

Inertial Alternating Direction Method of Multipliers for Non-Convex Non-Smooth Optimization

Le Thi Khanh Hien · Duy Nhat Phan ·
Nicolas Gillis

the date of receipt and acceptance should be inserted later

Abstract In this paper, we propose an algorithmic framework, dubbed inertial alternating direction methods of multipliers (iADMM), for solving a class of nonconvex nonsmooth multiblock composite optimization problems with linear constraints. Our framework employs the general minimization-majorization (MM) principle to update each block of variables so as to not only unify the convergence analysis of previous ADMM that use specific surrogate functions in the MM step, but also lead to new efficient ADMM schemes. To the best of our knowledge, in the *nonconvex nonsmooth* setting, ADMM used in combination with the MM principle to update each block of variables, and ADMM combined with *inertial terms for the primal variables* have not been studied in the literature. Under standard assumptions, we prove the subsequential convergence and global convergence for the generated sequence of iterates. We illustrate the effectiveness of iADMM on a class of nonconvex low-rank representation problems.

Keywords alternating direction methods of multipliers · majorization minimization · inertial block coordinate method · acceleration by extrapolation · low-rank representation

L. T. K. Hien and D. N. Phan contributed equally to this work.
L. T. K. Hien finished this work when she was at the University of Mons, Belgium.

L. T. K. Hien
Huawei Belgium Research Center, 3001 Leuven, Belgium
E-mail: let.hien@huawei.com

D. N. Phan
Dynamic Decision Making Laboratory, Carnegie Mellon University, USA
E-mail: dnphan@andrew.cmu.edu

N. Gillis
Department of Mathematics and Operational Research, Faculté Polytechnique, Université de Mons, Rue de Houdain 9, 7000 Mons, Belgium
E-mail: nicolas.gillis@umons.ac.be

1 Introduction

In this paper, we consider the following nonconvex minimization problem with linear constraints

$$\min_{x,y} F(x_1, \dots, x_s) + h(y) \quad \text{such that} \quad \sum_{i=1}^s \mathcal{A}_i x_i + \mathcal{B}y = b, \quad (1)$$

where $y \in \mathbb{R}^q$, $x_i \in \mathbb{R}^{n_i}$, $x := [x_1; \dots; x_s] \in \mathbb{R}^n$, $n = \sum_{i=1}^s n_i$, \mathcal{A}_i is a linear map from \mathbb{R}^{n_i} to \mathbb{R}^m , \mathcal{B} is a linear map from \mathbb{R}^q to \mathbb{R}^m , $b \in \mathbb{R}^m$, $h : \mathbb{R}^q \rightarrow \mathbb{R}$ is a differentiable function, and $F(x) = f(x) + \sum_{i=1}^s g_i(x_i)$, where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a nonconvex nonsmooth function and $g_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R} \cup \{+\infty\}$ are proper lower semi-continuous functions for $i = 1, 2, \dots, s$. We assume that F satisfies $\partial F(x) = \partial_{x_1} F(x) \times \dots \times \partial_{x_s} F(x)$, where ∂F denote the limiting subdifferential of F (see the definition in Appendix A). Note that this condition is satisfied when f is a sum of a continuously differentiable function and a block separable function; see [2, Proposition 2.1].

Notation. We denote $[s] := \{1, \dots, s\}$. For the \mathbf{p} -dimensional Euclidean space $\mathbb{R}^{\mathbf{p}}$, we use $\langle \cdot, \cdot \rangle$ to denote the inner product, and $\|\cdot\|$ to denote the corresponding induced norm. For a linear map \mathcal{M} , \mathcal{M}^* denotes the adjoint linear map with respect to the inner product, and $\|\mathcal{M}\|$ is the induced operator norm of \mathcal{M} . We use \mathcal{I} to denote the identity map. For a positive definite self-adjoint operator \mathcal{Q} , we denote $\|x\|_{\mathcal{Q}}^2 := \langle x, \mathcal{Q}x \rangle$. We denote the smallest eigenvalue of a symmetric linear self-map (that is, $\mathcal{M} = \mathcal{M}^*$) by $\lambda_{\min}(\mathcal{M})$. We use $Im(\mathcal{B})$ to denote the image of \mathcal{B} .

1.1 Nonconvex low-rank representation problem

Low-rank matrix approximations play a central role in various fields of computer science and applied mathematics, and are used in many applications, e.g., recommender systems [27], topic modeling [29], system identification [38], graph clustering [57], compression and denoising [55], to cite a few; see also below for other examples. Given a data matrix, D , the goal of low-rank matrix approximations is to find a nearby low-rank matrix, X . The low-rank assumption is valid in many applications as there are typically redundancy and correlations within large data sets; see, e.g., [47, 56] and the references therein.

In this paper, we will illustrate the use of (1) on the following generalized nonconvex low-rank representation problem: given a data matrix $D \in \mathbb{R}^{d \times n}$, solve

$$\min_{X,Y,Z} \sum_{i=1}^{\min(m,n)} r_1(\sigma(X)) + r_2(Y) + r_3(Z) \quad \text{such that} \quad D = A_1 X + Y A_2 + Z, \quad (2)$$

where $X \in \mathbb{R}^{m \times n}$, $Y \in \mathbb{R}^{d \times q}$, $Z \in \mathbb{R}^{d \times n}$, $A_1 \in \mathbb{R}^{d \times m}$, $A_2 \in \mathbb{R}^{q \times n}$, $\sigma(X)$ is the vector of singular values of X , r_1 is an increasing concave function to promote X to be of low rank (by promoting the sparsity of $\sigma(X)$), r_2 is a regularization function, and r_3 is a function that models some noise; e.g., taking $r_3(Z) = \frac{1}{2} \|Z\|_F^2$ when Z models Gaussian noise. Problem 2 generalize low-rank matrix approximations, taking A_1 as the identity matrix and $A_2 = 0$, so that $D = X + Z$ where X is low rank, and Z models the noise.

In particular, Problem (2) generalizes the following machine learning problems:

- (i) Let $r_1(t) = t^\chi$ with $0 < \chi \leq 1$, $r_2(Y) = \sum_{i=1}^{q-1} \|Y_i - Y_{i+1}\|$ where Y_i is the i -th column of Y , and let A_1 and A_2 be the identity matrices so that Problem (2) decomposes the data matrix D into the sum of three components, X , Y and Z . An application is video surveillance where each column of D is a vectorized image of a video frame, X is a low-rank matrix that plays the role of the background, Y is the foreground that has small variations between its columns (such as slowly moving objectives), and Z represents some noise [58].
- (ii) When A_1 and A_2 are identity matrices, $r_1(t) = t$, and $r_2(Y) = \lambda \|Y\|_1$ for some constant $\lambda > 0$, Problem (2) recovers the robust principal component analysis (robust PCA) model, see, e.g., [12]. Robust PCA decomposes the input matrix D as the sum of a low-rank matrix X , a sparse matrix Y modeling gross corruptions and outliers, and an additional noise matrix Z (e.g., $r_3(Z)$ is a multiple of $\|Z\|_F^2$ to model Gaussian noise). Robust PCA is also used for foreground-background separation in video surveillance.
- (iii) When $r_1(t) = t$ and $r_2(Y) = \|Y\|_*$, Problem (2) is the latent low-rank representation problem [34]. In [34], authors used $A_1 = DP_1$ and $A_2 = P_2^*D$, where P_1 and P_2 are computed by orthogonalizing the columns of D^* and D , respectively. We will use this application to illustrate the effectiveness of our proposed framework, iADMM, in Section 3.

Other applications of Problem (1) include statistical learning, see, e.g., [4, 59], and minimization on compact manifolds, see, e.g., [28, 60].

1.2 Motivation and related works

Let $\mathcal{A} := [\mathcal{A}_1 \dots \mathcal{A}_s]$ and $\mathcal{A}x := \sum_{i=1}^s \mathcal{A}_i x_i \in \mathbb{R}^m$. The augmented Lagrangian for Problem (1) is

$$\mathcal{L}(x, y, \omega) := F(x) + h(y) + \langle \omega, \mathcal{A}x + \mathcal{B}y - b \rangle + \frac{\beta}{2} \|\mathcal{A}x + \mathcal{B}y - b\|^2, \quad (3)$$

where $\beta > 0$ is a penalty parameter. ADMM was first introduced by [18] and [17]. It has recently become popular because of its efficacy in solving emerging large-scale problems in machine learning and computer vision [9, 50, 64, 65, 66]. For simplicity, let us describe the iteration scheme of a classical ADMM for solving Problem (1) with 2 blocks x and y :

$$x^{k+1} \in \underset{x}{\operatorname{argmin}} \mathcal{L}(x, y^k, \omega^k), \quad (4a)$$

$$y^{k+1} \in \underset{y}{\operatorname{argmin}} \mathcal{L}(x^{k+1}, y, \omega^k), \quad (4b)$$

$$\omega^{k+1} = \omega^k + \beta(\mathcal{A}x^{k+1} + \mathcal{B}y^{k+1} - b). \quad (4c)$$

For a multi-block problem, with $s > 1$, the scheme is similar, see, e.g., [58]. The update of x in (4a) (a similar discussion is applicable to (4b)) can be rewritten as $x^{k+1} \in \operatorname{argmin}_x F(x) + \varphi^k(x)$, where

$$\varphi^k(x) = \frac{\beta}{2} \|\mathcal{A}x + \mathcal{B}y^k - b\|^2 + \langle \omega^k, \mathcal{A}x + \mathcal{B}y^k - b \rangle. \quad (5)$$

Solving the subproblem (4a) is usually very expensive especially when F is not smooth. A remedy is minimizing a suitable surrogate function of $\mathcal{L}(\cdot, y^k, \omega^k)$ that allows a more efficient update for x . For example, since $\varphi^k(x)$ is upper bounded by

$$\hat{\varphi}(x) = \varphi^k(x^k) + \langle \nabla \varphi^k(x^k), x - x^k \rangle + \frac{\kappa\beta}{2} \|x - x^k\|^2, \quad (6)$$

where $\kappa \geq \|\mathcal{A}^* \mathcal{A}\|$ (because $\nabla \varphi^k(x)$ is $\beta \|\mathcal{A}^* \mathcal{A}\|$ -Lipschitz continuous), x can be updated by $x^{k+1} \in \operatorname{argmin}_x F(x) + \hat{\varphi}(x)$, which leads to the linearized ADMM method, see [32, 61]. This update has a closed form for some nonsmooth F ; see [43]. When $F = f + g$ and f is L_f -smooth then we can also use the upper bound $\hat{F}(x) = f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{L_f}{2} \|x - x^k\|^2 + g(x)$ of F to obtain $x^{k+1} \in \operatorname{argmin}_x \hat{F}(x) + \hat{\varphi}(x)$. This leads to the proximal linearized ADMM method, see [8, 35]. We note that $\mathcal{L}(\cdot, y^k, \omega^k)$ is always upper bounded by $\mathcal{L}(\cdot, y^k, \omega^k) + \mathbf{D}_\phi(x, x^k)$, where \mathbf{D}_ϕ is the Bregman distance associated with a continuously differentiable convex function ϕ on \mathbb{R}^n :

$$\mathbf{D}_\phi(a, b) := \phi(a) - \phi(b) - \langle \nabla \phi(b), a - b \rangle, \forall a, b \in \mathbb{R}^n. \quad (7)$$

For example, if $\phi(x) = \|x\|_{\mathcal{Q}}^2 = \langle x, \mathcal{Q}x \rangle$ then $\mathbf{D}_\phi(a, b) = \|a - b\|_{\mathcal{Q}}^2$. This upper bound leads to proximal ADMM, see [15, 30]. The above mentioned upper bound functions are specific examples of surrogate functions for $\mathcal{L}(\cdot, y^k, \omega^k)$ (see Definition 1 at page 6 for the definition of a surrogate function) while each method of updating x corresponds to a majorization-minimization (MM) step that minimizes the corresponding majorizer/surrogate function (see [52] for more specific examples of the MM procedure). In the convex setting (that is, $f(\cdot, \cdot)$ is convex), [13] and [23] use the MM principle to unify and generalize the convergence analysis of many ADMM for multi-blocks problems (that is, $s > 1$). However, ADMM with the MM principle has not been studied for the *nonconvex* problem (1), to the best of our knowledge.

When the linear coupling constraint is absent, the block coordinate descent (BCD) method is a standard approach to solve (1). Razaviyayn et al. [46] proposed the block successive upper-bound minimization (BSUM) framework that employs the MM principle in each block update. By employing suitable surrogate functions in each block update, BSUM recovers the typical BCD methods, for example of [19, 22, 45, 53, 5, 7, 54]. In the non-convex setting, BCD methods with inertial terms¹ have also been studied, and have showed significant improvement in their practical performance; see, e.g., [41] for inertial BCD methods with heavy-ball acceleration, [62, 63] for inertial BCD methods with Nesterov-type acceleration, and [44, 20] for inertial BCD methods that use two extrapolation points. Recently, the authors in [21] proposed a general inertial block MM framework for solving (1) without the linear coupling constraint. To the best of our knowledge, inertial ADMM with *Nesterov-type acceleration for the primal variables* have not been studied in the nonconvex case of (1) although some variants of ADMM with inertial terms for the primal variables have been analysed in the convex case (that is, when both F and h are convex); see e.g., [11, 31, 42].

¹ We use in this paper the terminology “inertial” to mean that an inertial term that involves the current iterate and the previous iterates is added to the objective of the subproblem to update each block, see [21].

Recently, [51] proposes ADMM with inertial term for the dual variable, see the description in [51, Expression (17)]. We would like to remark that we realize a gap² in the proof of [51, Lemma 5]. Let us also mention stochastic ADMM methods for solving Problem (1) in which the objective is in expectation formulation, see, e.g., [24, 25], which is out of the scope of this paper.

1.3 Contribution and Organization

In this paper, we propose iADMM, a framework of inertial alternating direction methods of multipliers, for solving the nonconvex nonsmooth problem (1). When no extrapolation is used, iADMM becomes a general ADMM framework that employs the minimization-majorization principle in each block update. For the first time in the *nonconvex* nonsmooth setting of Problem (1), we study ADMM and its inertial version combined with the MM principle when updating each block of variables. Moreover, our framework allows to use an over-relaxation parameter $\alpha \in (0, 2)$ to set $\alpha\beta$ as the constant stepsize for updating the dual variable ω . Note that $\alpha = 1$, see, e.g., [23, 30, 58], or $\alpha \in (0, \frac{1+\sqrt{5}}{2})$, see, e.g., [16, 65], are the standard choices in the nonconvex setting. Recently, [8] proposed proximal ADMM that use $\alpha \in (0, 2)$ for solving a special case of the nonconvex Problem (1) with $s = 1$ and $\mathcal{A} = -\mathcal{I}$.

Under standard assumptions and $\alpha \in (0, 2)$, we analyse the subsequential convergence for the generated sequence of iADMM and ADMM. When $F(x) + h(y)$ satisfies the Kurdyka-Lojasiewicz (KL) property and $\alpha = 1$, we prove the global convergence and provide the convergence rate for the generated sequence. We would like to emphasize that although proving convergence towards a critical point has become a typical task when considering the nonconvex nonsmooth Problem (1), see e.g., [8, 30, 58], the techniques to accomplish this task heavily depend on the considered algorithms and the involved assumptions. As far as we are aware of, this has not been done for ADMM used in combination with the MM principle and inertial terms for the primal variables.

Finally, we apply the proposed framework to solve a class of nonconvex low-rank representation to illustrate the efficacy of iADMM. More specifically, in order to illustrate the effect of MM procedure in Algorithm 1, we use suitable surrogate functions such that each block of variables has a close-form update rule (thus, we do not need to use an outer optimization solver to find a solution for the corresponding subproblem), see details in Section 3.1. In order to illustrate the acceleration effect of Algorithm 1, we also employ inertial terms and the extrapolation parameters are appropriately chosen to guarantee *a global convergence*, see details in Section 3.2. Indeed, the numerical results presented in Section 3.3 (see also Appendix C and Appendix D) empirically show the significant acceleration effect of using inertial terms.

The paper is organized as follows. In the next section, we describe the proposed method, iADMM, and analyse its convergence properties. In Section 3, we report the numerical results of iADMM on a class of nonconvex low-rank representation problems. We conclude the paper in Section 4. All the technical proofs are presented in Appendix B.

² Specifically, the second equality of [51, Expression (51)] is not correct.

2 An inertial ADMM framework

In this section, we describe the iADMM framework and prove its subsequential and global convergence. Throughout the paper, we make the following assumptions that are standard for studying Problem (1) and the convergence of ADMMs in the nonconvex setting, see for example [58, 8, 30].

Assumption 1 (i) $\sigma_{\mathcal{B}} := \lambda_{\min}(\mathcal{B}\mathcal{B}^*) > 0$.

(ii) $F(x) + h(y)$ is lower bounded.

(iii) The function h is L_h -smooth, that is, ∇h is L_h -Lipschitz continuous.

2.1 Description of iADMM

Let us first formally define a surrogate function. Some examples were given in the introduction. More examples can be found in [37, 46, 21].

Definition 1 (Surrogate function) Let $\mathcal{X} \subseteq \mathbb{R}^n$. A function $u : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is called a surrogate function of a function f on \mathcal{X} if the following two conditions are satisfied:

(a) $u(z, z) = f(z)$ for all $z \in \mathcal{X}$, and (b) $u(x, z) \geq f(x)$ for all $x, z \in \mathcal{X}$.

As we are considering multi-block problems, we need the following definition of a block surrogate function, which is a generalization of Definition 1.

Definition 2 (Block surrogate function) Let $\mathcal{X}_i \subseteq \mathbb{R}^{n_i}$, $\mathcal{X} \subseteq \mathbb{R}^n$. A function $u_i : \mathcal{X}_i \times \mathcal{X} \rightarrow \mathbb{R}$ is called a block i surrogate function of f on \mathcal{X} if the following conditions are satisfied:

(a) $u_i(z_i, z) = f(z)$ for all $z \in \mathcal{X}$,

(b) $u_i(x_i, z) \geq f(x_i, z_{\neq i})$ for all $x_i \in \mathcal{X}_i$ and $z \in \mathcal{X}$,

where $(x_i, z_{\neq i})$ denotes $(z_1, \dots, z_{i-1}, x_i, z_{i+1}, \dots, z_s)$. The block approximation error is defined as $e_i(x_i, z) := u_i(x_i, z) - f(x_i, z_{\neq i})$.

