

# Global Convergence of ADMM in Nonconvex Nonsmooth Optimization

Yu Wang · Wotao Yin · Jinshan Zeng<sup>†</sup>

November 24, 2016

**Abstract** In this paper, we analyze the convergence of the alternating direction method of multipliers (ADMM) for minimizing a nonconvex and possibly nonsmooth objective function,  $\phi(x_0, \dots, x_p, y)$ , subject to coupled linear equality constraints. Our ADMM updates each of the primal variables  $x_0, \dots, x_p, y$ , followed by updating the dual variable. We separate the variable  $y$  from  $x_i$ 's as it has a special role in our analysis.

The developed convergence guarantee covers a variety of nonconvex functions such as piecewise linear functions,  $\ell_q$  quasi-norm, Schatten- $q$  quasi-norm ( $0 < q < 1$ ), and SCAD. It also allows nonconvex constraints such as compact manifolds (e.g., spherical, Stiefel, and Grassman manifolds) and linear complementarity constraints. Also, the  $x_0$ -block can be almost any lower semi-continuous function.

By applying our analysis, we show, for the first time, that several ADMM algorithms applied to solve nonconvex models in statistical learning, optimization on manifold, optimization over complementarity constraints, and matrix decomposition are guaranteed to converge.

Our results provide sufficient conditions for ADMM to converge on (convex or nonconvex) monotropic programs with three or more blocks, as they are special cases of our model.

ADMM has been regarded as a variant to the augmented Lagrangian method (ALM). We present a simple example to illustrate how ADMM converges but ALM diverges with bounded penalty parameter  $\beta$ . Indicated by this example and other analysis in this paper, ADMM might be a better choice than ALM for some nonconvex *nonsmooth* problems, because ADMM is not only easier to implement, it is also more likely to converge for the concerned scenarios.

---

The work of W. Yin is supported in part by NSF grants DMS-1317602 and ECCS-1462398. The work of J. Zeng is supported in part by NSF grants 61603162, 11501440.

Y. Wang

Department of Statistics, University of California, Berkeley (UCB), Berkeley, CA 94704, USA  
E-mail: wang.yu@berkeley.edu

W. Yin

Department of Mathematics, University of California, Los Angeles (UCLA), Los Angeles, CA 90025, USA  
E-mail: wotaoyin@ucla.edu

Corresponding author: J. Zeng

College of Computer Information Engineering, Jiangxi Normal University, Nanchang, Jiangxi 330022, China  
E-mail: jsh.zeng@gmail.com

---

**Keywords** ADMM, nonconvex optimization, augmented Lagrangian method, block coordinate descent, sparse optimization

## 1 Introduction

In this paper, we consider the (possibly nonconvex and nonsmooth) optimization problem:

$$\begin{aligned} & \underset{x_0, x_1, \dots, x_p, y}{\text{minimize}} && \phi(x_0, x_1, \dots, x_p, y) \\ & \text{subject to} && A_0 x_0 + A_1 x_1 + \dots + A_p x_p + B y = b, \end{aligned} \tag{1}$$

where  $\phi : \mathbb{R}^{n_0} \times \dots \times \mathbb{R}^{n_p} \times \mathbb{R}^q \rightarrow \mathbb{R} \cup \{\infty\}$  is a continuous function,  $x_i \in \mathbb{R}^{n_i}$  are variables with their coefficient matrices  $A_i \in \mathbb{R}^{m \times n_i}$ ,  $i = 0, \dots, p$ , and  $y \in \mathbb{R}^q$  is the last variable with its coefficient matrix  $B \in \mathbb{R}^{m \times q}$ . The model remains the same without  $y$  and  $B y$ ; but we keep  $y$  and  $B$  to simplify the notation.

We set  $b = 0$  throughout the paper to simplify our analysis. All of our results still hold if  $b \neq 0$  is in the image of the matrix  $B$ , i.e.,  $b \in \text{Im}(B)$ .

Besides the linear constraints in (1), any constraint on each variable  $x_0, x_1, \dots, x_p$  and  $y$  can be treated as an indicator function and included in the objective function  $\phi$ .

