For the continuation versions of our algorithms, we set c = 2 and consider the following notations:

- BPALM-c: BPALM with the continuation;
- A-BPALM1-c: A-BPALM1 with the continuation;
- A-BPALM2-c: A-BPALM2 with the continuation.

Moreover, we implement the following baseline algorithms in our experiments:

- MU: the Multiplicative Update method proposed in [56],
- HALS: the Hierarchical Alternating Least Squares Algorithm proposed in [39];

¹ The codes are publicly available at https://github.com/MasoudAhoo/BPALM

For all of these algorithms, we use the same initialization, namely the successive projection algorithm (SPA) [4]. SPA is guaranteed to recover correctly the factors U and V, given that X = UV where $U \ge 0$, $V \ge 0$ and $VV^T = I_r$ (this is the noiseless case); in fact, ONMF is a special case of separable NMF in the absence of noise. Moreover, this recovery property holds true even in the presence of small bounded noise [35].

In the first experiment, we illustrate the evolution of the objective function (5.3) and the orthogonal error $O_{error} := ||I - V^k(V^k)^T||_F$ obtained by our proposed methods BPALM, A-BPALM1 and A-BPALM2. We generate two synthetic data sets with (p, q, r) = (500, 500, 10) and (p, q, r) = (500, 2000, 10) as follows. We use the MATLAB command rand to generate random nonnegative matrices $U \in \mathbb{R}_+^{p \times r}$ and $R \in \mathbb{R}_+^{p \times q}$, then we generate a random orthogonal nonnegative matrix $V \in \mathbb{R}_+^{r \times q}$. Next, we set X = UV to obtain the *p*-by-*q* orthogonal decomposable matrix *X*, and finally add 5% of noise by $X = X + 0.05 \frac{||X||_F}{||R||_F}R$. We run our algorithms with the fixed penalty parameter $\lambda = 100$ and stopped the algorithms after 15 s. The results are illustrated in (5.3). We make two observations: (i) A-BPALM1 and A-BPALM2 outperform BPALM, A-BPALM2 being the best among them; (ii) within the same running time, BPALM can make more iterations than its line search variants A-BPALM1 and A-BPALM2, and A-BPALM2 runs the least number of iterations compared to the others (Fig. 1).

In the second experiment, we compare our algorithms (with and without continuation) against MU and HALS. We generate 50 data sets with r = 10, and p and q being uniformly chosen within the range [200,1000]. For each (p, q, r), we generate a data set in the same way as in the first experiment. For each data set, we use the same SPA initialization for all algorithms and run each algorithm for 20 s. For the continuation versions BPALM-c, A-BPALM1-c and A-BPALM2-c, we start with $\lambda = 10$, stop the inner algorithms every 2 s to increase λ by factor c = 2; and for BPALM, A-BPALM1 and A-BPALM2, we use a fixed $\lambda = 5120$, which is the last value of λ used for continuation versions.

We report the mean and standard derivation of the final orthogonal error and the fitting error F_{error} := $\frac{\|X - U^k V^k\|_F}{\|X\|_F}$ obtained by each algorithm over 50 runs in Table 1.

We observe that, on average, the continuation versions outperform BPALM variants with fixed λ . Compared to HALS and MU, our algorithms A-BPALM2 and the continuation versions outperform HALS and MU in term of the orthogonal error and are competitive with HALS and MU in term of the fitting error. A-BPALM2 provides the best orthogonal error among the algorithms.

Finally, we report the performance of our algorithms on the Hubble telescope data set which is taken from [50]. In this problem, each row of the matrix X is a vectorized image of the Hubble telescope at a given wavelength for a total of p = 100 wavelengths. Each image contains $q = 128 \times 128$ pixels. Since each pixel in the image contains mostly a single material, it makes sense to use ONMF to cluster the pixel according to the material they contain (see Fig. 2 for an illustration). Since the continuation versions of our algorithms perform better, we here only apply the continuation versions of BPALM, A-BPALM1, and A-BPALM2. We use the SVD-based initialization as in



Fig. 1 A comparison among BPALM, A-BPALM1, and A-BPALM2: Subfigures **a** and **b** stand for function values vs. iterations for the 2 synthetic data with (p, q, r) = (500, 500, 10) (first dataset) and (p, q, r) = (500, 2000, 10) (second dataset), respectively; Subfigures **c** and **d** illustrate the orthogonal error vs. iterations for these data sets, respectively

[52]. We run each algorithm for 100 s. The final outputs of the algorithms, along with the ground true Hubble image, are illustrated in Fig. 2.