A separability condition is necessary in [13, Def. 3] for the surrogate function of f (i.e., when fixing z , the surrogate function u of f satisfies $\hat{u}(x) = \sum_{i=1}^s \hat{u}_i(x_i)$, where $\hat{u}(x) = u(x, z)$ and $\hat{u}_i(x_i) = u_i(x_i, z)$) while our upcoming analysis does not require such a condition.

The inertial alternating direction method of multipliers (iADMM) framework is described in Algorithm 1. iADMM cyclically update the blocks x_1, \dots, x_s and y . We use $x^{k,i}$ to denote $(x_1^{k+1}, \dots, x_i^{k+1}, x_{i+1}^k, \dots, x_s^k)$, let $x^{k,0} = x^k$ and $x^{k+1} = x^{k,s}$, where k is the outer iteration index, and i the cyclic inner iteration index ($i \in [s]$). The update of block x_i in (8) (note that $x_i^{k+1} = x_i^{k,i}$) means that iADMM chooses a surrogate function for $x_i \mapsto \mathcal{L}(x_i, x_{\neq i}^{k,i}, y^k, \omega^k)$, which is formed by summing a surrogate function of $x_i \mapsto f(x_i, x_{\neq i}^{k,i}) + g_i(x_i)$ and a surrogate function of $x_i \mapsto \varphi^k(x_i, x_{\neq i}^{k,i})$ where φ^k is defined in (5), then apply extrapolation to the latter surrogate function³. To update block y , as $h(y)$ is L_h -smooth, we apply Nesterov

³ It is important noting that it is possible to embed the general inertial term \mathcal{G}_i^k to the surrogate of $x_i \mapsto \mathcal{L}(x_i, x_{\neq i}^{k,i}, y^k, \omega^k)$ as in [21]. This inertial term may also lead to the extrapolation for the block surrogate function of $f(x)$ or for both the two block surrogates. However, to simplify our analysis, we only consider here the effect of the inertial term for the block surrogate of $\varphi^k(x)$.

Algorithm 1: iADMM, a general framework for solving Problem (1)

Choose $x^0 = x^{-1}$, $y^0 = y^{-1}$, ω^0 . Let u_i , $i \in [s]$, be a block i surrogate functions of $f(x)$ on \mathbb{R}^n .

For the choice of the extrapolation parameters, ζ_i^k and δ_k , and of the parameters κ_i , α and β , see the paragraph ‘‘Choosing parameters for iADMM’’ (page 8).

for $k = 0, \dots$ **do**

for $i = 1, \dots, s$ **do**

 Compute $\bar{x}_i^k = x_i^k + \zeta_i^k(x_i^k - x_i^{k-1})$, and update block x_i as follows

$$x_i^{k+1} \in \underset{x_i}{\operatorname{argmin}} \left\{ u_i(x_i, x^{k,i-1}) + g_i(x_i) + \langle \mathcal{A}_i^*(\omega^k + \beta(\mathcal{A}\bar{x}^{k,i-1} + \mathcal{B}y^k - b)), x_i \rangle + \frac{\kappa_i \beta}{2} \|x_i - \bar{x}_i^k\|^2 \right\}, \quad (8)$$

 where $\kappa_i \geq \|\mathcal{A}_i^* \mathcal{A}_i\|$, and $\bar{x}^{k,i-1} = (x_1^{k+1}, \dots, x_{i-1}^{k+1}, \bar{x}_i^k, x_{i+1}^k, \dots, x_s^k)$.

end for

 Compute $\hat{y}^k = y^k + \delta_k(y^k - y^{k-1})$, and update y as follows

$$y^{k+1} \in \underset{y}{\operatorname{argmin}} \left\{ \langle \mathcal{B}^* \omega^k + \nabla h(\hat{y}^k), y \rangle + \frac{\beta}{2} \|\mathcal{A}x^{k+1} + \mathcal{B}y - b\|^2 + \frac{Lh}{2} \|y - \hat{y}^k\|^2 \right\}. \quad (9)$$

 Update ω as follows

$$\omega^{k+1} = \omega^k + \alpha \beta (\mathcal{A}x^{k+1} + \mathcal{B}y^{k+1} - b). \quad (10)$$

end for

type acceleration on h as in (9). Together with Assumption 1, we make the following standard assumption for u_i throughout the paper.

Assumption 2 (i) *The block surrogate function $u_i(x_i, z)$ is continuous.*

(ii) *Given $z \in \mathbb{R}^n$, for $i \in [s]$, there exists a function $x_i \mapsto \bar{e}_i(x_i, z)$ such that $\bar{e}_i(\cdot, z)$ is continuously differentiable at z_i , $\bar{e}_i(z_i, z) = 0$, $\nabla_{x_i} \bar{e}_i(z_i, z) = 0$, and the block approximation error e_i satisfies*

$$e_i(x_i, z) \leq \bar{e}_i(x_i, z) \text{ for all } x_i. \quad (11)$$

Assumption 2 (ii) is satisfied when we simply choose $u_i(x_i, z) = f(x_i, z_{\neq i})$ (i.e., $f(x_i, z_{\neq i})$ is a surrogate function of itself), or when $e_i(\cdot, z)$ is continuously differentiable at z_i and $\nabla_{x_i} e_i(z_i, z) = 0$, or when $e_i(x_i, z) \leq c \|x_i - z_i\|^{1+\epsilon}$ for some $\epsilon > 0$ and $c > 0$; see [21, Lemma 3]. In the following, we provide some examples of block surrogate functions satisfying Assumption 2.

– The block proximal surrogate function, see, e.g., [1, 3, 20], has the following form

$$u_i(x_i, z) = f(x_i, z_{\neq i}) + \frac{\rho_i}{2} \|x_i - z_i\|^2,$$

where $\rho_i > 0$ is a scalar. We have $e_i(x_i, z) = \frac{\rho_i}{2} \|x_i - z_i\|^2$. In this case, $\bar{e}_i = e_i$.

– The Lipschitz gradient surrogate function, see, e.g., [62, 63, 20], has the form

$$u_i(x_i, z) = f(z) + \langle \nabla_i f(z), x_i - z_i \rangle + \frac{\kappa_i L_i^{(z)}}{2} \|x_i - z_i\|^2, \quad (12)$$

where $\kappa_i \geq 1$ and we assume $x_i \mapsto f(x_i, z_{\neq i})$ is differentiable and $\nabla_i f(x_i, z_{\neq i})$ is $L_i^{(z)}$ -Lipschitz continuous (we note that $L_i^{(z)}$ may depend on z). We have

$$\nabla_{x_i} e_i(x_i, z) = L_i^{(z)}(x_i - z_i) + \nabla_i f(z) - \nabla_i f(x_i, z_{\neq i}).$$

Hence $\nabla_{x_i} e_i(z_i, z) = 0$. In this case $\bar{e}_i = e_i$.

– The quadratic surrogate, see e.g., [14, 41], has the following form

$$u_i(x_i, z) = f(z) + \langle \nabla_i f(z), x_i - z_i \rangle + \frac{\kappa_i}{2} (x_i - z_i)^T H_i^{(z)} (x_i - z_i), \quad (13)$$

where $\kappa_i \geq 1$ and we assume f is twice differentiable, $H_i^{(z)}$ is a positive definite matrix such that $(H_i^{(z)} - \nabla_i^2 f(x_i, z_{\neq i}))$ is positive definite (we note that $H_i^{(z)}$ may depend on z). Similarly, we also have $\bar{e}_i = e_i$ in this case.

Choosing parameters for iADMM. The parameters of iADMM include: α in (10), κ_i and the extrapolation parameters ζ_i^k in the update (8) of block x_i , the extrapolation parameter δ_k in the update (9) of y , and the penalty parameter β . In the next section, Proposition 1 provides the formulas for η_i and γ_i^k that involve β , κ_i , and ζ_i^k , while Proposition 2 provides the formulas for η_y and γ_y^k that involve β and δ_k . To guarantee a subsequential convergence, we choose $\alpha \in (0, 2)$, and the parameters η_y , γ_y^k , η_i and γ_i^k satisfying the conditions of Proposition 4; see Theorem 1. To guarantee a global convergence, we choose $\alpha = 1$, use no extrapolation for y , and choose the other parameters to satisfy (21); see Theorem 2. It is important noting that the convexity of $x_i \mapsto u_i(x_i, z) + g_i(x_i)$ allows larger extrapolation parameters in the update of x_i (Proposition 1), while the convexity of h allows larger extrapolation parameters in the update of y (Proposition 2).

Remark 1 As we target Nesterov-type acceleration to update of y (h is assumed to be L_h -smooth), we analyse the update rule as in (9) for y . Updating y using $y^{k+1} \in \operatorname{argmin}_y \mathcal{L}(x^{k+1}, y, \omega^k)$ would work as well, and the convergence analysis of iADMM would be simplified by using the same rationale to obtain subsequential as well as global convergence. We hence omit this case in our analysis.

2.2 Convergence analysis

Assumptions. Throughout the paper we assume Assumption 1 and Assumption 2 hold, and $\alpha \in (0, 2)$.

Let $x^{k,i}$, y^k and ω^k be the iterates generated by iADMM. We define some additional notations as follows. We denote $\Delta x_i^k = x_i^k - x_i^{k-1}$, $\Delta y^k = y^k - y^{k-1}$, $\Delta \omega^k = \omega^k - \omega^{k-1}$, $\alpha_1 = \frac{|1-\alpha|}{\alpha \sigma_B (1-|1-\alpha|)}$, $\alpha_2 = \frac{3\alpha}{\sigma_B (1-|1-\alpha|)^2}$ and $\mathcal{L}^k = \mathcal{L}(x^k, y^k, \omega^k)$. We let ν_i , $i \in [s]$, and ν_y be arbitrary constants in $(0, 1)$. We take the following convention in the notation that allows us to analyse iADMM and its non-inertial version in parallel:

- If $\zeta_i^k = 0$ (i.e., there is no extrapolation in the update of x_i^k), then $\zeta_i^k / \nu_i = 0$ and $\nu_i = 0$.
- If $\delta_k = 0$ (i.e., there is no extrapolation in the update of y), then $\delta_k / \nu_y = 0$ and $\nu_y = 0$.

Now we present our main convergence results; see the proofs in Appendix B.

As iADMM allows to use extrapolation in the update of x_i^k and y^k , the Lagrangian is not guaranteed to decrease at each iteration. Instead, it has the following nearly sufficiently decreasing property as stated in the following Propositions 1 and 2.

Proposition 1 (i) Considering the update in (8), in general (when $x_i \mapsto u_i(x_i, z) + g_i(x_i)$ can be nonconvex), we choose $\kappa_i > \|A_i^* A_i\|$. Denote $a_i^k = \beta \zeta_i^k (\kappa_i + \|A_i^* A_i\|)$. Then we have

$$\mathcal{L}(x^{k,i}, y^k, \omega^k) + \eta_i \|\Delta x_i^{k+1}\|^2 \leq \mathcal{L}(x^{k,i-1}, y^k, \omega^k) + \gamma_i^k \|\Delta x_i^k\|^2, \quad (14)$$

where

$$\eta_i = \frac{(1 - \nu_i)(\kappa_i - \|A_i^* A_i\|)\beta}{2}, \quad \gamma_i^k = \frac{(a_i^k)^2}{2\nu_i(\kappa_i - \|A_i^* A_i\|)\beta}. \quad (15)$$

(ii) If $x_i \mapsto u_i(x_i, z) + g_i(x_i)$ is convex, we take $\kappa_i = \|A_i^* A_i\|$ (note that if $\|A_i^* A_i\| = 0$ then we can choose κ_i as in case (i)). Inequality (14) is then satisfied with

$$\gamma_i^k = \frac{\beta \|A_i^* A_i\| (\zeta_i^k)^2}{2}, \quad \eta_i = \frac{\beta \|A_i^* A_i\|}{2}. \quad (16)$$

Proposition 2 Considering the update in (9), we have

$$\mathcal{L}(x^{k+1}, y^{k+1}, \omega^k) + \eta_y \|\Delta y^{k+1}\|^2 \leq \mathcal{L}(x^{k+1}, y^k, \omega^k) + \gamma_y^k \|\Delta y^k\|^2,$$

where $\eta_y = \frac{(1 - \nu_y)(\beta \lambda_{\min}(\mathcal{B}^* \mathcal{B}) + L_h)}{2}$ and $\gamma_y^k = \frac{2L_h^2 \delta_k^2}{\nu_y(\beta \lambda_{\min}(\mathcal{B}^* \mathcal{B}) + L_h)}$ when $h(y)$ is nonconvex, and $\eta_y = \frac{L_h}{2}$ and $\gamma_y^k = \frac{L_h \delta_k^2}{2}$ when $h(y)$ is convex.

From Proposition 1 and Proposition 2, we obtain the following recursion for $\{\mathcal{L}^k\}$.

Proposition 3 We have

$$\begin{aligned} & \mathcal{L}^{k+1} + \eta_y \|\Delta y^{k+1}\|^2 + \sum_{i=1}^s \eta_i \|\Delta x_i^{k+1}\|^2 \\ & \leq \mathcal{L}^k + \sum_{i=1}^s \gamma_i^k \|\Delta x_i^k\|^2 + \gamma_y^k \|\Delta y^k\|^2 + \frac{\alpha_1}{\beta} (\|B^* \Delta \omega^k\|^2 - \|B^* \Delta \omega^{k+1}\|^2) \\ & \quad + \frac{\alpha_2}{\beta} L_h^2 \|\Delta y^{k+1}\|^2 + \frac{\alpha_2}{\beta} (\bar{\delta}_k L_h^2 \|\Delta y^k\|^2 + 4L_h^2 \delta_{k-1}^2 \|\Delta y^{k-1}\|^2), \end{aligned} \quad (17)$$

where $\bar{\delta}_k = 2$ if $\delta_k = 0$ for all k and $4(1 + \delta_k)^2$ otherwise.

Now we characterize the chosen parameters for Algorithm 1 in the following proposition.

Proposition 4 Let $\eta_y, \gamma_y^k, \eta_i,$ and $\gamma_i^k, i \in [s]$, be defined in Proposition 1 and Proposition 2. Denote $\mu = \eta_y - \frac{\alpha_2 L_h^2}{\beta}$. For $k \geq 1$, suppose the parameters are chosen such that $\mu > 0, \eta_i > 0$, and the following conditions are satisfied for some constants $0 < C_x, C_y < 1$:

$$\gamma_i^k \leq C_x \eta_i, \quad \frac{4\alpha_2 L_h^2 \delta_{k-1}^2}{\beta} \leq C_2 \mu, \quad \frac{\alpha_2 L_h^2 \bar{\delta}_k}{\beta} + \gamma_y^k \leq C_1 \mu, \quad (18)$$

where $\begin{cases} C_1 = C_y \text{ and } C_2 = 0 & \text{if } \delta_k = 0 \forall k, \\ 0 < C_1 < C_y \text{ and } C_2 = C_y - C_1 & \text{otherwise,} \end{cases}$ and $\bar{\delta}_k$ is defined in Proposition 3. Furthermore, suppose we use one of the following methods:

- we choose $\delta_k = 0$ for all k , that is, there is no extrapolation in the update of y ,
- we use extrapolation in the update of y and choose the parameters such that

$$\beta \geq \frac{4L_h\alpha}{\sigma_{\mathcal{B}}(1-|\alpha|)}, \quad \beta \geq \frac{6\alpha L_h^2}{\mu\sigma_{\mathcal{B}}(1-|\alpha|)} \max\left\{1, \frac{12\delta_k^2}{1-C_1}\right\}. \quad (19)$$

(i) For $K > 1$ we have

$$\begin{aligned} & \mathcal{L}^{K+1} + \mu\|\Delta y^{K+1}\|^2 + \sum_{i=1}^s \eta_i \|\Delta x_i^{K+1}\|^2 + \frac{\alpha_1}{\beta} \|\mathcal{B}^* \Delta w^{K+1}\|^2 + (1-C_1)\mu\|\Delta y^K\|^2 \\ & + \sum_{k=1}^{K-1} [(1-C_y)\mu\|\Delta y^k\|^2 + (1-C_x) \sum_{i=1}^s \eta_i \|\Delta x_i^{k+1}\|^2] \\ & \leq \mathcal{L}^1 + \frac{\alpha_1}{\beta} \|\mathcal{B}^* \Delta w^1\|^2 + C_x \sum_{i=1}^s \eta_i \|\Delta x_i^1\|^2 + \mu\|\Delta y^1\|^2 + C_2\mu\|\Delta y^0\|^2. \end{aligned} \quad (20)$$

(ii) The sequences $\{\Delta y^k\}$, $\{\Delta x_i^k\}$ and $\{\Delta w^k\}$ converge to 0.

We will assume that Algorithm 1 generates a bounded sequence in our subsequential and global convergence results. Let us provide a sufficient condition that guarantees this boundedness assumption.

Proposition 5 *If $b + \text{Im}(\mathcal{A}) \subseteq \text{Im}(\mathcal{B})$, $\lambda_{\min}(\mathcal{B}^*\mathcal{B}) > 0$ and $F(x) + h(y)$ is coercive over the feasible set $\{(x, y) : \mathcal{A}x + \mathcal{B}y = b\}$ then the sequences $\{x^k\}$, $\{y^k\}$ and $\{\omega^k\}$ generated by Algorithm 1 are bounded.*

It is important noting that the coercive condition of $F(x) + h(y)$ over the feasible set is weaker than the coercive condition of $F(x) + h(y)$ over $x \in \mathbb{R}^n, y \in \mathbb{R}^q$. Let us now present the subsequential convergence, the global convergence of the generated sequence and its convergence rate.

Theorem 1 (Subsequential convergence) *Suppose the parameters of Algorithm 1 are chosen satisfying the conditions in Proposition 4. If the generated sequence of Algorithm 1 is bounded, then every limit point of the generated sequence is a critical point of \mathcal{L} .*

Theorem 2 (Global convergence) *Suppose we do not use extrapolation to update y , that is, $\delta_k = 0$ for all k (note that extrapolation to update x_i is still applicable), and we take $\alpha = 1$. Then the conditions in (18) become*

$$\gamma_i^k \leq C_x \eta_i, \quad \frac{2\alpha_2 L_h^2}{\beta} \leq C_y \mu, \quad \text{for all } k \geq 0, i \in [s] \quad (21)$$

for some constants $0 < C_x, C_y < 1$. Furthermore, we assume that (i) for any $x, z \in \mathbb{R}^n$, $x_i \in \text{dom}(g_i)$, we have

$$\begin{aligned} \partial_{x_i}(f(x) + g_i(x_i)) &= \partial_{x_i} f(x) + \partial_{x_i} g_i(x_i), \\ \partial_{x_i}(u_i(x_i, z) + g_i(x_i)) &= \partial_{x_i} u_i(x_i, z) + \partial_{x_i} g_i(x_i), \end{aligned} \quad (22)$$

and (ii) for any x, z in a bounded subset of \mathbb{R}^n , if $s_i \in \partial_{x_i} u_i(x_i, z)$, there exists $\xi_i \in \partial_{x_i} f(x)$ such that

$$\|\xi_i - \mathbf{s}_i\| \leq L_i \|x - z\| \text{ for some constant } L_i. \quad (23)$$

If the generated sequence of Algorithm 1 is bounded and $F(x) + h(y)$ has the KL property (see Appendix A), then the whole generated sequence converges to a critical point of \mathcal{L} .