In spite of the success of ADMM on convex problems, the behavior of ADMM on nonconvex problems has been largely a mystery, especially when there are also nonsmooth functions and nonconvex sets in the problems. ADMM generally fails on nonconvexity problems, but it has found to not only work in some applications but often exhibit great performance! Indeed, successful examples include: matrix completion and separation [43, 45, 54, 56], asset allocation [49], tensor factorization [31], phase retrieval [50], compressive sensing [9], optimal power flow [57], direction fields correction [29], noisy color image restoration [29], image registration [6], network inference [36], and global conformal mapping [29]. In these applications, the objective function can be nonconvex, nonsmooth, or both. Examples include the piecewise linear function, the  $\ell_q$  quasi-norm for  $q \in (0, 1)$ , the Schatten- $q$  ( $0 < q < 1$ ) [52] quasi-norm  $f(X) = \sum_i \sigma_i(X)^q$  (where  $\sigma_i(X)$  denotes the  $i$ th largest singular value of  $X$ ), and the indicator function  $\iota_{\mathcal{B}}$ , where  $\mathcal{B}$  is a nonconvex set.

The success of these applications can be intriguing, since these applications are far beyond the scope of the theoretical conditions that ADMM is proved to converge. In fact, even the three-block ADMM can diverge on a simple convex problem [10]. Nonetheless, we still find that it works well in practice. This has motivated us to explore in the paper and respond to this question: when will the ADMM type algorithms converge if the objective function includes nonconvex nonsmooth functions?

We present our Algorithm 1, where  $L_\beta$  denotes the augmented Lagrangian, and show that it converges for a large class of problems. For simplicity, Algorithm 1 uses the standard ADMM subproblems, which minimize the augmented Lagrangian  $L_\beta$  with all but one variable fixed. It is possible to extend them to inexact, linearized, and/or prox-gradient subproblems as long as a few key principles (cf. §3.1) are preserved.

In this paper, under some assumptions on the objective and matrices, Algorithm 1 is proved to converge. Algorithm 1 is a generalization to the coordinate descent method. By setting  $A_0, A_1, \dots, A_p, B$  to 0, Algorithm 1 reduces to the *cyclic* coordinate descent method.

---

**Algorithm 1** Nonconvex ADMM for (1)

---

**Initialize**  $x_1^0, \dots, x_p^0, y^0, w^0$  such that  $B^T w^0 = -\nabla h(y^0)$   
**while** stopping criterion are not satisfied **do**  
  **for**  $i = 0, \dots, p$  **do**  
     $x_i^{k+1} \leftarrow \operatorname{argmin}_{x_i} L_\beta(x_{<i}^{k+1}, x_i, x_{>i}^k, y^k, w^k)$ ;  
  **end for**  
   $y^{k+1} \leftarrow \operatorname{argmin}_y L_\beta(\mathbf{x}^{k+1}, y, w^k)$ ;  
   $w^{k+1} \leftarrow w^k + \beta (\mathbf{A}\mathbf{x}^{k+1} + B y^{k+1})$ ;  
   $k \leftarrow k + 1$ ;  
**end while**  
**return**  $x_1^k, \dots, x_p^k$  and  $y^k$ .

---

### 1.1 Proposed algorithm

Denote the variable  $\mathbf{x} := [x_0; \dots; x_p] \in \mathbb{R}^n$  where  $n = \sum_{i=0}^p n_i$ . Let  $\mathbf{A} := [A_0 \ \dots \ A_p] \in \mathbb{R}^{m \times n}$  and  $\mathbf{A}\mathbf{x} := \sum_{i=0}^p A_i x_i \in \mathbb{R}^m$ . To present our algorithm, we define the augmented Lagrangian:

$$L_\beta(\mathbf{x}, y, w) := \phi(\mathbf{x}, y) + \langle w, \mathbf{A}\mathbf{x} + B y \rangle + \frac{\beta}{2} \|\mathbf{A}\mathbf{x} + B y\|^2. \quad (2)$$

The proposed Algorithm 1 extends the standard ADMM to multiple variable blocks. It also extends the *coordinate descent* algorithms to linear constraints. We let  $x_{<i} := [x_0; \dots; x_{i-1}] \in \mathbb{R}^{n_0+n_1+\dots+n_{i-1}}$  and  $x_{>i} := [x_{i+1}; \dots; x_p] \in \mathbb{R}^{n_{i+1}+\dots+n_p}$  (clearly,  $x_{<0}$  and  $x_{>p}$  are null variables, which may be used for notational ease). Subvectors  $x_{\leq i} := [x_{<i}, x_i]$  and  $x_{\geq i}$  are defined similarly. The convergence of Algorithm 1 will be given in Theorems 1 and 2.