We observe on Fig. 2 that A-BPALM1-c and A-BPALM2-c provide slightly better quality images compared to BPALM-c (look for example at the fifth basis image), while HALS and MU fail to cluster the pixels properly as their solutions have a too large orthogonal error O_{error} , more than 10 times larger than A-BPALM1-c and A-BPALM2-c.

Table 1 Mean \pm std of the final F_{error} and O_{error} obtained by each algorithm over 50 runs. In each column, the best results are displayed in bold	Algorithm	F _{error}	0 _{error}
	BPALM A-BPALM1	$2.679 \ 10^{-2} \pm 1.992 \ 10^{-4}$ $2.654 \ 10^{-2} \pm 2.370 \ 10^{-4}$	$\frac{1.26910^{0}\pm9.19110^{-1}}{1.33210^{-2}\pm2.61410^{-2}}$
	A-BPALM2 BPALM-c A-BPALM1-c	$2.573 \ 10^{-2} \pm 1.509 \ 10^{-4}$ $2.582 \ 10^{-2} \pm 2.079 \ 10^{-4}$ $2.568 \ 10^{-2} \pm 1.356 \ 10^{-4}$	$2.133 10^{-3} \pm 4.067 10^{-4}$ $2.395 10^{-3} \pm 3.787 10^{-4}$ $2 419 10^{-3} \pm 4 641 10^{-4}$
	A-BPALM2-c MU	$2.503 10^{-2} \pm 1.371 10^{-4}$ $2.503 10^{-2} \pm 1.540 10^{-4}$	$2.430 10^{-3} \pm 4.921 10^{-4}$ $9.566 10^{-3} \pm 5.528 10^{-4}$
	HALS	$2.41410^{-2}\pm1.43310^{-4}$	$1.22310^{-2}\pm 5.32410^{-4}$

6 Final remarks

We have analyzed two new alternating linearized minimization algorithms called BPALM and A-BPALM for solving the popular nonconvex nonsmooth optimization problem (1.1). To do so, we first introduced the notion of multi-block relative smoothness and verified the fundamental properties of the Bregman proximal alternating linearized mapping. Convergence analysis including the subsequential convergence, global convergence and convergence rate of the proposed algorithms is studied under the framework of multi-block relative smoothness and multi-block kernel functions. We employ BPALM and A-BPALM to solve the orthogonal nonnegative matrix factorization (ONMF) problem. We emphasize that, to the best of our knowledge, BPALM and A-BPALM are the first algorithms with rigorous convergence guarantee for solving ONMF in the literature. Some preliminary numerical tests are provided to illustrate the performance of our algorithms.

Applying (A-)BPALM on other problems and comparing them with state-ofthe-art algorithms is a topic for future work.

Appendix

Lemma 8.1 Let all assumptions of Theorem 4.4 be valid. Then, the following assertions hold:

- (i) $\lim_{k\to\infty} \operatorname{dist}(x^k,\omega(x^0)) = 0;$
- (ii) $\omega(\mathbf{x}^0)$ is a nonempty, compact, and connected set;
- (iii) the objective function φ is finite and constant on $\omega(\mathbf{x}^0)$.



(f) HALS, $F_{error} = 2.380 \times 10^{-2}$, $O_{error} = 1.274 \times 10^{-1}$

Fig. 2 Comparison of BPALM-c, A-BPALM1-c, A-BPALM2-c, MU and HALS on the Hubble image. Algorithms are run for 100 s. In each subfigure, each image corresponds to a row of V that has been reshaped as an image (since each entry corresponds to a pixel; see the description in the text)

Proof Lemma 8.1(i) is a direct consequence of Theorem 4.4, and Lemma 8.1(ii) and Lemma 8.1(iii) can be proved in the same way as [23, Lemma 5(iii)-(iv)].