We refer the readers to [49, Corollary 10.9] for a sufficient condition for (22) (see Appendix A for more details). Some specific examples that satisfy (22) include: (i) $g_i = 0$, (ii) the functions $x_i \mapsto f(x)$ and $x_i \mapsto u_i(x_i, z)$ are strictly differentiable (see [49, Exercise 10.10]), (iii) the functions $x_i \mapsto f(x)$ and $x_i \mapsto u_i(x_i, z)$ are convex and the relative interior qualification conditions are satisfied: $\text{ri}(\text{dom}(f(\cdot, x_{\neq i}))) \cap \text{ri}(\text{dom}g_i) \neq \emptyset$ and $\text{ri}(\text{dom}(g(\cdot, z))) \cap \text{ri}(\text{dom}g_i) \neq \emptyset$, where ri is short for relative interior. We note that although the condition in (23) is necessary for our convergence proof, the constant L_i does not influence how to choose the parameters in our framework. The condition in (23) is satisfied when both u_i and f are twice continuously differentiable and $\nabla_{x_i} e_i(x_i, x) = 0$ for all x (which implies that $\nabla_{x_i} u_i(x_i, x) = \nabla_{x_i} f(x)$ for all x). Indeed, in this case we have

$$\|\nabla_{x_i} u_i(x_i, z) - \nabla_{x_i} f(x)\| = \|\nabla_{x_i} u_i(x_i, z) - \nabla_{x_i} u_i(x_i, x)\| \leq L_i \|x - z\|$$

for some L_i because $\nabla_{x_i} u_i(x_i, z)$ is continuously differentiable and thus is Lipschitz continuous over any bounded subset. We note that all the examples given after Assumption 2 in Section 2 satisfy the condition in (23) when f is twice continuously differentiable.

Convergence rate A convergence rate for the generated sequence of iADMM can be derived using the same technique as in [1, Theorem 2]. To the best of our knowledge, in the nonconvex setting, the convergence rate for block coordinate methods (including inertial as well as non-inertial algorithms) appears to be the same in different papers in the literature since all papers use the technique in [1]. As it is similar to establish the rate for iADMM, we omit the details. Instead, we refer the readers to [62, Theorem 2.9] and [20, Theorem 3] for some examples of using this technique to establish the convergence rate. The type of the convergence rate depends on the value of the KL exponent, which is the coefficient \mathbf{a} such that $\mathcal{Y}(t)$ in Definition 6 (see Appendix A) equals $ct^{1-\mathbf{a}}$, where c is a constant. Specifically, when $\mathbf{a} = 0$, the algorithm converges after a finite number of steps, when $\mathbf{a} \in (0, 1/2]$, the algorithm has linear convergence, and when $\mathbf{a} \in (1/2, 1)$, the algorithm has sublinear convergence. Determining the value of the KL exponent is out of the scope of this paper, and is an active and challenging topic.

3 Numerical results

In this section, we apply iADMM to solve a latent low-rank representation problem. We consider Problem (2) with

$$- r_1(t) = \lambda_1 t \text{ to promote } X \text{ to be of low-rank, since } r_1(\sigma(X)) = \lambda_1 \sum_i \sigma_i(X) = \lambda_1 \|X\|_* \text{ is the nuclear norm [47].}$$

- $r_2(Y) = \lambda \sum_{i=1}^q \phi(\|Y_i\|_2)$, where $\phi(t) = 1 - \exp(-\theta t)$ is concave, $\theta > 0$ is a parameter, and Y_i is the i -th column of Y . This is a *nonconvex* regularization that promotes Y to be column sparse, that is, it promotes Y to have many columns equal to the zero vector [10]. In fact, $\phi(t) = 0$ when $t = 0$, while $\phi(t)$ quickly goes to 1 as t increases.
- $r_3(Z) = \frac{1}{2}\|Z\|^2$ to model Gaussian noise.
- $A_1 = DP_1$ and $A_2 = P_2^*D$, where P_1 and P_2 are computed by orthogonalizing the columns of D^* and D , respectively, as proposed in [34]. In [34], the authors showed that this resulting problem is a simpler equivalent form of the one in which D is considered as a dictionary, i.e. $A_1 = A_2 = D$. Hence, it can be scaled for data sets with a large number of observations.

In this scenario, Problem (2) takes the form of (1) with B being the identity operator, b being the data set D , x_1 and x_2 being the matrices X and Y , y being the matrix Z , $g_i = 0$, $h(Z) = \frac{1}{2}\|Z\|^2$ and $f(X, Y) = \lambda_1\|X\|_* + r_2(Y)$.

3.1 Surrogate functions and iADMM updates.

We choose $u_1(X, X^k, Y^k) = \lambda_1\|X\|_* + r_2(Y^k)$, and $u_2(Y, X^{k+1}, Y^k) = r_2(Y^k) + \sum_{i=1}^q \zeta_i^k (\|Y_i\|_2 - \|Y_i^k\|_2) + \lambda_1\|X^{k+1}\|_*$, where $\zeta_i^k = \lambda \nabla \phi(\|Y_i^k\|_2)$. The function u_1 satisfies Assumption 2, and u_2 satisfies Assumption 2 (i). Since ϕ is continuously differentiable with Lipschitz gradient on $[0, +\infty)$, and the Euclidean norm is Lipschitz continuous, it follows from Section 4.5 of [21] that u_2 also satisfies Assumption 2 (ii). We derive from [48, Corollary 5Q] that the condition in (23) is satisfied. According to the update (8), X^{k+1} is computed by solving the following nuclear norm problem

$$\min_X \lambda_1\|X\|_* + \left\langle A_1^* \left(\beta(A_1 \bar{X}^k + Y^k A_2 + Z^k - D) + W^k \right), X \right\rangle + \frac{\kappa_1 \beta}{2} \|X - \bar{X}^k\|^2, \quad (24)$$

where $\kappa_1 \geq \|A_1^* A_1\|$ and $\bar{X}^k = X^k + \zeta_1^k (X^k - X^{k-1})$. Let $\text{diag}(u)$ denote a diagonal matrix whose diagonal elements are the entries of u , and $[\cdot]_+$ denote the projection onto the nonnegative orthant. Problem (24) has a closed-form solution given by $X^{k+1} = US_{\lambda_1/(\kappa_1 \beta)} V^T$, where USV^T is the SVD of $\bar{X}^k - A_1^* (A_1 \bar{X}^k + Y^k A_2 + Z^k - D + W^k)/(\kappa_1 \beta)$ and $S_{\lambda_1/(\kappa_1 \beta)} = \text{diag}([S_{ii} - \lambda_1/(\kappa_1 \beta)]_+)$. Letting $\kappa_2 \geq \|A_2 A_2^*\|$ and $\bar{Y}^k = Y^k + \zeta_2^k (Y^k - Y^{k-1})$, the update (9) for Y is

$$Y^{k+1} \in \arg \min_Y \sum_{i=1}^q \zeta_i^k \|Y_i\|_2 + \langle (W^k + \beta(A_1 X^{k+1} + \bar{Y}^k A_2 + Z^k - D)) A_2^*, Y \rangle + \frac{\kappa_2 \beta}{2} \|Y - \bar{Y}^k\|^2.$$

It has a closed-form solution given by $Y_i^{k+1} = \left[\|P_i^k\| - \zeta_i^k / (\kappa_2 \beta) \right]_+ \frac{P_i^k}{\|P_i^k\|}$, where P_i^k is the i -th column of $\bar{Y}^k - (A_1 X^{k+1} + \bar{Y}^k A_2 + Z^k - D)/\kappa_2 - W^k / (\kappa_2 \beta)$. The update (9) for Z is $Z^{k+1} = -(W^k + \beta(A_1 X^{k+1} + Y^{k+1} A_2 - D))/(1 + \beta)$, and the update (10) for W is $W^{k+1} = W^k + \alpha \beta (A_1 X^{k+1} + Y^{k+1} A_2 + Z^{k+1} - D)$.

3.2 Choosing parameters

We have $L_h = 1$, $\sigma_B = 1$, and $\delta_k = 0$. As $h(Z)$ is convex and we do not apply extrapolation for Z , by Proposition 2, $\eta_y = \frac{1}{2}$ and $\gamma_y^k = 0$. Since $\|X\|_*$ and $\sum_{i=1}^q \zeta_i^k \|Y_i\|_2$ are convex, we choose $\kappa_1 = \|A_1^* A_1\|$, $\kappa_2 = \|A_2 A_2^*\|$, and the conditions in (21) become $\zeta_i^k \leq \sqrt{C_x}$ ($i = 1, 2$) and $\frac{(2+C_y)\alpha_2}{\beta} \leq \frac{C_y}{2}$. We take $C_x = 1 - 10^{-15}$, $\alpha = 1$, $C_y = 1 - 10^{-6}$, $\beta = 2(2 + C_y)\alpha_2/C_y$, $a_0 = 1$, $a_k = \frac{1}{2}(1 + \sqrt{1 + 4a_{k-1}^2})$, and $\zeta_i^k = \min\left\{\frac{a_k - 1 - 1}{a_k}, \sqrt{C_x}\right\}$. We set $\alpha = 1$ as we target global convergence. We have also conducted experiments with other values of α (namely 0.5, 1.4 and 1.8); see Appendix C.

3.3 Experiments

We compare the following three methods: (1) ADMM-mm: iADMM without extrapolation, (2) iADMM-mm: iADMM with extrapolation, (3) linearizedADMM: a linearized ADMM which is different from ADMM-mm for updating Y . linearizedADMM updates Y by solving $\min -\lambda \exp(\|Y_i\|_2) + \frac{\kappa_2 \beta}{2} \|Y_i - V_i^k\|^2$, where V_i^k is the i -th column of $X^k - (W^k + \beta(A_1 X^{k+1} + \bar{Y}^k A_2 + Z^k - D))A_2^*/(\kappa_2 \beta)$. Since these sub-problems do not have closed-form solutions, we employ an MM scheme to solve them. To examine the performance of the three algorithms, we consider subspace segmentation tasks. After obtaining a solution X^* , we follow the setting in [33] to construct the affinity matrix Q by $Q_{ij} = (\tilde{U}\tilde{U}^T)_{ij}$, where \tilde{U} is formed by $U^*(\Sigma^*)^{1/2}$ with normalized rows and $U^* \Sigma^* (V^*)^T$ being the SVD of X^* . Finally, we apply the Normalized Cuts [26] on Q to cluster the data into groups. The experiments are run on three data sets: Hopkins 155, extended Yale B and Umist. Hopkins 155 consists of 156 sequences, each of which has from 39 to 550 vectors drawn from two or three motions (one motion corresponds to one subspace). Each sequence is a sole segmentation task and thus there are 156 clustering tasks in total. Yale B contains 2414 frontal face images with 38 classes, and Umist contains 564 images with 20 classes. To avoid computational issues when computing the segmentation error rate, we construct clustering tasks by using the first 10 classes of these two data sets [36]. All tests are preformed using Matlab R2019a on a PC 2.3 GHz Intel Core i5 of 8GB RAM.

In our experiments, we choose $\theta = 5$, $\lambda_1 = \lambda = 0.01$ for Hopkins 155, and $\lambda_1 = \lambda = 1$ for the two other data sets. We set the initial points to zero, that is, $X^0 = 0$, $Y^0 = 0$, $Z^0 = 0$, $W^0 = 0$. We do not optimize numerical results by tweaking the parameters and initial points as this is beyond the scope of this work. It is important noting that we evaluate the algorithms on the same models with the same initializations. We run each algorithm 10, 300, and 500 seconds for each sequence of Hopkins 155, Umist10, and Yaleb10, respectively. Figure 1 displays the values of the segmentation error rate and the objective function versus the training time, and Table 1 reports the final values. Since there are 156 sequences (data sets) in Hopkins 155, we plot the average values, and report the final average results and standard deviation over these sequences. We observe that iADMM-mm converges the fastest on all the data sets, providing a significant acceleration of ADMM-mm. iADMM-mm achieves not only the best final objective function values

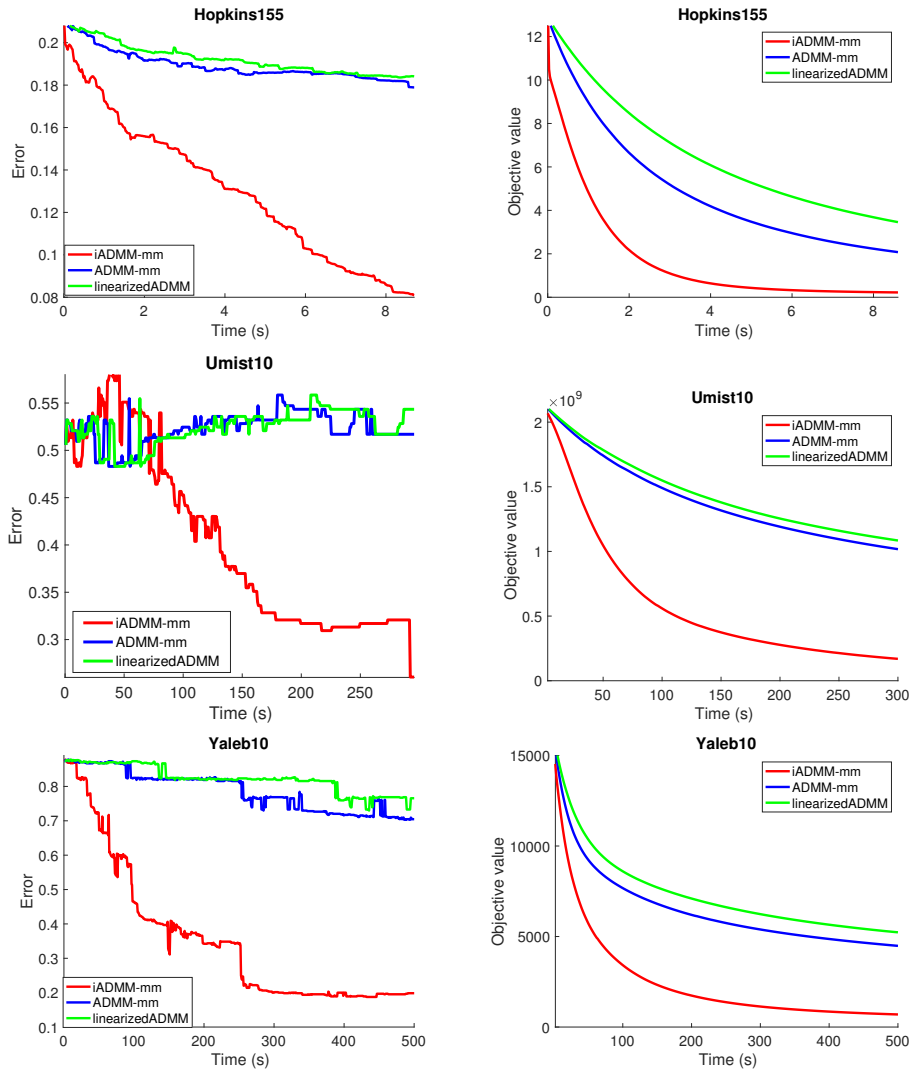


Fig. 1 Evolution of the segmentation error rate and the objective function value with respect to time. For Hopkins155, the results are the average values over 156 sequences.

but also the best segmentation error rates. This illustrates the usefulness of the acceleration technique. In addition, ADMM-mm outperforms linearizedADMM which illustrates the usefulness of properly choosing a proper surrogate function. The conclusions are the same for other values of α ; see Appendix C.

4 Conclusion

We have proposed and analysed iADMM, a framework of inertial alternating direction methods of multipliers, for solving a class of nonconvex nonsmooth

Table 1 Comparison of segmentation error rate and final objective function values obtained within the allotted time. Bold values indicate the best results.

	Method	Error mean \pm std	Obj. value mean \pm std
Hopkins	linearizedADMM	0.1579 \pm 0.1550	3.0254 \pm 2.4189
	ADMM-mm	0.1472 \pm 0.1513	1.8081 \pm 1.6674
	iADMM-mm	0.0562 \pm 0.1006	0.2023 \pm 0.1062
Urnist	linearizedADMM	0.5170	1.0838 $\times 10^9$
	ADMM-mm	0.5170	1.0167 $\times 10^9$
	iADMM-mm	0.2604	0.1694 $\times 10^9$
Yaleb	linearizedADMM	0.7656	5.2317 $\times 10^3$
	ADMM-mm	0.7047	4.4829 $\times 10^3$
	iADMM-mm	0.1984	0.6951 $\times 10^3$

optimization problem with linear constraints. The preliminary computational results in solving a class of nonconvex low-rank representation problems not only show the efficacy of using inertial terms for ADMM but also show the advantage of using suitable block surrogate functions that provide closed-form solutions in the block update of ADMM. We conclude the paper by mentioning two important questions that we consider as a future research directions: (i) Can we extend the cyclic update rule of iADMM to randomized/non-cyclic setting? (ii) To guarantee the global convergence, iADMM does not allow extrapolation in the update of y ; see Theorem 2. Can we extend the analysis to allow the extrapolation in the update of y ?

Declarations

Availability of data and material, and Code availability The data and code are available from <https://github.com/nhatpd/iADMM>.

Funding LTKH and NG acknowledge the support by the European Research Council (ERC starting grant no 679515), and by the Fonds de la Recherche Scientifique - FNRS and the Fonds Wetenschappelijk Onderzoek - Vlaanderen (FWO) under EOS Project no O005318F-RG47. NG also acknowledges the Francqui Foundation.

Conflicts of interest/Competing interests Not applicable.