### 1.2 Relation to the augmented Lagrangian method (ALM)

ALM is a widely-used method for solving constrained optimization models [22, 40]. It applies broadly to nonconvex nonsmooth problems. ADMM is an approximation to ALM by sequentially updating each of the primal variables.

ALM generally uses a sequence of penalty parameters  $\{\beta^k\}$ , which is nondecreasing and possibly unbounded. When  $\beta^k$  becomes large, the ALM subproblem becomes ill-conditioned. Therefore, using bounded  $\beta^k$  is practically desirable (see [12, Theorem 5.3], [3, Proposition 2.4], or [4, Chapter 7]). For general nonconvex and nonsmooth problems, it is well known that bounded  $\beta^k$  are not enough for the convergence of ALM. Proposition 1 below introduces a simple example on which ALM diverges with any bounded  $\beta^k$ . It is surprising, however, that ADMM converges in finite steps for any fixed  $\beta > 1$  on this example.

**Proposition 1** *Consider the problem*

$$\begin{aligned}
 & \underset{x, y \in \mathbb{R}}{\text{minimize}} && x^2 - y^2 \\
 & \text{subject to} && x = y, \ x \in [-1, 1].
 \end{aligned} \quad (3)$$

*It holds that*

1. *If  $\{\beta^k | k \in \mathbb{N}\}$  is bounded, ALM generates a divergent sequence;*

2. for any fixed  $\beta > 1$ , ADMM generates a convergent and finite sequence to a solution.

The proof is straightforward and included in the Appendix. ALM diverges because  $L_\beta(x, y, w)$  does not have a saddle point, and there is a non-zero duality gap. ADMM, however, is unaffected. As the proof shows, the ADMM sequence satisfies  $2y^k = -w^k, \forall k$ . By substituting  $w \equiv -2y$  into  $L_\beta(x, y, w)$ , we get a convex function in  $(x, y)$ ! Indeed,

$$\rho(x, y) := L_\beta(x, y, w)|_{w=-2y} = (x^2 - y^2) + \iota_{[-1,1]}(x) - 2y(x - y) + \frac{\beta}{2}|x - y|^2 = \frac{\beta + 2}{2}|x - y|^2 + \iota_{[-1,1]}(x),$$

where  $\iota_S$  is the indicator function of set  $S$  (that is,  $\iota_S(x) = 0$  if  $x \in S$ ; otherwise, equals infinity). It turns out that ADMM solves (3) by performing the following coordinate descent iteration to  $\rho(x, y)$ :

$$\begin{cases} x^{k+1} = \operatorname{argmin}_x \rho(x, y^k), \\ y^{k+1} = y^k - \frac{\beta}{(\beta+2)^2} \frac{d}{dy} \rho(x^{k+1}, y^k). \end{cases}$$

Our analysis for the general case will show that the primal variable  $y$  somehow “controls” the dual variable  $w$  and reduces ADMM to an iteration that is similar to coordinate descent.

### 1.3 Related literature

The original ADMM was proposed in [20,18]. For convex problems, its convergence was established firstly in [19] and its convergence rates given in [21,15,16] in different settings. When the objective function is nonconvex, the recent results [54,26,34] directly make assumptions on the iterates  $(\mathbf{x}^k, y^k, w^k)$ . Hong et al. [23] deals with the nonconvex separable objective functions for some specific  $A_i$ , which forms the sharing and consensus problem. Li and Pong [30] studied the convergence of ADMM for some special nonconvex models, where one of the matrices  $A$  and  $B$  is an identity matrix. Wang et al. [46,47] studied the convergence of the nonconvex Bregman ADMM algorithm, which includes ADMM as a special case. We review their results and compare to ours in §4 below.