Lemma 8.2 Let all assumptions of Theorem 4.5 is satisfied. If $\varphi(\mathbf{x}^k) > \varphi^*$, there exists $\varepsilon, \eta > 0$ and the desingularizing function ψ such that

$$\psi'(\varphi(\mathbf{x}^k) - \varphi^*) \operatorname{dist}(0, \partial \varphi(\mathbf{x}^k)) \ge 1 \quad \text{for } k \ge k_0.$$
(8.1)

Proof From Lemma 8.1(ii), the set of limit points $\omega(\mathbf{x}^0)$ of $(\mathbf{x}^k)_{k \in \mathbb{N}}$ is nonempty and compact and φ is finite and constant on $\omega(\mathbf{x}^0)$ due to Lemma 8.1(iii). Moreover,

 $\varphi(\mathbf{x}^k) > \varphi^*$ and the sequence $(\varphi(\mathbf{x}^k))_{k \in \mathbb{N}}$ is decreasing (Proposition 4.1(i)), i.e., there exist $\eta > 0$ and $k_1 \in \mathbb{N}$ such that $\varphi^* < \varphi(\mathbf{x}^k) < \varphi^* + \eta$ for all $k \ge k_1$. For $\varepsilon > 0$, Proposition 4.1(i) implies that there exists $k_2 \in \mathbb{N}$ such that $\operatorname{dist}(\mathbf{x}^k, \omega(\mathbf{x}^0)) < \varepsilon$ for $k \ge k_2$. Setting $k_0 := \max\{k_1, k_2\}$ and according to Fact 2.2, there exist $\varepsilon, \eta > 0$ and a desingularization function ψ such that for any element in

$$\{\mathbf{x}^{k} \mid \operatorname{dist}(\mathbf{x}^{k}, \boldsymbol{\omega}(\mathbf{x}^{0})) < \varepsilon\} \cap [\boldsymbol{\varphi}^{\star} < \boldsymbol{\varphi}(\mathbf{x}^{k}) < \boldsymbol{\varphi}^{\star} + \eta] \quad \text{for } k \ge k_{0},$$

the inequality (8.1) is valid.

We next present the proof of Theorem 4.7.

Proof of Theorem 4.7. The proof has two key parts.

In the first part, we show that there exist c > 0 and $\overline{k} \in \mathbb{N}$ such that for all $k \ge \overline{k}$ the following inequalities hold for i = 1, ..., N:

$$\|x_i^k - x_i^\star\| \le \begin{cases} c \max\{1, \frac{\kappa}{1-\theta}\}\sqrt{\mathcal{S}_{k-1}} & \text{if } \theta \in (0, 1/2], \\ c \frac{\kappa}{1-\theta} \mathcal{S}_{k-1}^{1-\theta} & \text{if } \theta \in (1/2, 1). \end{cases}$$
(8.2)

Let $\varepsilon > 0$ be as described in (4.16) and $x^k \in \mathbf{B}(x^*;\varepsilon)$ for all $k \ge \tilde{k}$ and $\tilde{k} \in \mathbb{N}$. By the definitions of a_k and b_k in (4.15) and using (4.14), we get $a_{k+1} \le \frac{1}{2}a_k + b_k$ for all $k \ge \tilde{k}$. Since $(\varphi(\mathbf{x}^k))_{k\in\mathbb{N}}$ is nonincreasing,

$$\sum_{j=k}^{\infty} a_{j+1} \le \frac{1}{2} \sum_{j=k}^{\infty} (a_j - a_{j+1} + a_{j+1}) + \frac{\widehat{c}N}{2} \sum_{j=k}^{\infty} \left(\Delta_j - \Delta_{j+1} \right) = \frac{1}{2} \sum_{j=k}^{\infty} a_{j+1} + \frac{1}{2} a_k + \frac{\widehat{c}N}{2} \Delta_k.$$