References

1. Attouch, H., Bolte, J.: On the convergence of the proximal algorithm for nonsmooth functions involving analytic features. *Mathematical Programming* **116**(1), 5–16 (2009)
2. Attouch, H., Bolte, J., Redont, P., Soubeyran, A.: Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Łojasiewicz inequality. *Mathematics of Operations Research* **35**(2), 438–457 (2010)
3. Attouch, H., Bolte, J., Svaiter, B.F.: Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized gauss-seidel methods. *Mathematical Programming* **137**(1), 91–129 (2013)
4. Bach, F., Jenatton, R., Mairal, J., Obozinski, G.: Optimization with sparsity-inducing penalties. *Foundations and Trends in Machine Learning* **4** (2011). DOI 10.1561/22000000015

5. Beck, A., Tetrushvili, L.: On the convergence of block coordinate descent type methods. *SIAM Journal on Optimization* **23**, 2037–2060 (2013)
6. Bochnak, J., Coste, M., Roy, M.F.: *Real Algebraic Geometry*. Springer (1998)
7. Bolte, J., Sabach, S., Teboulle, M.: Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming* **146**(1), 459–494 (2014)
8. Bot, R.I., Nguyen, D.K.: The proximal alternating direction method of multipliers in the nonconvex setting: Convergence analysis and rates. *Mathematics of Operations Research* **45**(2), 682–712 (2020)
9. Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J.: Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.* **3**(1), 1–122 (2011)
10. Bradley, P.S., Mangasarian, O.L.: Feature selection via concave minimization and support vector machines. In: *Proceeding of international conference on machine learning ICML’98* (1998)
11. Buccini, A., Dell’Acqua, P., Donatelli, M.: A general framework for admm acceleration. *Numerical Algorithms* **85** (2020). DOI 10.1007/s11075-019-00839-y
12. Candès, E.J., Li, X., Ma, Y., Wright, J.: Robust principal component analysis? *J. ACM* **58**(3) (2011)
13. Canyi, L., Feng, J., Yan, S., Lin, Z.: A unified alternating direction method of multipliers by majorization minimization. *IEEE transactions on pattern analysis and machine intelligence* **40**, 527 – 541 (2018). DOI 10.1109/TPAMI.2017.2689021
14. Chouzenoux, E., Pesquet, J.C., Repetti, A.: A block coordinate variable metric forward–backward algorithm. *Journal of Global Optimization* **66**, 457–485 (2016)
15. Deng, W., Yin, W.: On the global and linear convergence of the generalized alternating direction method of multipliers. *Rice CAAM tech report TR12-14* **66** (2012)
16. Fazel, M., Pong, T.K., Sun, D., Tseng, P.: Hankel matrix rank minimization with applications to system identification and realization. *SIAM Journal on Matrix Analysis and Applications* **34**(3), 946–977 (2013)
17. Gabay, D., Mercier, B.: A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers & Mathematics with Applications* **2**(1), 17 – 40 (1976)
18. Glowinski, R., Marroco, A.: Sur l’approximation, par éléments finis d’ordre un, et la résolution, par pénalisation-dualité d’une classe de problèmes de dirichlet non linéaires. *ESAIM: Mathematical Modelling and Numerical Analysis - Modélisation Mathématique et Analyse Numérique* **9**(R2), 41–76 (1975)
19. Grippo, L., Sciandrone, M.: On the convergence of the block nonlinear gauss–seidel method under convex constraints. *Operations Research Letters* **26**(3), 127 – 136 (2000)
20. Hien, L.T.K., Gillis, N., Patrinos, P.: Inertial block proximal method for non-convex nonsmooth optimization. In: *Thirty-seventh International Conference on Machine Learning ICML 2020* (2020)
21. Hien, L.T.K., Phan, D.N., Gillis, N.: Inertial block majorization minimization framework for nonconvex nonsmooth optimization (2020). ArXiv:2010.12133
22. Hildreth, C.: A quadratic programming procedure. *Naval Research Logistics Quarterly* **4**(1), 79–85 (1957)
23. Hong, M., Chang, T.H., Wang, X., Razaviyayn, M., Ma, S., Luo, Z.Q.: A block successive upper-bound minimization method of multipliers for linearly constrained convex optimization. *Mathematics of Operations Research* **45**(3), 833–861 (2020)
24. Huang, F., Chen, S., Huang, H.: Faster stochastic alternating direction method of multipliers for nonconvex optimization. In: K. Chaudhuri, R. Salakhutdinov (eds.) *Proceedings of the 36th International Conference on Machine Learning, Proceedings of Machine Learning Research*, vol. 97, pp. 2839–2848. PMLR (2019). URL <http://proceedings.mlr.press/v97/huang19a.html>
25. Huang, F., Chen, S., Lu, Z.: Stochastic alternating direction method of multipliers with variance reduction for nonconvex optimization (2016). ArXiv:1610.02758
26. Jianbo Shi, Malik, J.: Normalized cuts and image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* **22**(8), 888–905 (2000)
27. Koren, Y., Bell, R., Volinsky, C.: Matrix factorization techniques for recommender systems. *Computer* **42**(8), 30–37 (2009)
28. Lai, R., Osher, S.: A splitting method for orthogonality constrained problems. *Journal of Scientific Computing* **58** (2014). DOI 10.1007/s10915-013-9740-x

29. Lee, D.D., Seung, H.S.: Learning the parts of objects by non-negative matrix factorization. *Nature* **401**(6755), 788–791 (1999)
30. Li, G., Pong, T.K.: Global convergence of splitting methods for nonconvex composite optimization. *SIAM Journal on Optimization* **25**(4), 2434–2460 (2015). DOI 10.1137/140998135
31. Li, H., Lin, Z.: Accelerated alternating direction method of multipliers: An optimal $\mathcal{O}(1/k)$ nonergodic analysis. *Journal of Scientific Computing* **79**, 671–699 (2019)
32. Lin, Z., Liu, R., Su, Z.: Linearized alternating direction method with adaptive penalty for low-rank representation. In: J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, K.Q. Weinberger (eds.) *Advances in Neural Information Processing Systems*, vol. 24, pp. 612–620. Curran Associates, Inc. (2011)
33. Liu, G., Lin, Z., Yan, S., Sun, J., Yu, Y., Ma, Y.: Robust recovery of subspace structures by low-rank representation. *IEEE Trans. Pattern Anal. Mach. Intell.* **35**(1), 171–184 (2013)
34. Liu, G., Yan, S.: Latent low-rank representation for subspace segmentation and feature extraction. 2011 International Conference on Computer Vision pp. 1615–1622 (2011)
35. Liu, Q., Shen, X., Gu, Y.: Linearized admm for nonconvex nonsmooth optimization with convergence analysis. *IEEE Access* **7**, 76,131–76,144 (2019)
36. Lu, C., Tang, J., Yan, S., Lin, Z.: Nonconvex nonsmooth low rank minimization via iteratively reweighted nuclear norm. *IEEE Transactions on Image Processing* **25**(2), 829–839 (2016)
37. Mairal, J.: Optimization with first-order surrogate functions. In: *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28, ICML’13*, pp. 783–791. JMLR.org (2013)
38. Markovskiy, I.: *Low rank approximation: algorithms, implementation, applications*, vol. 906. Springer (2012)
39. Melo, J.G., Monteiro, R.D.C.: Iteration-complexity of a jacobi-type non-euclidean admm for multi-block linearly constrained nonconvex programs (2017)
40. Nesterov, Y.: *Introductory lectures on convex optimization: A basic course*. Kluwer Academic Publ. (2004)
41. Ochs, P.: Unifying abstract inexact convergence theorems and block coordinate variable metric ipiano. *SIAM Journal on Optimization* **29**(1), 541–570 (2019)
42. Ouyang, Y., Chen, Y., Lan, G., Pasiliao, E.: An accelerated linearized alternating direction method of multipliers. *SIAM Journal on Imaging Sciences* **8**(1), 644–681 (2015)
43. Parikh, N., Boyd, S.: Proximal algorithms. *Foundations and Trends in Optimization* **1**(3), 127–239 (2014)
44. Pock, T., Sabach, S.: Inertial proximal alternating linearized minimization (iPALM) for nonconvex and nonsmooth problems. *SIAM Journal on Imaging Sciences* **9**(4), 1756–1787 (2016)
45. Powell, M.J.D.: On search directions for minimization algorithms. *Mathematical Programming* **4**(1), 193–201 (1973)
46. Razaviyayn, M., Hong, M., Luo, Z.: A unified convergence analysis of block successive minimization methods for nonsmooth optimization. *SIAM Journal on Optimization* **23**(2), 1126–1153 (2013)
47. Recht, B., Fazel, M., Parrilo, P.A.: Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review* **52**(3), 471–501 (2010)
48. Rockafellar, R.T.: *The Theory Of Subgradients And Its Applications To Problems Of Optimization - Convex And Nonconvex Functions*. Heldermann, Heidelberg, Berlin (1981)
49. Rockafellar, R.T., Wets, R.J.B.: *Variational Analysis*. Springer Verlag, Heidelberg, Berlin, New York (1998)
50. Scheinberg, K., Ma, S., Goldfarb, D.: Sparse inverse covariance selection via alternating linearization methods. In: J.D. Lafferty, C.K.I. Williams, J. Shawe-Taylor, R.S. Zemel, A. Culotta (eds.) *Advances in Neural Information Processing Systems 23*, pp. 2101–2109. Curran Associates, Inc. (2010)
51. Sun, T., Barrio, R., Rodríguez, M., Jiang, H.: Inertial nonconvex alternating minimizations for the image deblurring. *IEEE Transactions on Image Processing* **28**(12), 6211–6224 (2019)
52. Sun, Y., Babu, P., Palomar, D.P.: Majorization-minimization algorithms in signal processing, communications, and machine learning. *IEEE Transactions on Signal Processing* **65**(3), 794–816 (2017). DOI 10.1109/TSP.2016.2601299
53. Tseng, P.: Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications* **109**(3), 475–494 (2001)

54. Tseng, P., Yun, S.: A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming* **117**(1), 387–423 (2009)
55. Udell, M., Horn, C., Zadeh, R., Boyd, S.: Generalized low rank models. *Foundations and Trends in Machine Learning* **9**(1), 1–118 (2016)
56. Udell, M., Townsend, A.: Why are big data matrices approximately low rank? *SIAM Journal on Mathematics of Data Science* **1**(1), 144–160 (2019)
57. Von Luxburg, U.: A tutorial on spectral clustering. *Statistics and computing* **17**(4), 395–416 (2007)
58. Wang, Y., Yin, W., Zeng, J.: Global convergence of admm in nonconvex nonsmooth optimization. *Journal of Scientific Computing* **78**, 29–63 (2019). DOI 10.1007/s10915-018-0757-z
59. Wang, Y., Zeng, J., Peng, Z., Chang, X., Xu, Z.: Linear convergence of adaptively iterative thresholding algorithms for compressed sensing. *IEEE Transactions on Signal Processing* **63**(11), 2957–2971 (2015)
60. Wen, Z., Yin, W.: A feasible method for optimization with orthogonality constraints. *Mathematical Programming* **142** (2010)
61. Xu, M., Wu, T.: A class of linearized proximal alternating direction methods. *J. Optimization Theory and Applications* **151**, 321–337 (2011). DOI 10.1007/s10957-011-9876-5
62. Xu, Y., Yin, W.: A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on Imaging Sciences* **6**(3), 1758–1789 (2013). URL <https://doi.org/10.1137/120887795>
63. Xu, Y., Yin, W.: A globally convergent algorithm for nonconvex optimization based on block coordinate update. *Journal of Scientific Computing* **72**(2), 700–734 (2017)
64. Yang, J., Zhang, Y., Yin, W.: An efficient TVL1 algorithm for deblurring multichannel images corrupted by impulsive noise. *SIAM Journal on Scientific Computing* **31**(4), 2842–2865 (2009)
65. Yang, L., Pong, T.K., Chen, X.: Alternating direction method of multipliers for a class of nonconvex and nonsmooth problems with applications to background/foreground extraction. *SIAM Journal on Imaging Sciences* **10**(1), 74–110 (2017). DOI 10.1137/15M1027528
66. Yin, W., Osher, S., Goldfarb, D., Darbon, J.: Bregman iterative algorithms for $l(1)$ -minimization with applications to compressed sensing. *Siam Journal on Imaging Sciences* **1**, 143–168 (2008)

APPENDIX

A Preliminaries of non-convex non-smooth optimization

In this appendix, we recall some basic definitions and results, namely directional derivative and subdifferentials in Definition 3, critical point in Definition 4, the subdifferential of a sum of function in Proposition 6, and KL functions in Definition 6.

Let $g : \mathbb{E} \rightarrow \mathbb{R} \cup \{+\infty\}$ be a proper lower semicontinuous function.

Definition 3 [49, Definition 8.3]

- (i) For any $x \in \text{dom } g$, and $d \in \mathbb{E}$, we denote the directional derivative of g at x in the direction d by

$$g'(x; d) = \liminf_{\tau \downarrow 0} \frac{g(x + \tau d) - g(x)}{\tau}.$$

- (ii) For each $x \in \text{dom } g$, we denote $\hat{\partial}g(x)$ as the Frechet subdifferential of g at x which contains vectors $v \in \mathbb{E}$ satisfying

$$\liminf_{y \neq x, y \rightarrow x} \frac{1}{\|y - x\|} (g(y) - g(x) - \langle v, y - x \rangle) \geq 0.$$

If $x \notin \text{dom } g$, then we set $\hat{\partial}g(x) = \emptyset$.

- (iii) The limiting-subdifferential $\partial g(x)$ of g at $x \in \text{dom } g$ is defined as follows:

$$\partial g(x) := \left\{ v \in \mathbb{E} : \exists x^{(k)} \rightarrow x, g(x^{(k)}) \rightarrow g(x), v^{(k)} \in \hat{\partial}g(x^{(k)}), v^{(k)} \rightarrow v \right\}.$$

- (iv) The horizon subdifferential $\partial^\infty g(x)$ of g at x is defined as follows:

$$\partial^\infty g(x) := \left\{ v \in \mathbb{E} : \exists \lambda^{(k)} \rightarrow 0, \lambda^{(k)} \geq 0, \lambda^{(k)} x^{(k)} \rightarrow x, g(x^{(k)}) \rightarrow g(x), v^{(k)} \in \hat{\partial}g(x^{(k)}), v^{(k)} \rightarrow v \right\}.$$

Definition 4 We call $x^* \in \text{dom } F$ a critical point of F if $0 \in \partial F(x^*)$.

Definition 5 [49, Definition 7.5] A function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is called subdifferentially regular at \bar{x} if $f(\bar{x})$ is finite and the epigraph of f is Clarke regular at $(\bar{x}, f(\bar{x}))$ as a subset of $\mathbb{R}^n \times \mathbb{R}$ (see [49, Definition 6.4] for the definition of Clarke regularity of a set at a point).

Proposition 6 [49, Corollary 10.9] Suppose $f = f_1 + \dots + f_m$ for proper lower semi-continuous function $f_i : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ and let $\bar{x} \in \text{dom } f$. Suppose each function f_i is subdifferentially regular at \bar{x} , and the condition that the only combination of vector $\nu_i \in \partial^\infty f_i(\bar{x})$ with $\nu_1 + \dots + \nu_m = 0$ is $\nu_i = 0$ for $i \in [m]$. Then we have

$$\partial f(\bar{x}) = \partial f_1(\bar{x}) + \dots + \partial f_m(\bar{x}).$$

To obtain a global convergence, we need the following Kurdyka-Lojasiewicz (KL) property for $F(x) + h(y)$.

Definition 6 A function $\phi(\cdot)$ is said to have the KL property at $\bar{\mathbf{x}} \in \text{dom } \partial \phi$ if there exists $\varsigma \in (0, +\infty]$, a neighborhood U of $\bar{\mathbf{x}}$ and a concave function $\mathcal{Y} : [0, \varsigma) \rightarrow \mathbb{R}_+$ that is continuously differentiable on $(0, \varsigma)$, continuous at 0, $\mathcal{Y}(0) = 0$, and $\mathcal{Y}'(t) > 0$ for all $t \in (0, \eta)$, such that for all $\mathbf{x} \in U \cap [\phi(\bar{\mathbf{x}}) < \phi(\mathbf{x}) < \phi(\bar{\mathbf{x}}) + \varsigma]$, we have

$$\mathcal{Y}'(\phi(\mathbf{x}) - \phi(\bar{\mathbf{x}})) \text{dist}(0, \partial \phi(\mathbf{x})) \geq 1, \quad (25)$$

where $\text{dist}(0, \partial \phi(\mathbf{x})) = \min \{\|\mathbf{z}\| : \mathbf{z} \in \partial \phi(\mathbf{x})\}$. If $\phi(\mathbf{x})$ has the KL property at each point of $\text{dom } \partial \phi$ then ϕ is a KL function.

When $\mathcal{Y}(t) = ct^{1-\alpha}$, where c is a constant, we call α the KL coefficient.

Many non-convex non-smooth functions in practical applications belong to the class of KL functions, for examples, real analytic functions, semi-algebraic functions, and locally strongly convex functions, see for example [6, 7].

B Proofs

In this appendix, we provide the proofs of all propositions and theorems of our paper. Before that, let us give some preliminary results. We use x, z to denote vectors in \mathbb{R}^n .

Lemma 1 [21, Lemma 2.8] *If the function $x_i \mapsto \Theta(x_i, z)$ is ρ -strongly convex, differentiable at z_i , and $\nabla_{x_i} \Theta(z_i, z) = 0$ then we have*

$$\Theta(x_i, z) \geq \frac{\rho}{2} \|x_i - z_i\|^2.$$

We recall the notation $(x_i, z_{\neq i}) = (z_1, \dots, z_{i-1}, x_i, z_{i+1}, \dots, z_s)$. Suppose we are trying to solve

$$\min_x \Psi(x) := \Phi(x) + \sum_{i=1}^s g_i(x_i).$$

Proposition 7 [21, Theorem 2.7] *Suppose $\mathcal{G}_i^k : \mathbb{R}^{n_i} \times \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{n_i}$ be some extrapolation operator that satisfies $\mathcal{G}_i^k(x_i^k, x_i^{k-1}) \leq a_i^k \|x_i^k - x_i^{k-1}\|$. Let $u_i(x_i, z)$ is a block surrogate function of $\Phi(x)$. We assume one of the following conditions holds:*

- $x_i \mapsto u_i(x_i, z) + g_i(x_i)$ is ρ_i -strongly convex,
- the approximation error $\Theta(x_i, z) := u_i(x_i, z) - \Phi(x_i, z_{\neq i})$ satisfying $\Theta(x_i, z) \geq \frac{\rho_i}{2} \|x_i - z_i\|^2$ for all x_i .