### 1.4 Contribution and novelty

The main contribution of this paper is the establishment of the global convergence of Algorithm 1 under certain assumptions given in Theorems 1 and 2 below. The assumptions apply to largely many nonconvex and nonsmooth objective functions. The developed theoretical results can be extended to the case where subproblems are solved inexactly with summable errors. We also allow the primal block variables  $x_1, \dots, x_p$  to be updated in an arbitrary order as long as  $x_0$  is updated first and  $y$  is updated last (just before the  $w$ -update). The novelty of this paper can be summarized as follows:

- (1) **Weaker assumptions.** Compared to the related works [54,26,34,23,30,46,47], the convergence conditions in this paper are weaker, extending the ADMM theory to significantly more nonconvex functions and nonconvex sets. See Table 1. In addition, we allow the primal variables  $x_1, \dots, x_p$  to be updated in an arbitrary order at each iteration<sup>1</sup>, which is new in the ADMM literature. We show that most of our assumptions are necessary by providing counter examples. We also give the first example that causes ADMM to converge but ALM to diverge.

<sup>1</sup> This is the best that one hope (except for very specific problems) since [55, Section 1] shows a convex 2-block problem, which ADMM fails to converge.

- 
- (2) **New examples.** By applying our main theorems, we prove convergence for the nonconvex ADMM applied to the following problems:
- statistical regression based on nonconvex regularizer such as MCP, SCAD, and  $\ell_q$  quasi-norm;
  - minimizing smooth functions subject to norm or Stiefel/Grassmannian manifold constraints;
  - matrix decomposition using nonconvex Schatten- $q$  regularizer;
  - smooth minimization subject to complementarity constraints.
- (3) **Novel techniques.** We improve upon the existing analysis techniques and introduce new ones.
- (a) *An induction technique for nonconvex, nonsmooth case.* The analysis uses the augmented Lagrangian as the Lyapunov function: Algorithm 1 produces a sequence of points whose augmented Lagrangian function values are decreasing and lower bounded. This technique appeared first in [23] and also in [30, 46]. However, it has trouble handling nonsmooth functions. An induction technique is introduced to overcome this difficulty and extend the current framework to nonconvex, nonsmooth, multi-block cases. The technique is used in the proof of Lemma 9.
  - (b) *Restricted prox-regularity.* Most of the convergence analysis of nonconvex optimization either assumes or proves the sufficient descent and bounded subgradient properties (c.f., [1, 23]). This property is easily obtainable if the objective is smooth. However, some nonconvex and nonsmooth objectives (e.g. nonconvex  $\ell_q$  quasi-norm) violate these properties. We overcome this challenge with the introduced *restricted prox-regularity property* (Definition 2). If the objective satisfies such a property, we prove that the sequence enjoy sufficient descent and bounded subgradients after a finite number of iterations.
  - (c) *More general linear mappings.* Most nonconvex ADMM analysis is applied to the primal variables  $\mathbf{x}$  and  $y$  directly. This requires the matrices  $A_0, A_1, \dots, A_p, B$  to either identity or have full column/row rank. In this paper, we introduce techniques to work with possibly rank-deficient  $A_0, A_1, \dots, A_p, B$  (see, for example, Lemma 5). This allows us to ensure convergence of ADMM on some important applications in signal processing and statistical learning (see §5).

In addition, we use several other techniques that are tailored to relax our convergence assumptions as much as possible.

## 1.5 Organization

The remainder of this paper is organized as follows. Section 2 presents the main convergence analysis. Section 3 gives the detailed proofs. Section 4 discusses the tightness of the assumptions, the primal variable update order, and inexact minimization issues. Section 5 applies the developed theorem in some typical applications and obtains novel convergence results. Finally, section 6 concludes this paper.

## 2 Main results

### 2.1 Definitions

In these definitions,  $\partial f$  denotes the set of general subgradients of  $f$  in [41, Definition 8.3]. We call a function *Lipschitz differentiable* if it is differentiable and the gradient is Lipschitz continuous. The functions given in the next two definitions are permitted in our model.

**Definition 1 (Piecewise linear function)** A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is *piecewise linear* if there exist polyhedra  $U_1, \dots, U_K \subset \mathbb{R}^n$ , vectors  $a_1, \dots, a_K \in \mathbb{R}^n$ , and points  $b_1, \dots, b_K \in \mathbb{R}$  such that  $\bigcup_{i=1}^K \overline{U_i} = \mathbb{R}^n$ ,  $U_i \cap U_j = \emptyset$  ( $\forall i \neq j$ ), and  $f(x) = a_i^T x + b_i$  when  $x \in U_i$ ,  $i = 1, \dots, K$ .