Together with the arithmetic and quadratic mean inequalities, $\psi(S_k) \le \psi(S_{k-1})$, and Proposition 4.1(i), this lead to

$$\sum_{j=k}^{\infty} a_{j+1} \leq a_{k} + \widehat{c}N\Delta_{k} = \sum_{i=1}^{N} \|x_{i}^{k} - x_{i}^{k-1}\| + \widehat{c}N\psi(\mathcal{S}_{k}) \leq \sqrt{N}\sqrt{\sum_{i=1}^{N} \|x_{i}^{k} - x_{i}^{k-1}\|^{2}} + \widehat{c}N\psi(\mathcal{S}_{k})$$

$$\leq \sqrt{2N}\max\{\frac{1}{\sqrt{\sigma_{1}}}, \dots, \frac{1}{\sqrt{\sigma_{N}}}\}\sqrt{\sum_{i=1}^{N} \mathbf{D}_{h}(\mathbf{x}^{k-1,i}, \mathbf{x}^{k-1,i-1})} + \widehat{c}N\psi(\mathcal{S}_{k})$$

$$\leq \sqrt{\frac{2N}{\rho}}\max\{\frac{1}{\sqrt{\sigma_{1}}}, \dots, \frac{1}{\sqrt{\sigma_{N}}}\}\sqrt{\mathcal{S}_{k-1} - \mathcal{S}_{k}} + \widehat{c}N\psi(\mathcal{S}_{k-1}).$$
(8.3)

On the other hand, for i = 1, ..., N, we have

$$\|x_i^k - x_i^{\star}\| \le \|x_i^{k+1} - x_i^k\| + \|x_i^{k+1} - x_i^{\star}\| \le \dots \le \sum_{j=k}^{\infty} \|x_i^{j+1} - x_i^j\|.$$

This inequality, together with (8.3), yields

$$\sum_{i=1}^{N} \|x_i^k - x_i^\star\| \le \sqrt{\frac{2N}{\rho}} \max\{\frac{1}{\sqrt{\sigma_1}}, \dots, \frac{1}{\sqrt{\sigma_N}}\} \sqrt{\mathcal{S}_{k-1} - \mathcal{S}_k} + \widehat{c} N \psi(\mathcal{S}_{k-1}),$$

leading to

$$\|x_{i}^{k} - x_{i}^{\star}\| \le c \max\{\sqrt{\mathcal{S}_{k-1}}, \psi(\mathcal{S}_{k-1})\} \quad i = 1, \dots, N,$$
(8.4)

Deringer

712

where $c := \sqrt{\frac{2N}{\rho}} \max\left\{\frac{1}{\sqrt{\sigma_1}}, \dots, \frac{1}{\sqrt{\sigma_N}}\right\} + \hat{c}N$ and $\psi(s) := \frac{\kappa}{1-\theta}s^{1-\theta}$. Let us consider the nonlinear equation

$$\sqrt{\mathcal{S}_{k-1}} - \frac{\kappa}{1-\theta} \mathcal{S}_{k-1}^{1-\theta} = 0,$$

which has a solution at $S_{k-1} = \left(\frac{(1-\theta)}{\kappa}\right)^{\frac{2}{1-2\theta}}$. Form the monotonicity of S_k , there exists $\hat{k} \in \mathbb{N}$ such that for $k \ge \hat{k}$ (8.4) holds and

$$S_{k-1} \leq \left(\frac{1-\theta}{\kappa}\right)^{\frac{2}{1-2\theta}}.$$

We now consider two cases: (a) $\theta \in (0, 1/2]$; (b) $\theta \in (1/2, 1)$. In Case (a), if $\theta \in (0, 1/2)$, then $\psi(S_{k-1}) \leq \sqrt{S_{k-1}}$. If $\theta = 1/2$, then $\psi(S_{k-1}) = \frac{\kappa}{1-\theta}\sqrt{S_{k-1}}$, i.e.,

$$\max\{\sqrt{\mathcal{S}_{k-1}}, \psi(\mathcal{S}_{k-1})\} = \max\{1, \frac{\kappa}{1-\theta}\}\sqrt{\mathcal{S}_{k-1}}$$