Note that ρ_i may depend on z . Let

$$x_i^{k+1} = \operatorname{argmin}_{x_i} u_i(x_i, x^{k,i-1}) + g_i(x_i) - \langle \mathcal{G}_i^k(x_i^k, x_i^{k-1}), x_i \rangle.$$

Then we have

$$\Psi(x^{k,i-1}) + \gamma_i^k \|x_i^k - x_i^{k-1}\|^2 \geq \Psi(x^{k,i}) + \eta_i^k \|x_i^{k+1} - x_i^k\|^2, \quad (26)$$

where

$$\gamma_i^k = \frac{(a_i^k)^2}{2\nu\rho_i}, \quad \eta_i^k = \frac{(1-\nu)\rho_i}{2},$$

and $0 < \nu < 1$ is a constant. If we do not apply extrapolation, that is $a_i^k = 0$, then (26) is satisfied with $\gamma_i^k = 0$ and $\eta_i^k = \rho_i/2$.

The following proposition is derived from [20, Remark 3] and [62, Lemma 2.1].

Proposition 8 *Suppose $x_i \mapsto \Phi(x)$ is a L_i -smooth convex function and $g_i(x_i)$ is convex. Define $\bar{x}^{k,i-1} = (x_1^{k+1}, \dots, x_{i-1}^{k+1}, \bar{x}_i^k, x_{i+1}^k, \dots, x_s^k)$, $\hat{x}_i^k = x_i^k + \alpha_i^k(x_i^k - x_i^{k-1})$ and $\bar{x}_i^k = x_i^k + \beta_i^k(x_i^k - x_i^{k-1})$. Let $x_i^{k+1} = \operatorname{argmin}_{x_i} \langle \nabla \Phi(\bar{x}^{k,i-1}), x_i \rangle + g_i(x_i) + \frac{L_i}{2} \|x_i - \hat{x}_i^k\|^2$. Then we have Inequality (26) is satisfied with*

$$\gamma_i^k = \frac{L_i}{2} \left((\beta_i^k)^2 + \frac{(\gamma_i^k - \alpha_i^k)^2}{\nu} \right), \quad \eta_i^k = \frac{(1-\nu)L_i}{2}.$$

If $\alpha_i^k = \beta_i^k$ then we have Inequality (26) is satisfied with

$$\gamma_i^k = \frac{L_i}{2} (\beta_i^k)^2, \quad \eta_i^k = \frac{L_i}{2}.$$

B.1 Proof of Proposition 1

(i) Suppose we are updating x_i^k . Let us recall that

$$\mathcal{L}(x, y, \omega) := f(x) + \sum_{i=1}^s g_i(x_i) + h(y) + \varphi(x, y, \omega),$$

where

$$\varphi(x, y, \omega) = \frac{\beta}{2} \|\mathcal{A}x + \mathcal{B}y - b\|^2 + \langle \omega, \mathcal{A}x + \mathcal{B}y - b \rangle. \quad (27)$$

Denote $\mathbf{u}_i(x_i, z, y, \omega) = u_i(x_i, z) + h(y) + \hat{\varphi}_i(x_i, z, y, \omega)$, where

$$\hat{\varphi}_i(x_i, z, y, \omega) = \varphi(z, y, \omega) + \langle \mathcal{A}_i^*(\omega + \beta(\mathcal{A}z + \mathcal{B}y - b)), x_i - z_i \rangle + \frac{\kappa_i \beta}{2} \|x_i - z_i\|^2.$$

We see that $\hat{\varphi}_i(x_i, z, y, \omega)$ is a block surrogate function of $x \mapsto \varphi(x, y, \omega)$ with respect to block x_i , and $\mathbf{u}_i(x_i, z, y, \omega)$ is a block surrogate function of $x \mapsto f(x) + h(y) + \varphi(x, y, \omega)$ with respect to block x_i . The update in (8) can be rewritten as follows.

$$x_i^{k+1} = \underset{x_i}{\operatorname{argmin}} \mathbf{u}_i(x_i, x^{k,i-1}, y^k, \omega^k) + g_i(x_i) - \langle \mathcal{G}_i^k(x_i^k, x_i^{k-1}), x_i \rangle, \quad (28)$$

where

$$\mathcal{G}_i^k(x_i^k, x_i^{k-1}) = \beta \mathcal{A}_i^* \mathcal{A}(x^{k,i-1} - \bar{x}^{k,i-1}) + \kappa_i \beta \zeta_i^k (x_i^k - x_i^{k-1}). \quad (29)$$

The block approximation error function between $\mathbf{u}_i(x_i, z, y, \omega)$ and $x \mapsto f(x) + h(y) + \varphi(x, y, \omega)$ is defined as

$$\begin{aligned} \mathbf{e}_i(x_i, z, y, \omega) &= \mathbf{u}_i(x_i, z, y, \omega) - (f(x_i, z_{\neq i}) + h(y) + \varphi((x_i, z_{\neq i}), y, \omega)) \\ &= u_i(x_i, z) - f(x_i, z_{\neq i}) + \hat{\varphi}_i(x_i, z, y, \omega) - \varphi((x_i, z_{\neq i}), y, \omega) \\ &\geq \theta_i(x_i, z, y, \omega) := \end{aligned} \quad (30)$$

$$\varphi(z, y, \omega) - \varphi((x_i, z_{\neq i}), y, \omega) + \langle \mathcal{A}_i^*(\omega + \beta(\mathcal{A}z + \mathcal{B}y - b)), x_i - z_i \rangle + \frac{\kappa_i \beta}{2} \|x_i - z_i\|^2.$$

We have $\nabla_{x_i} \theta_i(x_i, z, y, \omega) = \kappa_i \beta (x_i - z_i) + \nabla_{x_i} \varphi(z, y, \omega) - \nabla_{x_i} \varphi((x_i, z_{\neq i}), y, \omega)$. So $\nabla_{x_i} \theta_i(z_i, z, y, \omega) = 0$. On the other hand, note that $x_i \mapsto \varphi((x_i, z_{\neq i}), y^k, \omega^k)$ is $\beta \|\mathcal{A}_i^* \mathcal{A}_i\|$ -smooth. So, $x_i \mapsto \theta_i(x_i, z, y, \omega)$ is a $\beta(\kappa_i - \|\mathcal{A}_i^* \mathcal{A}_i\|)$ -strongly convex function. From Lemma 1 we have $\theta_i(x_i, z, y, \omega) \geq \frac{\beta(\kappa_i - \|\mathcal{A}_i^* \mathcal{A}_i\|)}{2} \|x_i - z_i\|^2$. The result follows from (28), (30) and Proposition (7).

(ii) When $x_i \mapsto u_i(x_i, z) + g_i(x_i)$ is convex and we apply the update as in (8), it follows from Proposition 8 (see also [21, Remark 4.1]) that

$$\begin{aligned} u_i(x_i^k, x^{k,i-1}) + g_i(x_i^k) + \varphi(x^{k,i-1}, y^k, \omega^k) &+ \frac{\beta \|\mathcal{A}_i^* \mathcal{A}_i\|}{2} (\zeta_i^k)^2 \|x_i^k - x_i^{k-1}\|^2 \\ &\geq u_i(x_i^{k+1}, x^{k,i-1}) + g_i(x_i^{k+1}) + \varphi(x^{k,i}, y^k, \omega^k) + \frac{\beta \|\mathcal{A}_i^* \mathcal{A}_i\|}{2} \|x_i^{k+1} - x_i^k\|^2. \end{aligned} \quad (31)$$

On the other hand, note that $u_i(x_i^k, x^{k,i-1}) = f(x^{k,i-1})$ and $u_i(x_i^{k+1}, x^{k,i-1}) \geq f(x^{k,i})$. The result follows then.

B.2 Proof of Proposition 2

Denote

$$\hat{h}(y, y') = h(y') + \langle \omega, \mathcal{A}x + \mathcal{B}y' - b \rangle + \langle \mathcal{B}^* \omega + \nabla h(y'), y - y' \rangle + \frac{L_h}{2} \|y - y'\|^2.$$

Then we have $\hat{h}(y, y') + \frac{\beta}{2} \|\mathcal{A}x + \mathcal{B}y - b\|^2$ is a surrogate function of $y \mapsto h(y) + \varphi(x, y, \omega)$. Note that the function $y \mapsto \hat{h}(y, y') + \frac{\beta}{2} \|\mathcal{A}x + \mathcal{B}y - b\|^2$ is $(L_h + \beta \lambda_{\min}(\mathcal{B}^* \mathcal{B}))$ -strongly convex. The result follows from Proposition 7 (see also [21, Section 4.2.1]).

Suppose $h(y)$ is convex. We note that $y \mapsto \frac{\beta}{2} \|\mathcal{A}x + \mathcal{B}y - b\|^2$ is also convex and plays the role of g_i in Proposition 8. The result follows from Proposition 8.

B.3 Proof of Proposition 3

Note that

$$\mathcal{L}(x^{k+1}, y^{k+1}, \omega^{k+1}) = \mathcal{L}(x^{k+1}, y^{k+1}, \omega^k) + \frac{1}{\alpha\beta} \langle \omega^{k+1} - \omega^k, \omega^{k+1} - \omega^k \rangle \quad (32)$$

From the optimality condition of (9) we have

$$\nabla h(\hat{y}^k) + L_h(y^{k+1} - \hat{y}^k) + \mathcal{B}^* \omega^k + \beta \mathcal{B}^*(Ax^{k+1} + \mathcal{B}y^{k+1} - b) = 0.$$

Together with (10) we obtain

$$\nabla h(\hat{y}^k) + L_h(\Delta y^{k+1} - \delta_k \Delta y^k) + \mathcal{B}^* \omega^k + \frac{1}{\alpha} \mathcal{B}^*(w^{k+1} - w^k) = 0. \quad (33)$$

Hence,

$$\mathcal{B}^* w^{k+1} = (1 - \alpha) \mathcal{B}^* \omega^k - \alpha (\nabla h(\hat{y}^k) + L_h(\Delta y^{k+1} - \delta_k \Delta y^k)), \quad (34)$$

which implies that

$$\mathcal{B}^* \Delta w^{k+1} = (1 - \alpha) \mathcal{B}^* \Delta w^k - \alpha \Delta z^{k+1}, \quad (35)$$

where $\Delta z^{k+1} = z^{k+1} - z^k$ and $z^{k+1} = \nabla h(\hat{y}^k) + L_h(\Delta y^{k+1} - \delta_k \Delta y^k)$. We now consider 2 cases.

Case 1: $0 < \alpha \leq 1$. From the convexity of $\|\cdot\|$ we have

$$\|\mathcal{B}^* \Delta w^{k+1}\|^2 \leq (1 - \alpha) \|\mathcal{B}^* \Delta w^k\|^2 + \alpha \|\Delta z^{k+1}\|^2 \quad (36)$$

Case 2: $1 < \alpha < 2$. We rewrite (35) as $\mathcal{B}^* \Delta w^{k+1} = -(\alpha - 1) \mathcal{B}^* \Delta w^k - \frac{\alpha}{2 - \alpha} (2 - \alpha) \Delta z^{k+1}$. Hence

$$\|\mathcal{B}^* \Delta w^{k+1}\|^2 \leq (\alpha - 1) \|\mathcal{B}^* \Delta w^k\|^2 + \frac{\alpha^2}{(2 - \alpha)} \|\Delta z^{k+1}\|^2 \quad (37)$$

Combine (36) and (37) we obtain

$$\|\mathcal{B}^* \Delta w^{k+1}\|^2 \leq |1 - \alpha| \|\mathcal{B}^* \Delta w^k\|^2 + \frac{\alpha^2}{1 - |1 - \alpha|} \|\Delta z^{k+1}\|^2, \quad (38)$$

which implies

$$(1 - |1 - \alpha|) \|\mathcal{B}^* \Delta w^{k+1}\|^2 \leq |1 - \alpha| (\|\mathcal{B}^* \Delta w^k\|^2 - \|\mathcal{B}^* \Delta w^{k+1}\|^2) + \frac{\alpha^2}{1 - |1 - \alpha|} \|\Delta z^{k+1}\|^2. \quad (39)$$

On the other hand, when we use extrapolation for the update of y we have

$$\begin{aligned} \|\Delta z^{k+1}\|^2 &= \|\nabla h(\hat{y}^k) - \nabla h(\hat{y}^{k-1}) + L_h(\Delta y^{k+1} - \delta_k \Delta y^k) - L_h(\Delta y^k - \delta_{k-1} \Delta y^{k-1})\|^2 \\ &\leq 3L_h^2 \|\hat{y}^k - \hat{y}^{k-1}\|^2 + 3L_h^2 \|\Delta y^{k+1}\|^2 + 3\|(1 + \delta_k)L_h \Delta y^k - L_h \delta_{k-1} \Delta y^{k-1}\|^2 \\ &\leq 6L_h^2 [(1 + \delta_k)^2 \|\Delta y^k\|^2 + \delta_{k-1}^2 \|\Delta y^{k-1}\|^2] + 3L_h^2 \|\Delta y^{k+1}\|^2 \\ &\quad + 6(1 + \delta_k)^2 L_h^2 \|\Delta y^k\|^2 + 6L_h^2 \delta_{k-1}^2 \|\Delta y^{k-1}\|^2 \\ &= 3L_h^2 \|\Delta y^{k+1}\|^2 + 12(1 + \delta_k)^2 L_h^2 \|\Delta y^k\|^2 + 12L_h^2 \delta_{k-1}^2 \|\Delta y^{k-1}\|^2. \end{aligned} \quad (40)$$

If we do not use extrapolation for y then we have

$$\begin{aligned} \|\Delta z^{k+1}\|^2 &= \|\nabla h(y^k) - \nabla h(y^{k-1}) + L_h \Delta y^{k+1} - L_h \Delta y^k\|^2 \\ &\leq 3L_h^2 \|\Delta y^k\|^2 + 3L_h^2 \|\Delta y^{k+1}\|^2 + 3L_h^2 \|\Delta y^k\|^2 = 6L_h^2 \|\Delta y^k\|^2 + 3L_h^2 \|\Delta y^{k+1}\|^2. \end{aligned} \quad (41)$$

Furthermore, note that $\sigma_{\mathcal{B}} \|\Delta w^{k+1}\|^2 \leq \|\mathcal{B}^* \Delta w^{k+1}\|^2$. Therefore, it follows from (39) that

$$\begin{aligned} \|\Delta w^{k+1}\|^2 &\leq \frac{|1 - \alpha|}{\sigma_{\mathcal{B}}(1 - |1 - \alpha|)} (\|\mathcal{B}^* \Delta w^k\|^2 - \|\mathcal{B}^* \Delta w^{k+1}\|^2) \\ &\quad + \frac{\alpha^2 3L_h^2}{\sigma_{\mathcal{B}}(1 - |1 - \alpha|)^2} (\|\Delta y^{k+1}\|^2 + \bar{\delta}_k \|\Delta y^k\|^2 + 4\delta_{k-1}^2 \|\Delta y^{k-1}\|^2). \end{aligned} \quad (42)$$

The result is obtained from (42), (32) and Proposition 1.

B.4 Proof of Proposition 4

(i) From Inequality (17) and the conditions in (18),

$$\begin{aligned} & \mathcal{L}^{k+1} + \mu \|\Delta y^{k+1}\|^2 + \sum_{i=1}^s \eta_i \|\Delta x_i^{k+1}\|^2 + \frac{\alpha_1}{\beta} \|\mathcal{B}^* \Delta w^{k+1}\|^2 \\ & \leq \mathcal{L}^k + C_1 \mu \|\Delta y^k\|^2 + C_2 \mu \|\Delta y^{k-1}\|^2 + C_x \sum_{i=1}^s \eta_i \|\Delta x_i^k\|^2 + \frac{\alpha_1}{\beta} \|\mathcal{B}^* \Delta w^k\|^2. \end{aligned} \quad (43)$$

By summing from $k = 1$ to K Inequality (43) and noting that $C_1 + C_2 = C_y$ we obtain (20).

(ii) Let us prove $\{\Delta y^k\}$ and $\{\Delta x_i^k\}$ converge to 0. Let us first prove the second situation, that is we use extrapolation for the update of y and Inequality (19) is satisfied. From (34) we have $\alpha \mathcal{B}^* w^{k+1} = -(1 - \alpha) \mathcal{B}^* \Delta \omega^{k+1} - \alpha z^{k+1}$, where $z^{k+1} = \nabla h(\hat{y}^k) + L_h(\Delta y^{k+1} - \delta_k \Delta y^k)$. Using the same technique that derives Inequality (38), we obtain the following

$$\alpha \sigma_{\mathcal{B}} \|w^{k+1}\|^2 \leq \alpha \|\mathcal{B}^* w^{k+1}\|^2 \leq |1 - \alpha| \|\mathcal{B}^* \Delta \omega^{k+1}\|^2 + \frac{\alpha^2}{1 - |1 - \alpha|} \|z^{k+1}\|^2. \quad (44)$$

On the other hand, we have

$$\mathcal{L}^k = F(x^k) + h(y^k) + \frac{\beta}{2} \|\mathcal{A}x^k + \mathcal{B}y^k - b + \frac{\omega^k}{\beta}\|^2 - \frac{1}{2\beta} \|\omega^k\|^2 \geq F(x^k) + h(y^k) - \frac{1}{2\beta} \|\omega^k\|^2.$$

Together with (44) and

$$\begin{aligned} \|z^k\|^2 &= \|\nabla h(\hat{y}^{k-1}) - \nabla h(y^k) + \nabla h(y^k) + L_h(\Delta y^k - \delta_{k-1} \Delta y^{k-1})\|^2 \\ &\leq 4\|\nabla h(\hat{y}^{k-1}) - \nabla h(y^k)\|^2 + 4\|\nabla h(y^k)\|^2 + 4L_h^2 \|\Delta y^k\|^2 + 4L_h^2 \delta_{k-1}^2 \|\Delta y^{k-1}\|^2 \\ &\leq 12L_h^2 \|\Delta y^k\|^2 + 12\delta_{k-1}^2 \|\Delta y^{k-1}\|^2 + 4\|\nabla h(y^k)\|^2. \end{aligned}$$

we obtain

$$\begin{aligned} \mathcal{L}^k &\geq F(x^k) + h(y^k) - \frac{1}{2\alpha\beta\sigma_{\mathcal{B}}} (|1 - \alpha| \|\mathcal{B}^* \Delta \omega^k\|^2 + \frac{\alpha^2}{1 - |1 - \alpha|} \|z^k\|^2) \\ &\geq F(x^k) + h(y^k) - \frac{|1 - \alpha|}{2\alpha\beta\sigma_{\mathcal{B}}} \|\mathcal{B}^* \Delta \omega^k\|^2 \\ &\quad - \frac{\alpha}{2\beta\sigma_{\mathcal{B}}(1 - |1 - \alpha|)} (12L_h^2 \|\Delta y^k\|^2 + 12\delta_{k-1}^2 \|\Delta y^{k-1}\|^2 + 4\|\nabla h(y^k)\|^2) \end{aligned} \quad (45)$$