**Table 1** Conditions for ADMM convergence (note:  $f_0, f_1, \dots, f_p$  are not required to exist)

	Scenario 1		Scenario 2
model	minimize $\phi(\mathbf{x}, y) = g(\mathbf{x}) + \sum_{i=0}^p f_i(x_i) + h(y)$ subject to $\mathbf{A}\mathbf{x} + B y = b$		minimize $\phi(\mathbf{x}, y)$ subject to $\mathbf{A}\mathbf{x} + B y = b$
$\phi$	coercive over the feasible set $\{(\mathbf{x}, y) : \mathbf{A}\mathbf{x} + B y = b\}$		
$g, h$	Lipschitz differentiable		$\phi$ Lipschitz differentiable
	Scenario 1a	Scenario 1b	
$f_0$	lower semi-continuous	$\partial f_0$ bounded in any bounded set	
$f_1, \dots, f_p$	restricted prox-regular	piecewise linear	
$\mathbf{A}, B$	$\text{Im}(\mathbf{A}) \subseteq \text{Im}(B)$		
	solution to each ADMM sub-problem is Lipschitz w.r.t. input (Assumption A3)		

**Definition 2 (Restricted prox-regularity)** Let  $M \in \mathbb{R}_+$ ,  $f : \mathbb{R}^N \rightarrow \mathbb{R} \cup \{\infty\}$ , and define the exclusion set

$$S_M := \{x \in \text{dom}(f) : \|d\| > M \text{ for all } d \in \partial f(x)\}.$$

$f$  is called *restricted prox-regular* if, for any  $M > 0$  and bounded set  $T \subseteq \text{dom} f$ , there exists  $\gamma > 0$  such that

$$f(y) + \frac{\gamma}{2}\|x - y\|^2 \geq f(x) + \langle d, y - x \rangle, \quad \forall x \in T \setminus S_M, y \in T, d \in \partial f(x), \|d\| \leq M. \quad (4)$$

(If  $T \setminus S_M$  is empty, (4) is satisfied.)

Throughout the paper,  $\|\cdot\|$  represents the Euclidean norm. Definition 2 is related to, but different from, the concepts *prox-regularity* [39], *hypomonotonicity* [41, Example 12.28] and *semi-convexity* [35, 25, 28, 37], all of which impose global conditions. Definition 2 only requires (4) to hold over a subset. Functions such as  $\ell_q$  quasi-norms ( $0 < q < 1$ ), Schatten- $q$  quasi-norms ( $0 < q < 1$ ), and indicator functions of compact smooth manifolds are examples of Definition 2 but *not* prox-regular, hypomonotone or semiconvex.

Definition 2 introduces functions that do not satisfy (4) globally *only because* they are asymptotically “steep” in the exclusion set  $S_M$ . Such functions include  $|x|^q$  ( $0 < q < 1$ ), for which  $S_M$  has the form  $(-\epsilon_M, 0) \cup (0, \epsilon_M)$ ; the Schatten- $q$  quasi-norm ( $0 < q < 1$ ), for which  $S_M = \{X : \exists i, \sigma_i(X) < \epsilon_M\}$  as well as  $\log(x)$ , for which  $S_M = (0, \epsilon_M)$ , where  $\epsilon_M$  is a constant depending on  $M$ . We only need (4) because the iterates  $x_i^k$  of Algorithm 1, for all large  $k$ , never enter the exclusion set  $S_M$ .

## 2.2 Main theorems

To ensure the boundedness of the sequence  $(\mathbf{x}^k, y^k, w^k)$ , we only need the coercivity of the objective function within the feasible set.

**A1 (coercivity)** Define the feasible set  $\mathcal{F} := \{(\mathbf{x}, y) \in \mathbb{R}^{n+a} : \mathbf{A}\mathbf{x} + B y = 0\}$ . The objective function  $\phi(\mathbf{x}, y)$  is coercive over this set, that is,  $\phi(\mathbf{x}, y) \rightarrow \infty$  if  $(\mathbf{x}, y) \in \mathcal{F}$  and  $\|(\mathbf{x}, y)\| \rightarrow \infty$ ;

If the feasible set of  $(\mathbf{x}, y)$  is bounded, then A1 holds trivially for any continuous objective function. Therefore, A1 is much weaker than assuming that the objective function is coercive over the entire space  $\mathbb{R}^{n+q}$ . Assumption A1 can be dropped if the boundedness of the sequence can be deduced from other means.