Therefore, it holds that $\max\{\sqrt{S_{k-1}}, \psi(S_{k-1})\} \le \max\{1, \frac{\kappa}{1-\theta}\}\sqrt{S_{k-1}}$. In Case (b), we have that

$$\psi(\mathcal{S}_{k-1}) \geq \sqrt{\mathcal{S}_{k-1}},$$

i.e., $\max\{\sqrt{S_{k-1}}, \psi(S_{k-1})\} = \frac{\kappa}{1-\theta} S_{k-1}^{1-\theta}$. Then, it follows from (8.4) that (8.2) holds for all $k \ge k := \max\{\hat{k}, \hat{k}\}$.

In the second part of the proof, we will show the assertions in the statement of the theorem. For $(\mathcal{G}_i^k, \ldots, \mathcal{G}_N^k) \in \partial \varphi(\mathbf{x}^k)$ as defined in Proposition 4.3, by Proposition 4.1(i), we infer

$$\begin{split} \mathcal{S}_{k-1} - \mathcal{S}_{k} &= \varphi(x^{k-1}) - \varphi(x^{k}) \ge \rho \sum_{i=1}^{N} \mathbf{D}_{h}(x^{k-1,i}, x^{k-1,i-1}) \ge \frac{\rho}{2} \sum_{i=1}^{N} \sigma_{i} \|x_{i}^{k} - x_{i}^{k-1}\|^{2} \\ &\ge \frac{\rho}{2N} \min\{\sigma_{1}, \dots, \sigma_{N}\} \left(\sum_{i=1}^{N} \|x_{i}^{k} - x_{i}^{k-1}\| \right)^{2} \ge \frac{\rho}{2N\overline{c}^{2}} \min\{\sigma_{1}, \dots, \sigma_{N}\} \|(\mathcal{G}_{i}^{k}, \dots, \mathcal{G}_{N}^{k})\|^{2} \\ &\ge \frac{\rho}{2N\overline{c}^{2}} \min\{\sigma_{1}, \dots, \sigma_{N}\} \mathbf{dist}(0, \partial\varphi(x^{k}))^{2} \ge \frac{\rho}{2N\overline{c}^{2}\kappa^{2}} \min\{\sigma_{1}, \dots, \sigma_{N}\} \mathcal{S}_{k-1}^{2\theta} = \widetilde{c} \mathcal{S}_{k-1}^{2\theta}, \end{split}$$

with $\tilde{c} := \frac{\rho}{2N\bar{c}^2\kappa^2} \min\{\sigma_1, \dots, \sigma_N\}$ and for all $k \ge \bar{k}$. Hence, all assumptions of Fact 4.6 hold with $\alpha = 2\theta$. Therefore, our results follows from this fact and (8.2).

Acknowledgements We would like to thank the anonymous reviewers for their insightful comments that helped improve the paper; in particular, one of the reviewers gave a suggestion that leads to the kernel function in Proposition 5.1. The first author is grateful to Andreas Themelis for his useful comments and discussions on the paper. MA and PP acknowledge the support by the Research Foundation Flanders (FWO) research projects G086518N andG086318N; Research Council KU Leuven C1 project No. C14/18/068; Fonds de la Recherche Scientifique - FNRS andthe Fonds Wetenschappelijk Onderzoek - Vlaanderen (FWO) under EOS project no 30468160 (SeLMA). LTKH andNG also acknowledge the support by the European Research Council (ERC starting grant no 679515).