Since $h(y)$ is L_h -smooth, for all $y \in \mathbb{R}^q$ and $\alpha_L > 0$ we have, (see [40])

$$h(y - \alpha_L \nabla f(y)) \leq h(y) - \alpha_L (1 - \frac{L_h \alpha_L}{2}) \|\nabla h(y)\|^2.$$

Let us choose α_L such that $\alpha_L (1 - \frac{L_h \alpha_L}{2}) = \frac{4\alpha}{2\beta\sigma_{\mathcal{B}}(1 - |1 - \alpha|)}$. Note that this equation always has a positive solution when $\beta \geq \frac{4L_h \alpha}{\sigma_{\mathcal{B}}(1 - |1 - \alpha|)}$. Then we have

$$h(y^k) - \frac{4\alpha}{2\beta\sigma_{\mathcal{B}}(1 - |1 - \alpha|)} \|\nabla h(y^k)\|^2 \geq h(y^k - \alpha_L \nabla f(y^k)).$$

Together with (45) we get

$$\begin{aligned} \mathcal{L}^k &\geq F(x^k) + h(y^k - \alpha_L \nabla f(y^k)) - \frac{|1 - \alpha|}{2\alpha\beta\sigma_{\mathcal{B}}} \|\mathcal{B}^* \Delta \omega^k\|^2 \\ &\quad - \frac{\alpha}{2\beta\sigma_{\mathcal{B}}(1 - |1 - \alpha|)} (12L_h^2 \|\Delta y^k\|^2 + 12\delta_{k-1}^2 \|\Delta y^{k-1}\|^2). \end{aligned} \quad (46)$$

So from $\frac{\alpha_1}{\beta} \geq \frac{|1-\alpha|}{2\alpha\beta\sigma_{\mathcal{B}}}$, $\mu \geq \frac{\alpha 12L_h^2}{2\beta\sigma_{\mathcal{B}}(1-|\alpha|)}$, $(1-C_1)\mu \geq \frac{\alpha 12L_h^2 12\delta_k^2}{2\beta\sigma_{\mathcal{B}}(1-|\alpha|)}$ we have

$$\begin{aligned} & \mathcal{L}^{K+1} + \mu \|\Delta y^{K+1}\|^2 + \frac{\alpha_1}{\beta} \|B^* \Delta w^{K+1}\|^2 + (1-C_1)\mu \|\Delta y^K\|^2 \\ & \geq F(x^{K+1}) + h(y^{K+1} - \alpha_L \nabla f(y^{K+1})). \end{aligned} \quad (47)$$

Hence $\mathcal{L}^{K+1} + \mu \|\Delta y^{K+1}\|^2 + \frac{\alpha_1}{\beta} \|B^* \Delta w^{K+1}\|^2 + (1-C_1)\mu \|\Delta y^K\|^2$ is lower bounded.

Furthermore, since η_i and μ are positive numbers we derive from Inequality (20) that $\sum_{k=1}^{\infty} \|\Delta y^k\|^2 < +\infty$ and $\sum_{k=1}^{\infty} \|\Delta x_i^k\|^2 < +\infty$. Therefore, $\{\Delta y^k\}$ and $\{\Delta x_i^k\}$ converge to 0.

Let us now consider the first situation when $\delta_k = 0$ for all k .

From Inequality (17) and the conditions in (18) we have

$$\begin{aligned} & \mathcal{L}^{k+1} + \mu \|\Delta y^{k+1}\|^2 + \sum_{i=1}^s \eta_i \|\Delta x_i^{k+1}\|^2 + \frac{\alpha_1}{\beta} \|B^* \Delta w^{k+1}\|^2 \\ & \leq \mathcal{L}^k + C_y \mu \|\Delta y^k\|^2 + C_x \sum_{i=1}^s \eta_i \|\Delta x_i^k\|^2 + \frac{\alpha_1}{\beta} \|B^* \Delta w^k\|^2. \end{aligned} \quad (48)$$

By summing Inequality (48) from $k = 1$ to K we obtain

$$\begin{aligned} & \mathcal{L}^{K+1} + C_y \mu \|\Delta y^{K+1}\|^2 + C_x \sum_{i=1}^s \eta_i \|\Delta x_i^{K+1}\|^2 + \frac{\alpha_1}{\beta} \|B^* \Delta w^{K+1}\|^2 \\ & \quad + \sum_{k=1}^K [(1-C_y)\mu \|\Delta y^{k+1}\|^2 + (1-C_x) \sum_{i=1}^s \eta_i \|\Delta x_i^{k+1}\|^2] \\ & \leq \mathcal{L}^1 + \frac{\alpha_1}{\beta} \|B^* \Delta w^1\|^2 + \sum_{i=1}^s \eta_i^0 \|\Delta x_i^1\|^2 + C\mu \|\Delta y^1\|^2. \end{aligned} \quad (49)$$

Denote the value of the right side of Inequality (48) by $\hat{\mathcal{L}}^k$. Note that $0 < C_x, C_y < 1$, then from (48) we have the sequence $\{\hat{\mathcal{L}}^k\}$ is non-increasing. It follows from [39, Lemma 2.9] that $\hat{\mathcal{L}}^k \geq \vartheta$ for all k , where ϑ is the lower bound of $F(x^k) + h(y^k)$. For completeness, let us provide the proof in the following. We have

$$\begin{aligned} \hat{\mathcal{L}}^k & \geq \mathcal{L}^k = F(x^k) + h(y^k) + \frac{\beta}{2} \|Ax^k + By^k - b\|^2 + \frac{1}{\alpha\beta} \langle \omega^k, \omega^k - \omega^{k-1} \rangle \\ & \geq \vartheta + \frac{1}{2\alpha\beta} (\|\omega^k\|^2 - \|\omega^{k-1}\|^2 + \|\Delta\omega^k\|^2) \geq \vartheta + \frac{1}{2\alpha\beta} (\|\omega^k\|^2 - \|\omega^{k-1}\|^2), \end{aligned} \quad (50)$$

Assume that there exists k_0 such that $\hat{\mathcal{L}}^k < \vartheta$ for all $k \geq k_0$. As $\hat{\mathcal{L}}^k$ is non-increasing we have

$$\sum_{k=1}^K (\hat{\mathcal{L}}^k - \vartheta) \leq \sum_{k=1}^{k_0} (\hat{\mathcal{L}}^k - \vartheta) + (K - k_0)(\hat{\mathcal{L}}^k - \vartheta)$$

Hence $\sum_{k=1}^{\infty} (\hat{\mathcal{L}}^k - \vartheta) = -\infty$. However, from (50) we have

$$\sum_{k=1}^K (\hat{\mathcal{L}}^k - \vartheta) \geq \sum_{k=1}^K \frac{1}{2\alpha\beta} \|\omega^k\|^2 - \frac{1}{2\alpha\beta} \|\omega^{k-1}\|^2 \geq \frac{1}{2\alpha\beta} (-\|\omega^0\|^2),$$

which gives a contradiction.

Since $\hat{\mathcal{L}}^K \geq \vartheta$ and η_i and μ are positive numbers we derive from Inequality (20) that $\sum_{k=1}^{\infty} \|\Delta y^k\|^2 < +\infty$ and $\sum_{k=1}^{\infty} \|\Delta x_i^k\|^2 < +\infty$. Therefore, $\{\Delta y^k\}$ and $\{\Delta x_i^k\}$ converge to 0.

Now we prove $\{\Delta\omega^k\}$ goes to 0. Since $\sum_{k=1}^{\infty} \|\Delta y^k\|^2 < +\infty$, we derive from (40) that $\sum_{k=1}^{\infty} \|\Delta z^k\|^2 < +\infty$. Summing up Equality (38) from $k = 1$ to K we have

$$(1 - |1 - \alpha|) \sum_{k=1}^K \|B^* \Delta\omega^k\|^2 + \|B^* \Delta\omega^{K+1}\|^2 \leq \|B^* \Delta\omega^1\|^2 + \frac{\alpha^2}{1 - |1 - \alpha|} \sum_{k=1}^K \|\Delta z^{k+1}\|^2,$$

which implies that $\sum_{k=1}^{\infty} \|B^* \Delta\omega^k\|^2 < +\infty$. Hence, $\|B^* \Delta\omega^k\|^2 \rightarrow 0$. Since $\sigma_{\mathcal{B}} > 0$ we have $\{\Delta\omega^k\}$ goes to 0.

B.5 Proof of Proposition 5

We remark that we use the idea in the proof of [58, Lemma 6] to prove the proposition. However, our proof is more complicated since in our framework $\alpha \in (0, 2)$, the function h is linearized and we use extrapolation for y .

Note that as $\sigma_{\mathcal{B}} > 0$ we have \mathcal{B} is a surjective. Together with the assumption $b + \text{Im}(\mathcal{A}) \subseteq \text{Im}(\mathcal{B})$ we have there exist \bar{y}^k such that $\mathcal{A}x^k + \mathcal{B}\bar{y}^k - b = 0$.

Now we have

$$\begin{aligned} \mathcal{L}^k &= F(x^k) + h(y^k) + \frac{\beta}{2} \|\mathcal{A}x^k + \mathcal{B}y^k - b\|^2 + \langle \omega^k, \mathcal{A}x^k + \mathcal{B}y^k - b \rangle \\ &= F(x^k) + h(y^k) + \frac{\beta}{2} \|\mathcal{A}x^k + \mathcal{B}y^k - b\|^2 + \langle \mathcal{B}^* \omega^k, y^k - \bar{y}^k \rangle \end{aligned} \quad (51)$$

From (33) we have

$$\begin{aligned} \langle \mathcal{B}^* \omega^k, y^k - \bar{y}^k \rangle &= \langle \nabla h(\hat{y}^k) + L_h(\Delta y^{k+1} - \delta_k \Delta y^k) + \frac{1}{\alpha} \mathcal{B}^*(w^{k+1} - w^k), \bar{y}^k - y^k \rangle \\ &\geq \langle \nabla h(y^k), \bar{y}^k - y^k \rangle - (\|\nabla h(y^k) - \nabla h(\hat{y}^k)\| + L_h \|\Delta y^{k+1}\| + L_h \delta_k \|\Delta y^k\| \\ &\quad + \frac{1}{\alpha} \|\mathcal{B}^* \Delta \omega^{k+1}\|) \|\bar{y}^k - y^k\|. \end{aligned}$$

Therefore, it follows from (51) and L_h -smooth property of h that

$$\mathcal{L}^k \geq F(x^k) + h(\bar{y}^k) - \frac{L_h}{2} \|y^k - \bar{y}^k\|^2 - (2L_h \delta_k \|\Delta y^k\| + L_h \|\Delta y^{k+1}\| + \frac{1}{\alpha} \|\mathcal{B}^* \Delta \omega^{k+1}\|) \|y^k - \bar{y}^k\|. \quad (52)$$

On the other hand, we have

$$\|y^k - \bar{y}^k\|^2 \leq \frac{1}{\lambda_{\min}(\mathcal{B}^* \mathcal{B})} \|\mathcal{B}(\bar{y}^k - y^k)\|^2 = \frac{1}{\lambda_{\min}(\mathcal{B}^* \mathcal{B})} \|\mathcal{A}x^k + \mathcal{B}y^k - b\|^2 = \frac{1}{\lambda_{\min}(\mathcal{B}^* \mathcal{B})} \left\| \frac{1}{\alpha \beta} \Delta \omega^k \right\|^2. \quad (53)$$

We have proved in Proposition 4 that $\|\Delta \omega^k\|$, $\|\Delta x^k\|$ and $\|\Delta y^k\|$ converge to 0. Furthermore, from Proposition 4 we have \mathcal{L}^k is upper bounded. Therefore, from (52), (53) and (20) we have $F(x^k) + h(\bar{y}^k)$ is upper bounded. So $\{x^k\}$ is bounded. Consequently, $\mathcal{A}x^k$ is bounded.

Furthermore, we have

$$\|y^k\|^2 \leq \frac{1}{\lambda_{\min}(\mathcal{B}^* \mathcal{B})} \|\mathcal{B}y^k\|^2 = \frac{1}{\lambda_{\min}(\mathcal{B}^* \mathcal{B})} \left\| \frac{1}{\alpha \beta} \Delta \omega^k - \mathcal{A}x^k - b \right\|^2.$$

Therefore, $\{y^k\}$ is bounded, which implies that $\|\nabla h(\hat{y}^k)\|$ is also bounded. Finally, from (33) and the assumption $\lambda_{\min}(\mathcal{B}\mathcal{B}^*) > 0$ we also have $\{\omega^k\}$ is bounded.

B.6 Proof of Theorem 1

Suppose $(x^{k_n}, y^{k_n}, \omega^{k_n})$ converges to (x^*, y^*, ω^*) . Since $\Delta x_i^{k_i}$ goes to 0, we have $x_i^{k_n+1}$ and $x_i^{k_n-1}$ also converge to x_i^* for all $i \in [s]$. From (28), for all x_i ,

$$\mathbf{u}_i(x_i^{k+1}, x^{k,i-1}, y^k, \omega^k) + g_i(x_i^{k+1}) \leq \mathbf{u}_i(x_i, x^{k,i-1}, y^k, \omega^k) + g_i(x_i) - \langle \mathcal{G}_i^k(x_i^k, x_i^{k-1}), x_i - x_i^{k+1} \rangle. \quad (54)$$

Choosing $x_i = x_i^*$ and $k = k_n - 1$ in (54) and noting that $\mathbf{u}_i(x_i, z)$ is continuous by Assumption 2 (i), we have $\limsup_{n \rightarrow \infty} \mathbf{u}_i(x_i^*, x^*, y^*, \omega^*) + g_i(x_i^{k_n}) \leq \mathbf{u}_i(x_i^*, x^*, y^*, \omega^*) + g_i(x_i^*)$. On the other hand, as $g_i(x_i)$ is lower semi-continuous. Hence, $g_i(x_i^{k_n})$ converges to $g_i(x_i^*)$. Now we choose $k = k_n \rightarrow \infty$ in (54) for all x_i we obtain

$$\begin{aligned} L_0(x^*, y^*, \omega^*) + g_i(x_i^*) &\leq \mathbf{u}_i(x_i, x^*, y^*, \omega^*) + g_i(x_i) \\ &= L_0(x_i, x_{\neq i}^*, y^*, \omega^*) + \mathbf{e}_i(x_i, x^*, y^*, \omega^*) + g_i(x_i), \end{aligned} \quad (55)$$

where $L_0(x, y, \omega) = f(x) + h(y) + \varphi(x, y, \omega)$ and \mathbf{e}_i is the approximation error defined in (30). We have

$$\begin{aligned} \mathbf{e}_i(x_i, x^*, y^*, \omega^*) &= u_i(x_i, x^*) - f(x_i, x_{\neq i}^*) + \hat{\varphi}_i(x_i, x^*, y^*, \omega^*) - \varphi((x_i, x_{\neq i}^*), y^*, \omega^*) \\ &\leq \bar{e}_i(x_i, x^*) + \hat{\varphi}_i(x_i, x^*, y^*, \omega^*) - \varphi((x_i, x_{\neq i}^*), y^*, \omega^*). \end{aligned}$$

Note that $\bar{e}_i(x_i^*, x^*) = 0$ by Assumption 2. From (55) we have x_i^* is a solution of

$$\min_{x_i} L(x_i, x_{\neq i}^*, y^*, \omega^*) + \bar{e}_i(x_i, x^*) + \hat{\varphi}_i(x_i, x^*, y^*, \omega^*) - \varphi((x_i, x_{\neq i}^*), y^*, \omega^*).$$

Writing the optimality condition for this problem we obtain $0 \in \partial_{x_i} \mathcal{L}(x^*, y^*, \omega^*)$. Totally similarly we can prove that $0 \in \partial_y \mathcal{L}(x^*, y^*, \omega^*)$. On the other hand, we have

$$\Delta \omega^k = \omega^k - \omega^{k-1} = \alpha \beta (\mathcal{A}x^k + \mathcal{B}y^k - b) \rightarrow 0.$$

Hence, $\partial_\omega \mathcal{L}(x^*, y^*, \omega^*) = \mathcal{A}x^* + \mathcal{B}y^* - b = 0$.

As we assume $\partial F(x) = \partial_{x_1} F(x) \times \dots \times \partial_{x_s} F(x)$, we have

$$\begin{aligned} \partial \mathcal{L}(x, y, \omega) &= \partial F(x) + \nabla \left(h(y) + \langle \omega, \mathcal{A}x + \mathcal{B}y - b \rangle + \frac{\beta}{2} \|\mathcal{A}x + \mathcal{B}y - b\|^2 \right) \\ &= \partial_{x_1} \mathcal{L}(x, y, \omega) \times \dots \times \partial_{x_s} \mathcal{L}(x, y, \omega) \times \partial_y \mathcal{L}(x, y, \omega) \times \partial_\omega \mathcal{L}(x, y, \omega). \end{aligned}$$

So $0 \in \partial \mathcal{L}(x^*, y^*, \omega^*)$.

B.7 Proof of Theorem 2

Note that we assume the generated sequence of Algorithm 1 is bounded. The following analysis is considered in the bounded set that contains the generated sequence of Algorithm 1. We first prove some preliminary results.