References

- Ahookhosh, M.: Accelerated first-order methods for large-scale convex optimization: nearly optimal complexity under strong convexity. Math. Methods Oper. Res. 89(3), 319–353 (2019)
- Ahookhosh, M., Hien, L.T.K., Gillis, N., Patrinos, P.: A block inertial bregman proximal algorithm for nonsmooth nonconvex problems with application to symmetric nonnegative matrix tri-factorization. J. Optim. Theory Appl. (2021)
- Ahookhosh, M., Themelis, A., Patrinos, P.: A Bregman forward-backward linesearch algorithm for nonconvex composite optimization: superlinear convergence to nonisolated local minima. SIAM J. Optim. 31(1), 653–685 (2021)
- Araújo, U., Saldanha, B., Galvão, R., Yoneyama, T., Chame, H., Visani, V.: The successive projections algorithm for variable selection in spectroscopic multicomponent analysis. Chemometr. Intell. Lab. Syst. 57(2), 65–73 (2001)
- Armijo, L.: Minimization of functions having Lipschitz continuous first partial derivatives. Pac. J. Math. 16(1), 1–3 (1966)
- Attouch, H., Bolte, J., Redont, P., Soubeyran, A.: Alternating proximal algorithms for weakly coupled convex minimization problems. applications to dynamical games and PDE's. J. Convex Anal. 15(3), 485 (2008)
- 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. Math. Oper. Res. 35(2), 438–457 (2010)
- 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. Math. Program. 137(1), 91–129 (2013)
- Attouch, H., Redont, P., Soubeyran, A.: A new class of alternating proximal minimization algorithms with costs-to-move. SIAM J. Optim. 18(3), 1061–1081 (2007)
- Attouch, H., Soubeyran, A.: Inertia and reactivity in decision making as cognitive variational inequalities. J. Conv. Anal. 13(2), 207 (2006)
- 11. Auslender, A.: Optimisation méthodes numériques. Mason, Paris (1976)
- 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. J. Optim. Theory Appl. 182, 1068–1087 (2019)
- 13. Bauschke, H.H., Bolte, J., Teboulle, M.: A descent lemma beyond Lipschitz gradient continuity: first-order methods revisited and applications. Math. Oper. Res. **42**(2), 330–348 (2016)
- Bauschke, H.H., Dao, M.N., Lindstrom, S.B.: Regularizing with Bregman–Moreau envelopes. SIAM J. Optim. 28(4), 3208–3228 (2018)
- Beck, A., Pauwels, E., Sabach, S.: The cyclic block conditional gradient method for convex optimization problems. SIAM J. Optim. 25(4), 2024–2049 (2015)
- 16. Beck, A., Sabach, S., Teboulle, M.: An alternating semiproximal method for nonconvex regularized structured total least squares problems. SIAM J. Matrix Anal. Appl. **37**(3), 1129–1150 (2016)
- Beck, A., Teboulle, M.: A fast iterative shrinkage-thresholding algorithm for linear inverse problems. SIAM J. Imag. Sci. 2(1), 183–202 (2009)
- Beck, A., Tetruashvili, L.: On the convergence of block coordinate descent type methods. SIAM J. Optim. 23(4), 2037–2060 (2013)
- 19. Bertsekas, D.P., Tsitsiklis, J.N.: Parallel and Distributed Computation: Numerical Methods. Prentice-Hall, Inc., Hoboken (1989)
- Bolte, J., Daniilidis, A., Lewis, A.: The Łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. SIAM J. Optim. 17(4), 1205–1223 (2007)
- Bolte, J., Daniilidis, A., Lewis, A., Shiota, M.: Clarke subgradients of stratifiable functions. SIAM J. Optim. 18(2), 556–572 (2007)
- Bolte, J., Daniilidis, A., Ley, O., Mazet, L.: Characterizations of Łojasiewicz inequalities: subgradient flows, talweg, convexity. Trans. Am. Math. Soc. 362(6), 3319–3363 (2010)
- 23. Bolte, J., Sabach, S., Teboulle, M.: Proximal alternating linearized minimization for nonconvex and nonsmooth problems. Math. Program. **146**(1–2), 459–494 (2014)
- Bolte, J., Sabach, S., Teboulle, M., Vaisbourd, Y.: First order methods beyond convexity and Lipschitz gradient continuity with applications to quadratic inverse problems. SIAM J. Optim. 28(3), 2131–2151 (2018)

- Boţ, R.I., Csetnek, E.R.: An inertial Tseng's type proximal algorithm for nonsmooth and nonconvex optimization problems. J. Optim. Theory Appl. 171(2), 600–616 (2016)
- Bot, R.I., Nguyen, D.K.: The proximal alternating direction method of multipliers in the nonconvex setting: convergence analysis and rates. Math. Oper. Res. 45(2), 682–712 (2020)
- 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 Comput. Math. Math. Phys. 7(3), 200–217 (1967)
- Chen, G., Teboulle, M.: Convergence analysis of a proximal-like minimization algorithm using Bregman functions. SIAM J. Optim. 3(3), 538–543 (1993)
- 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, Hoboken (2009)
- Combettes, P.L., Pesquet, J.C.: Stochastic quasi-Fejér block-coordinate fixed point iterations with random sweeping. SIAM J. Optim. 25(2), 1221–1248 (2015)
- 31. Van den Dries, L.: Tame Topology and o-Minimal Structures, vol. 248. Cambridge University Press, Cambridge (1998)
- Fercoq, O., Bianchi, P.: A coordinate-descent primal-dual algorithm with large step size and possibly nonseparable functions. SIAM J. Optim. 29(1), 100–134 (2019)
- Fu, X., Huang, K., Sidiropoulos, N.D., Ma, W.K.: Nonnegative matrix factorization for signal and data analytics: Identifiability, algorithms, and applications. IEEE Signal Process. Mag. 36(2), 59–80 (2019)
- 34. Gillis, N.: The why and how of nonnegative matrix factorization. Regular. Optim. Kernels Support Vector Mach. **12**(257), 257–291 (2014)
- Gillis, N., Vavasis, S.A.: Fast and robust recursive algorithmsfor separable nonnegative matrix factorization. IEEE Trans. Pattern Anal. Mach. Intell. 36(4), 698–714 (2013)
- Grippo, L., Sciandrone, M.: On the convergence of the block nonlinear Gauss-Seidel method under convex constraints. Operat. Res. Lett. 26(3), 127–136 (2000)
- 37. Hanzely, F., Richtárik, P.: Fastest rates for stochastic mirror descent methods. arXiv:1803.07374 (2018)
- 38. Hanzely, F., Richtarik, P., Xiao, L.: Accelerated bregman proximal gradient methods for relatively smooth convex optimization. Comput Optim Appl **22**, 1–36 (2021)
- Kimura, K., Tanaka, Y., Kudo, M.: A fast hierarchical alternating least squares algorithm for orthogonal nonnegative matrix factorization. In: D. Phung, H. Li (eds.) Proceedings of the Sixth Asian Conference on Machine Learning, *Proceedings of Machine Learning Research*, vol. 39, pp. 129–141. PMLR, Nha Trang City, Vietnam (2015). http://proceedings.mlr.press/v39/kimura14.html
- Kurdyka, K.: On gradients of functions definable in o-minimal structures. Annales de l'institut Fourier 48(3), 769–783 (1998)
- 41. Latafat, P., Freris, N.M., Patrinos, P.: A new randomized block-coordinate primal-dual proximal algorithm for distributed optimization. IEEE Trans. Autom. Cont. **64**(10), 4050–4065 (2019)
- Latafat, P., Themelis, A., Patrinos, P.: Block-coordinate and incremental aggregated proximal gradient methods for nonsmooth nonconvex problems. Math. Program. 1–30. arxiv.org/abs/1906.10053 (2021)
- Li, Q., Zhu, Z., Tang, G., Wakin, M.B.: Provable Bregman-divergence based methods for nonconvex and non-Lipschitz problems. arXiv preprint arXiv:1904.09712 (2019)
- 44. Łojasiewicz, S.: Une propriété topologique des sous-ensembles analytiques réels. Les équations aux dérivées partielles pp. 87–89 (1963)
- 45. Łojasiewicz, S.: Sur la géométrie semi- et sous- analytique. Annales de l'institut Fourier 43(5), 1575– 1595 (1993)
- 46. Lu, H., Freund, R.M., Nesterov, Y.: Relatively smooth convex optimization by first-order methods, and applications. SIAM J. Optim. **28**(1), 333–354 (2018)
- Mukkamala, M.C., Ochs, P., Pock, T., Sabach, S.: Convex-concave backtracking for inertial bregman proximal gradient algorithms in nonconvex optimization. SIAM J. Math. Data Sci. 2(3), 658–682 (2020)
- Nesterov, Y.: Efficiency of coordinate descent methods on huge-scale optimization problems. SIAM J. Optim. 22(2), 341–362 (2012)
- 49. Nesterov, Y.: Gradient methods for minimizing composite functions. Math. Program. **140**(1), 125–161 (2013)
- Pauca, V.P., Piper, J., Plemmons, R.J.: Nonnegative matrix factorization for spectral data analysis. Linear Algebra Appl. 416(1), 29–47 (2006)