(A) The optimality condition of (28) gives us

$$\begin{aligned} \mathcal{G}_i^k(x_i^k - x_i^{k-1}) - \mathcal{A}_i^*(\omega^k + \beta(\mathcal{A}x^{k,i-1} + \mathcal{B}y^k - b)) - \kappa_i \beta(x_i^{k+1} - x_i^k) \\ \in \partial_{x_i}(u_i(x_i^{k+1}, x^{k,i-1}) + g_i(x_i^{k+1})). \end{aligned} \quad (56)$$

As (22) holds, there exists $\mathbf{s}_i^{k+1} \in \partial u_i(x_i^{k+1}, x^{k,i-1})$ and $\mathbf{t}_i^{k+1} \in \partial g_i(x_i^{k+1})$ such that

$$\mathcal{G}_i^k(x_i^k - x_i^{k-1}) - \mathcal{A}_i^*(\omega^k + \beta(\mathcal{A}x^{k,i-1} + \mathcal{B}y^k - b)) - \kappa_i \beta(x_i^{k+1} - x_i^k) = \mathbf{s}_i^{k+1} + \mathbf{t}_i^{k+1} \quad (57)$$

As (23) holds, there exists $\xi_i^{k+1} \in \partial_{x_i} f(x^{k+1})$ such that

$$\|\xi_i^{k+1} - \mathbf{s}_i^{k+1}\| \leq L_i \|x^{k+1} - x^{k,i-1}\|. \quad (58)$$

Denote $\tau_i^{k+1} := \xi_i^{k+1} + \mathbf{t}_i^{k+1} \in \partial_{x_i} F(x^{k+1})$ (as (22) holds). Then, from (57) we have

$$\tau_i^{k+1} = \xi_i^{k+1} + \mathcal{G}_i^k(x_i^k - x_i^{k-1}) - \mathcal{A}_i^*(\omega^k + \beta(\mathcal{A}x^{k,i-1} + \mathcal{B}y^k - b)) - \kappa_i \beta(x_i^{k+1} - x_i^k) - \mathbf{s}_i^{k+1}. \quad (59)$$

On the other hand, we note that

$$\partial_{x_i} \mathcal{L}(x^{k+1}, y^{k+1}, \omega^{k+1}) = \partial_{x_i} F(x^{k+1}) + \mathcal{A}_i^*(\omega^{k+1} + \beta(\mathcal{A}x^{k+1} + \mathcal{B}y^{k+1} - b)). \quad (60)$$

Let $d_i^{k+1} := \tau_i^{k+1} + \mathcal{A}_i^*(\omega^{k+1} + \beta(\mathcal{A}x^{k+1} + \mathcal{B}y^{k+1} - b)) \in \partial_{x_i} \mathcal{L}(x^{k+1}, y^{k+1}, \omega^{k+1})$. From (59),

$$\begin{aligned} \|d_i^{k+1}\| &= \left\| \xi_i^{k+1} + \mathcal{G}_i^k(x_i^k - x_i^{k-1}) - \mathcal{A}_i^*(\omega^k + \beta(\mathcal{A}x^{k,i-1} + \mathcal{B}y^k - b)) - \kappa_i \beta(x_i^{k+1} - x_i^k) \right. \\ &\quad \left. - \mathbf{s}_i^{k+1} + \mathcal{A}_i^*(\omega^{k+1} + \beta(\mathcal{A}x^{k+1} + \mathcal{B}y^{k+1} - b)) \right\| \end{aligned} \quad (61)$$

Together with (58) we obtain

$$\begin{aligned} \|d_i^{k+1}\| &\leq a_i^k \|\Delta x_i^k\| + \beta \|\mathcal{A}_i^* A\| \|x^{k+1} - x^{k,i-1}\| + \beta \|\mathcal{A}_i^* \mathcal{B}\| \|\Delta y^{k+1}\| + \|\mathcal{A}_i^*\| \|\Delta \omega^{k+1}\| \\ &\quad + \kappa_i \beta \|\Delta x_i^{k+1}\| + L_i \|x^{k+1} - x^{k,i-1}\|. \end{aligned} \quad (62)$$

It follows from (9) that

$$\mathcal{B}^* \omega^k + \nabla h(\hat{y}^k) + \beta \mathcal{B}^* (\mathcal{A}x^{k+1} + \mathcal{B}y^{k+1} - b) + L_h (y^{k+1} - \hat{y}^k) = 0.$$

Let $d_y^{k+1} := \nabla h(y^{k+1}) + \mathcal{B}^* (\omega^{k+1} + \beta (\mathcal{A}x^{k+1} + \mathcal{B}y^{k+1} - b))$. Then $d_y^{k+1} \in \partial_y \mathcal{L}(x^{k+1}, y^{k+1}, \omega^{k+1})$ and

$$\begin{aligned} \|d_y^{k+1}\| &= \|\nabla h(y^{k+1}) - \nabla h(\hat{y}^k) + \mathcal{B}^* (\omega^{k+1} - \omega^k) - L_h (y^{k+1} - \hat{y}^k)\| \\ &\leq 2L_h \|y^{k+1} - \hat{y}^k\| + \|\mathcal{B}^*\| \|\Delta \omega^{k+1}\| \leq 2L_h (\|\Delta y^{k+1}\| + \delta_k \|\Delta y^k\|) + \|\mathcal{B}^*\| \|\Delta \omega^{k+1}\|. \end{aligned}$$

Let $d_\omega^{k+1} := \mathcal{A}x^{k+1} + \mathcal{B}y^{k+1} - b$. We have $d_\omega^{k+1} \in \partial_\omega \mathcal{L}(x^{k+1}, y^{k+1}, \omega^{k+1})$ and

$$d_\omega^{k+1} = (\omega^{k+1} - \omega^k) / (\alpha\beta) = \Delta \omega^{k+1} / (\alpha\beta).$$

(B) Let us now prove $F(x^{k_n})$ converges to $F(x^*)$. This implies $\mathcal{L}(x^{k_n}, y^{k_n}, \omega^{k_n})$ converges to $\mathcal{L}(x^*, y^*, \omega^*)$ since \mathcal{L} is differentiable in y and ω . We have

$$F(x^{k_n}) = f(x^{k_n}) + \sum_{i=1}^s g_i(x_i^{k_n}) = u_s(x_s^{k_n}, x^{k_n}) + \sum_{i=1}^s g_i(x_i^{k_n}).$$

So $F(x^{k_n})$ converges to $u_s(x_s^*, x^*) + \sum_{i=1}^s g_i(x_i^*) = F(x^*)$.

We now proceed to prove the global convergence. Denote $\mathbf{z} = (x, y, \omega)$, $\tilde{\mathbf{z}} = (\tilde{x}, \tilde{y}, \tilde{\omega})$, and $\mathbf{z}^k = (x^k, y^k, \omega^k)$. We consider the following auxiliary function

$$\tilde{\mathcal{L}}(\mathbf{z}, \tilde{\mathbf{z}}) = \mathcal{L}(x, y, \omega) + \sum_{i=1}^s \frac{\eta_i + C_x \eta_i}{2} \|x_i - \tilde{x}_i\|^2 + \frac{(1 + C_y)\mu}{2} \|y - \tilde{y}\|^2 + \frac{\alpha_1}{\beta} \|B^*(\omega - \tilde{\omega})\|^2.$$

The auxiliary sequence $\tilde{\mathcal{L}}(\mathbf{z}^k, \mathbf{z}^{k-1})$ has the following properties.

1. **Sufficient decreasing property.** From (48) we have

$$\begin{aligned} \tilde{\mathcal{L}}(\mathbf{z}^{k+1}, \mathbf{z}^k) &+ \sum_{i=1}^s \frac{\eta_i - C_x \eta_i}{2} (\|x_i^{k+1} - x_i^k\|^2 + \|x_i^k - x_i^{k-1}\|^2) \\ &+ \frac{(1 - C_y)\mu}{2} (\|y^{k+1} - y^k\|^2 + \|y^k - y^{k-1}\|^2) \leq \tilde{\mathcal{L}}(\mathbf{z}^k, \mathbf{z}^{k-1}). \end{aligned}$$

2. **Boundedness of subgradient.** In the proof (A) above, we have proved that

$$\|d^{k+1}\| \leq a_1 (\|x^{k+1} - x^k\| + \|x^k - x^{k-1}\| + \|y^{k+1} - y^k\| + \|\omega^{k+1} - \omega^k\|)$$

for some constant a_1 and $d^{k+1} \in \partial \mathcal{L}(\mathbf{z}^{k+1})$. On the other hand, as we use $\alpha = 1$, from (35) we obtain

$$\begin{aligned} \sqrt{\sigma_B} \|\omega^{k+1} - \omega^k\| &\leq \|B^*(\omega^{k+1} - \omega^k)\| = \|\Delta z^{k+1}\| \\ &= \|\nabla h(y^k) - \nabla h(y^{k-1}) + L_h(\Delta y^{k+1} - \Delta y^k)\| \leq 2L_h \|y^k - y^{k-1}\| + L_h \|y^{k+1} - y^k\|. \end{aligned} \quad (63)$$

Hence,

$$\|d^{k+1}\| \leq a_2 (\|x^{k+1} - x^k\| + \|x^k - x^{k-1}\| + \|y^{k+1} - y^k\| + \|y^k - y^{k-1}\|)$$

for some constant a_2 . Note that

$$\partial \tilde{\mathcal{L}}(\mathbf{z}, \tilde{\mathbf{z}}) = \partial \mathcal{L}(\mathbf{z}, \tilde{\mathbf{z}}) + \partial \left(\sum_{i=1}^s \frac{\eta_i + C_x \eta_i}{2} \|x_i - \tilde{x}_i\|^2 + \frac{(1 + C_y)\mu}{2} \|y - \tilde{y}\|^2 + \frac{\alpha_1}{\beta} \|B^*(\omega - \tilde{\omega})\|^2 \right).$$

Hence, it is not difficult to show that

$$\|\mathbf{d}^{k+1}\| \leq a_3 (\|x^{k+1} - x^k\| + \|x^k - x^{k-1}\| + \|y^{k+1} - y^k\| + \|y^k - y^{k-1}\|)$$

for some constant a_3 and $\mathbf{d}^{k+1} \in \partial \tilde{\mathcal{L}}(\mathbf{z}^{k+1}, \mathbf{z}^k)$.

3. **KL property.** Since $F(x) + h(y)$ has KL property, then $\tilde{\mathcal{L}}(\mathbf{z}, \bar{\mathbf{z}})$ also has KL property.
4. **A continuity condition.** Suppose \mathbf{z}^{k_n} converges to (x^*, y^*, ω^*) . In the proof (B) above, we have proved that $\mathcal{L}(\mathbf{z}^{k_n})$ converges to $\mathcal{L}(x^*, y^*, \omega^*)$. Furthermore, from Proposition 4 we proved that $\|\mathbf{z}^{k+1} - \mathbf{z}^k\|$ goes to 0. Hence we have \mathbf{z}^{k_n-1} converges to (x^*, y^*, ω^*) . So, $\tilde{\mathcal{L}}(\mathbf{z}^{k+1}, \mathbf{z}^k)$ converges to $\tilde{\mathcal{L}}(\mathbf{z}^*, \mathbf{z}^*)$.

Using the same technique as in [7, Theorem 1], see also [20, 41], we can prove that

$$\sum_{k=1}^{\infty} (\|x^{k+1} - x^k\| + \|x^k - x^{k-1}\| + \|y^{k+1} - y^k\| + \|y^k - y^{k-1}\|) < \infty.$$

which implies $\{(x^k, y^k)\}$ converges to (x^*, y^*) . From (63) we obtain

$$\sum_{k=1}^{\infty} \|\omega^{k+1} - \omega^k\| \leq \sum_{k=1}^{\infty} (\|y^{k+1} - y^k\| + \|y^k - y^{k-1}\|) < \infty.$$

Hence, $\{\omega^k\}$ also converges to ω^* .

C Additional Experiment for different values of α

In this experiment, we rerun the experiments from Section 3 with other values for α , namely 0.5, 1.4 and 1.8; see Figures 2-4 (on pages 31-33). The penalty parameter β is computed by $\beta = 2(2 + C_y)\alpha_2/C_y$, where $C_y = 1 - 10^{-6}$ and $\alpha_2 = \frac{3\alpha}{(1-|1-\alpha|)^2}$. Although the segmentation errors and objective function values differ for different values of α , we observe that, in all cases, iADMM-mm outperforms ADMM-mm which outperforms linearizedADMM. This confirms our observations from Section 3. On the other hand, we observe that the performances of ADMM-mm and linearizedADMM are similar for different values of α ; however, the performances of iADMM-mm (that is, ADMM-mm with inertial terms) for $\alpha = 0.5$ and $\alpha = 1.4$ are slightly worse than for $\alpha = 1$, and the value $\alpha = 1.8$ leads to significantly worse performances for iADMM-mm. It is known that, in the convex setting, the ADMM variants often perform better for $\alpha > 1$. However, in our experiments, $\alpha = 1$ provides the best performance for iADMM-mm. A possible reason is that the global convergence of iADMM-mm has been established only for the case $\alpha = 1$ (see Theorem 2) while $\alpha \in (0, 2)$ only guarantees a subsequential convergence (see Theorem 1).

D Additional experiments for a regularized nonnegative matrix factorization problem

In the previous example, the function $f(X, Y) = \lambda_1 \|X\|_* + r_2(Y)$ was separable while our framework allows non-separable functions; see (1) and the discussion that follows. To illustrate the use and effectiveness of iADMM on a non-separable case, let us consider the following regularized nonnegative matrix factorization (NMF) problem

$$\min_{W \in \mathbb{R}_+^{n \times r}, H \in \mathbb{R}_+^{r \times m}} 1/2 \|X - WH\|^2 + c_1 \|W\|_F^2 + c_2 \|H\|_F^2, \quad (64)$$

where $X \in \mathbb{R}^{n \times m}$ is a given nonnegative matrix, and $c_1 > 0$ and $c_2 > 0$ are regularized parameters. Problem (64) can be rewritten in the form of (1) as follows:

$$\min_{W \in \mathbb{R}_+^{n \times r}, H \in \mathbb{R}_+^{r \times m}} 1/2 \|X - WH\|^2 + c_1 \|W\|_F^2 + c_2 \|Y\|_F^2, \quad (65)$$

$$\text{such that } H - Y = 0.$$

In this case, $x_1 = W$, $x_2 = H$, $y = Y$, $f(W, H) = \frac{1}{2} \|X - WH\|^2 + \beta \|W\|_F^2$, $g_1(W)$ and $g_2(H)$ are indicator functions of $\mathbb{R}_+^{n \times r}$ and $\mathbb{R}_+^{r \times m}$ respectively, $h(Y) = c_2 \|Y\|_F^2$, $\mathcal{A}_1 = 0$, $\mathcal{A}_2 = \mathcal{I}$,

$\mathcal{B} = -\mathcal{I}$ (where \mathcal{I} is identity operator), and $b = 0$. As $W \mapsto f(W, H)$ is L_W -Lipschitz smooth and $H \mapsto f(W, H)$ is L_H -Lipschitz smooth, where $L_W = \|HH^\top\| + 2c_1$ and $L_H = \|W^\top W\|$, we use the Lipschitz gradient surrogate for block W and H as in (12), and apply the inertial term as in the footnote 3 (that is, we apply inertial terms that also lead to the extrapolation for the block surrogate of f). The augmented Lagrangian for (65) is

$$\mathcal{L}(W, H, Y, \omega) = f(W, H) + h(Y) + \langle H - Y, \omega \rangle + \frac{\beta}{2} \|H - Y\|^2.$$

Applying iADMM for solving (65), the update of W is

$$\begin{aligned} W^{k+1} &\in \arg \min_{W \in \mathbb{R}_+^{n \times r}} \langle -(X - \bar{W}^k H^k)(H^k)^\top + 2c_1 \bar{W}^k, W \rangle + \frac{L_W(H^k)}{2} \|W - \bar{W}^k\|^2 \\ &= \max \left\{ \bar{W}^k - \frac{1}{L_W(H^k)} \left(-(X - \bar{W}^k H^k)(H^k)^\top + 2c_1 \bar{W}^k \right), 0 \right\}, \end{aligned} \quad (66)$$

where $\bar{W}^k = W^k + \zeta_1^k (W^k - W^{k-1})$. Note that we have used extrapolation for the surrogate of $W \mapsto f(W, H)$. The update of H is

$$\begin{aligned} H^{k+1} &\in \arg \min_{H \in \mathbb{R}_+^{r \times m}} \langle -(W^{k+1})^\top (X - W^{k+1} \bar{H}^k) + \omega^k + \beta(\bar{H}^k - Y^k), H \rangle \\ &\quad + \frac{\beta + L_H(W^{k+1})}{2} \|H - \bar{H}^k\|^2 \\ &= \max \left\{ \bar{H}^k - \frac{1}{\beta + L_H(W^{k+1})} \left(-(W^{k+1})^\top (X - W^{k+1} \bar{H}^k) + \omega^k + \beta(\bar{H}^k - Y^k) \right), 0 \right\}, \end{aligned} \quad (67)$$

where $\bar{H}^k = H^k + \zeta_2^k (H^k - H^{k-1})$. We do not use extrapolation for Y (that is, $\delta_k = 0$), and simply choose $\alpha = 1$. The update of Y is

$$\begin{aligned} Y^{k+1} &\in \arg \min_Y \langle -\omega^k + 2c_2 Y^k, Y \rangle + \frac{\beta}{2} \|Y - H^{k+1}\|^2 + c_2 \|Y - Y^k\|^2 \\ &= \frac{1}{\beta + 2c_2} (\beta H^{k+1} + \omega^k), \end{aligned} \quad (68)$$

while the update of ω is

$$\omega^{k+1} = \omega^k + \beta(H^{k+1} - Y^{k+1}).$$

Choosing parameters By Proposition 8, the update of W in (66) implies that Inequality (14) is satisfied:

$$\mathcal{L}(W^{k+1}, H^k, Y^k, \omega^k) + \eta_1^k \|W^{k+1} - W^k\|^2 \leq \mathcal{L}(W^k, H^k, Y^k, \omega^k) + \gamma_1^k \|W^k - W^{k-1}\|^2,$$

where

$$\eta_1^k = \frac{L_W(H^k)}{2}, \quad \gamma_1^k = \frac{L_W(H^k)}{2} (\zeta_1^k)^2.$$

Note that we use η_1^k instead of η_1 as this value varies along with the update of H (because we used the extrapolation for the surrogate of $W \mapsto f(W, H)$). Similarly, the update of H in (67) implies that Inequality (14) is satisfied:

$$\mathcal{L}(W^{k+1}, H^{k+1}, Y^k, \omega^k) + \eta_2^k \|H^{k+1} - H^k\|^2 \leq \mathcal{L}(W^{k+1}, H^k, Y^k, \omega^k) + \gamma_2^k \|H^k - H^{k-1}\|^2,$$

where

$$\eta_2^k = \frac{L_H(W^{k+1}) + \beta}{2}, \quad \gamma_2^k = \frac{L_H(W^{k+1}) + \beta}{2} (\zeta_2^k)^2.$$

Because of the update of Y in (68), the inequality in Proposition (2) is satisfied:

$$\mathcal{L}(W^{k+1}, H^{k+1}, Y^{k+1}, \omega^k) + \eta_y \|Y^{k+1} - Y^k\|^2 \leq \mathcal{L}(W^{k+1}, H^{k+1}, Y^k, \omega^k) + \gamma_y^k \|Y^k - Y^{k-1}\|^2,$$

where $\eta_y = c_2$ and $\gamma_y^k = 0$. Following the same rationale that leads to Theorem 1, we obtain, as in (18),

$$\gamma_i^k \leq C_x \eta_i^{k-1}, \frac{2\alpha_2(2c_2)^2}{\beta} \leq C_y(\eta_y - \frac{\alpha_2(2c_2)^2}{\beta}),$$

where $\alpha_2 = \frac{3\alpha}{\sigma_B(1-|1-\alpha|)^2} = 3$ and $0 < C_x, C_y < 1$. In our experiments, we choose

$$\zeta_1^k = \min \left\{ \frac{a_{k-1} - 1}{a_k}, \sqrt{C_x \frac{L_W(H^{k-1})}{L_W(H^k)}} \right\}, \zeta_2^k = \min \left\{ \frac{a_{k-1} - 1}{a_k}, \sqrt{C_x \frac{L_H(W^{k+1}) + \beta}{L_H(W^k) + \beta}} \right\},$$

where $a_0 = 1$, $a_k = \frac{1}{2}(1 + \sqrt{1 + 4a_{k-1}^2})$, and $\beta \geq 4c_2 \frac{(6+3C_y)}{C_y}$.