- Pock, T., Sabach, S.: Inertial proximal alternating linearized minimization (iPALM) for nonconvex and nonsmooth problems. SIAM J. Imag .Sci. 9(4), 1756–1787 (2016)
- 52. Pompili, F., Gillis, N., Absil, P.A., Glineur, F.: Two algorithms for orthogonal nonnegative matrix factorization with application to clustering. Neurocomputing **141**, 15–25 (2014)
- Razaviyayn, M., Hong, M., Luo, Z.Q.: A unified convergence analysis of block successive minimization methods for nonsmooth optimization. SIAM J. Optim. 23(2), 1126–1153 (2013)
- Richtárik, P., Takáč, M.: Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. Math. Program. 144(1–2), 1–38 (2014)
- Rockafellar, R.T., Wets, R.J.B.: Variational Analysis, vol. 317. Springer Science & Business Media, Berlin (2011)
- Choi, S.: Algorithms for orthogonal nonnegative matrix factorization. In: 2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence), pp. 1828–1832 (2008)
- 57. Shefi, R., Teboulle, M.: On the rate of convergence of the proximal alternating linearized minimization algorithm for convex problems. EURO J. Comput. Optim. 4(1), 27–46 (2016)
- Tam, M.K.: Regularity properties of non-negative sparsity sets. J. Math. Anal. Appl. 447(2), 758– 777 (2017)
- 59. Teboulle, M.: A simplified view of first order methods for optimization. Math. Program. **170**(1), 67–96 (2018)
- Themelis, A., Ahookhosh, M., Patrinos, P.: On the acceleration of forward-backward splitting via an inexact Newton method. In: Luke, R., Bauschke, H., Burachik, R. (eds.) Splitting Algorithms, Modern Operator Theory, and Applications, pp. 363–412. Springer, Berlin (2019)
- Tseng, P.: Convergence of a block coordinate descent method for nondifferentiable minimization. J. Optim. Theory Appl. 109(3), 475–494 (2001)
- 62. Tseng, P., Yun, S.: A coordinate gradient descent method for nonsmooth separable minimization. Math. Program. **117**(1-2), 387-423 (2009)
- 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)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Authors and Affiliations

Masoud Ahookhosh¹ · Le Thi Khanh Hien² · Nicolas Gillis² · Panagiotis Patrinos³

Le Thi Khanh Hien ThiKhanhHien.LE@umons.ac.be

Nicolas Gillis nicolas.gillis@umons.ac.be

Panagiotis Patrinos panos.patrinos@esat.kuleuven.be

- ¹ Department of Mathematics, University of Antwerp, Middelheimlaan 1, 2020 Antwerp, Belgium
- ² Department of Mathematics and Operational Research, Faculté polytechnique, Université de Mons. Rue de Houdain 9, 7000 Mons, Belgium
- ³ Department of Electrical Engineering (ESAT-STADIUS), KU Leuven, Kasteelpark Arenberg 10, 3001 Leuven, Belgium