Experiments We will compare iADMM with (i) ADMM (that is iADMM without using the inertial terms: $\zeta_1^k = \zeta_2^k = 0$), and (ii) TITAN - the inertial block majorization minimization proposed in [21] that directly solves Problem (64) and competes favorably with the state of the art on the NMF problem (see [20] which is a special case of TITAN). In our implementation, we use Lipschitz gradient surrogate for W and H and use default parameter setting for TITAN.

In the following experiments, we set the parameters c_1 and c_2 of Problem (64) to be $c_1 = 0.001$ and $c_2 = 0.01$.

In the first experiment, we generate 2 synthetic low-rank data sets X with $(n, m, r) = (500, 200, 20)$ and $(n, m, r) = (500, 500, 20)$: we generate U and V by using the MATLAB command `rand(n,r)` and `rand(r,m)` respectively, and then let $X=U*V$. For each data set, we run each algorithm with the same 30 random initial points $W_0=\text{rand}(n,r)$, $H_0=\text{rand}(r,m)$ (for iADMM and ADMM we let $Y_0=H_0$ and $\omega_0=\text{zeros}(r,m)$), and for each initial point we run each algorithm for 15 seconds. We report the evolution of the average objective function values of Problem (64) with respect to time in Figure 5 and the mean \pm std of the final objective function values in Table 2. We observe that iADMM outperforms ADMM which illustrates the acceleration effect. Among the algorithms, TITAN converges the fastest, but only slightly faster than iADMM. However, iADMM provides the best final objective function values on average.

In the second experiment, we test the algorithms on 4 image data sets CBCL⁴ (2429 images of dimension 19×19), ORL⁵ (400 images of dimension 92×112), Frey⁶ (1965 images of dimension 28×20), and Umist⁷ (565 images of dimension 92×112). For each data set, we run each algorithm with the same 20 random initial points. We run each algorithm 100 seconds for the data sets Umist and ORL and 30 seconds for the data sets CBCL and Frey. We draw the evolution of the average objective functions values with respect to time in Figure 6 and the mean \pm std of the final objective function values in Table 3.

Once again we observe that although iADMM converges slightly slower than TITAN, iADMM always produces the best objective function values among the three algorithms. On the other hand, ADMM also outperforms TITAN in term of the final objective function values. This means that, for some reason, ADMM and iADMM are able to avoid spurious local minima more effectively than TITAN.

Table 2 Mean and standard deviation of the objective function value over 30 random initializations on the synthetic data sets. The best result is highlighted in bold.

(n, m, r)	iADMM	ADMM	TITAN
(500, 200, 20)	16.443 \pm 5.015 $\times 10^{-1}$	35.873 \pm 2.299	17.751 \pm 1.092
(500, 500, 20)	34.289 \pm 4.492	135.037 \pm 6.409	35.799 \pm 1.525

⁴ <http://cbcl.mit.edu/software-datasets/heisele/facerecognition-database.html>

⁵ <https://cam-orl.co.uk/facedatabase.html>

⁶ <https://cs.nyu.edu/~roweis/data.html>

⁷ <https://cs.nyu.edu/~roweis/data.html>

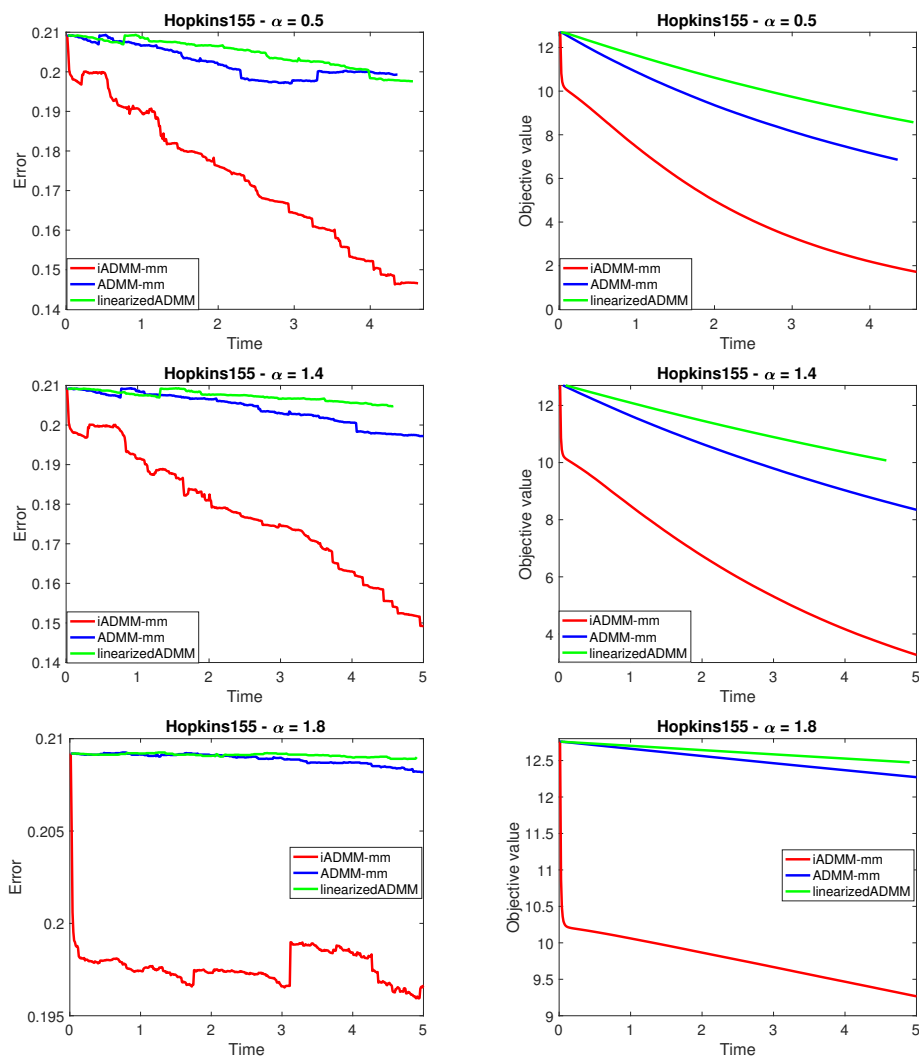


Fig. 2 Evolution of the average value of the segmentation error rate and the objective function value with respect to time on Hopkins155.

Table 3 Mean and standard deviation of the objective function value over 20 random initializations on the image data sets. The best result is highlighted in bold.

Data set	iADMM	ADMM	TITAN
CBCL	1 659.323 ± 2.514	$1\,800.626 \pm 1.156 \times 10^1$	$3\,321.104 \pm 7.271$
ORL	8 409.274 ± 7.688	$13\,825.606 \pm 1.312 \times 10^2$	$16\,844.426 \pm 1.439 \times 10^1$
Frey	1 525.242 ± 3.555	$1\,706.385 \pm 7.380$	$3\,048.246 \pm 4.737$
Umist	14 635.222 ± 1.444 × 10¹	$18\,195.557 \pm 1.059 \times 10^2$	$29\,316.019 \pm 3.433 \times 10^1$

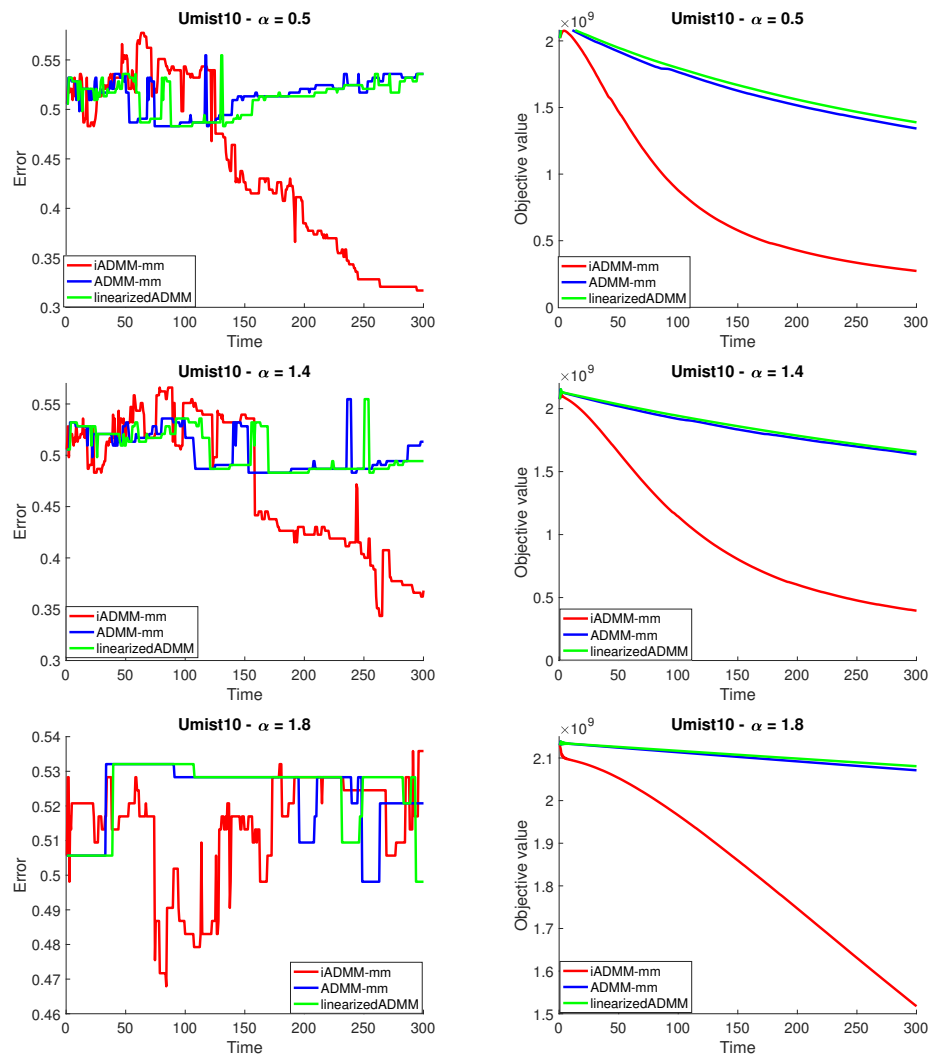


Fig. 3 Evolution of the segmentation error rate and the objective function value with respect to time on Umist10.

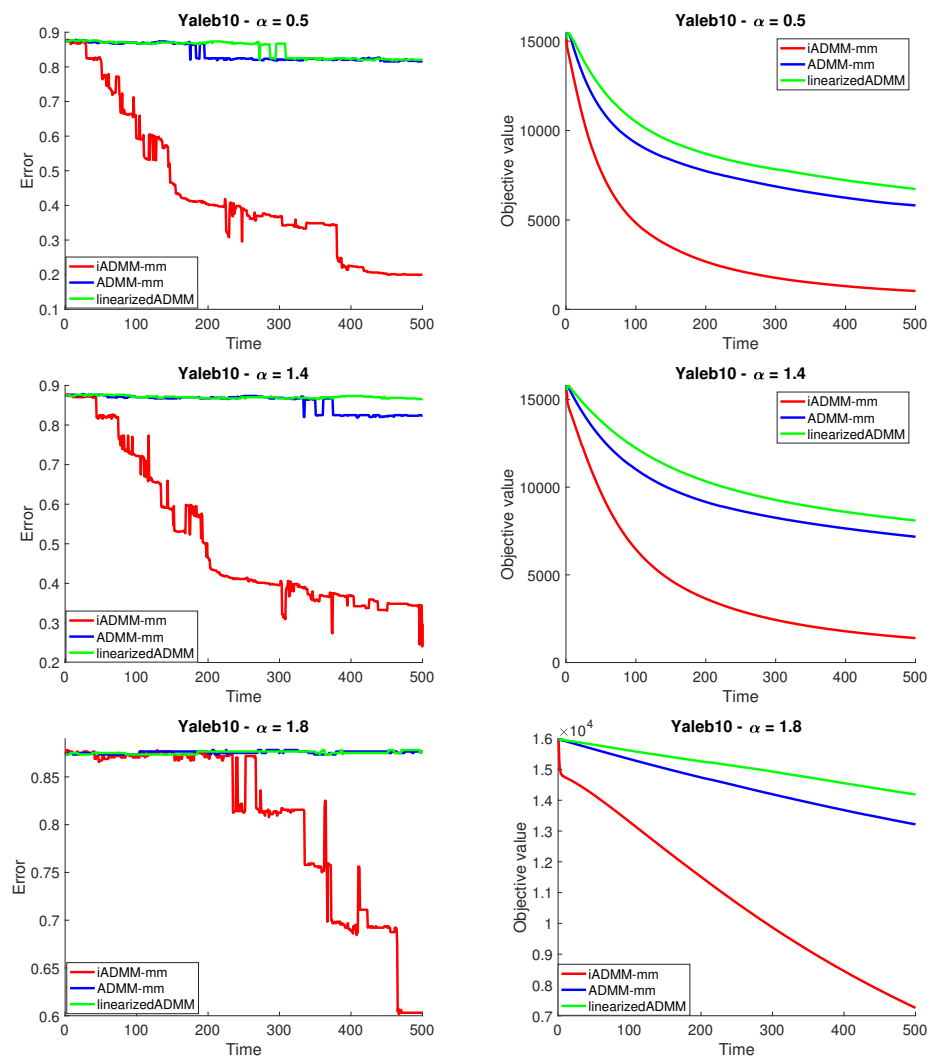


Fig. 4 Evolution of the segmentation error rate and the objective function value with respect to time on Yaleb10.

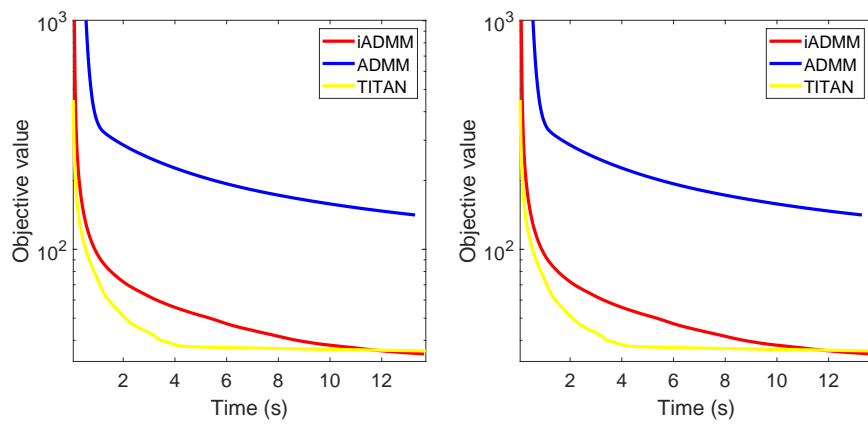


Fig. 5 Evolution of the average value of the objective function value of Problem (64) with respect to time on synthetic data sets with $(n, m, r) = (500, 200, 20)$ (left) and $(n, m, r) = (500, 500, 20)$ (right).

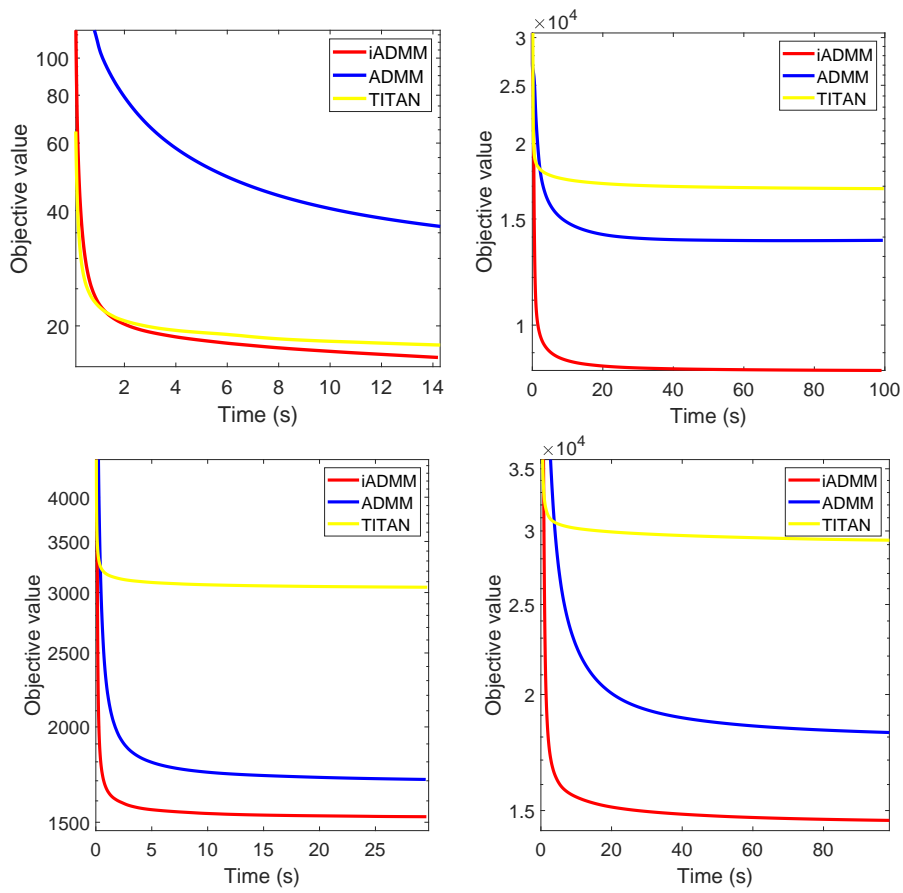


Fig. 6 Evolution of the average value of the objective function value of Problem (64) with respect to time on the image data sets CBCL (top left), ORL (top right), Frey (bottom left) and Umist (bottom right).