Within the proof,  $A_i x_i^k$  and  $By^k$  often appear in the first order conditions (e.g. see equations (12), (13)). In order to have a reverse control, i.e. controlling  $x_i^k, y^k$  based on  $A_i x_i^k, By^k$ , we need the following two assumptions on matrices  $A_i$  and  $B$ .

**A2 (feasibility)**  $\text{Im}(\mathbf{A}) \subseteq \text{Im}(B)$ , where  $\text{Im}(\cdot)$  returns the image of a matrix;

**A3 (Lipschitz sub-minimization paths)**

- (a) For any  $\mathbf{x}$ , there exists a Lipschitz continuous map  $H : \text{Im}(B) \rightarrow \mathbb{R}^q$  obeying  $H(u) = \text{argmin}_y \{\phi(\mathbf{x}, y) : By = u\}$ ,
- (b) For  $i = 0, \dots, p$  and any  $x_{<i}, x_{>i}$  and  $y$ , there exists a Lipschitz continuous map  $F_i : \text{Im}(A_i) \rightarrow \mathbb{R}^{n_i}$  obeying  $F_i(u) = \text{argmin}_{x_i} \{\phi(x_{<i}, x_i, x_{>i}, y) : A_i x_i = u\}$ ,  
and that the above  $F_i$  and  $H$  have a universal Lipschitz constant  $\bar{M} > 0$ .

These two assumptions allow us to control  $x_i^k, y^k$  by  $A_i x_i^k, By^k$  as in Lemma 1.

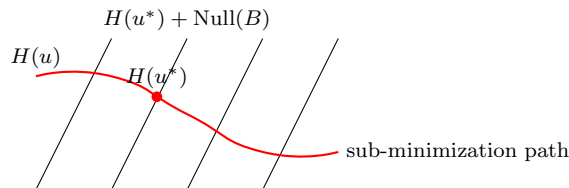
**Lemma 1** *It holds that,  $\forall k_1, k_2 \in \mathbb{N}$ ,*

$$\|y^{k_1} - y^{k_2}\| \leq \bar{M} \|By^{k_1} - By^{k_2}\|, \quad (5)$$

$$\|x_i^{k_1} - x_i^{k_2}\| \leq \bar{M} \|A_i x_i^{k_1} - A_i x_i^{k_2}\|, \quad i = 0, 1, \dots, p, \quad (6)$$

where  $\bar{M}$  is given in A3.

They weaken the full column rank assumption typically assumed for matrices  $A_i$  and  $B$ . When  $A_i$  and  $B$  have full column rank, their null spaces are trivial and, therefore,  $F_i, H$  reduce to linear operators and satisfy A3. However, assumption A3 allows non-trivial null spaces and holds for more functions. For example, if a function  $f$  is a  $C^2$  with its Hessian matrix  $H$  bounded everywhere  $\sigma_1 I \succeq H \succeq \sigma_2 I$  ( $\sigma_1 > \sigma_2 > 0$ ), then  $F$  satisfies A3 for any matrix  $A$ . Also note that we write  $H(u) = \text{argmin}\{\phi(\mathbf{x}, y) : By = u\}$  instead of  $H(u) \in \text{argmin}\{\phi(\mathbf{x}, y) : By = u\}$ , so the unique minimizer is a part of the assumption. If the uniqueness fails to hold, i.e., there exists  $y_1, y_2$  such that  $By_1 = By_2$  and  $\phi(\mathbf{x}, y_1) = \phi(\mathbf{x}, y_2)$ , then the augmented Lagrangian cannot distinguish them, causing troubles to the boundedness of the sequence.



**Fig. 1** Illustration of assumption A3, which assume that  $H(u) = \text{argmin}\{h(y) : By = u\}$  is Lipschitz [42].

As for the objective function, we consider two different scenarios:

- Theorem 1 considers the scenario where  $\mathbf{x}$  and  $y$  are decoupled in the objective function;
- Theorem 2 considers the scenario where  $\mathbf{x}$  and  $y$  are possibly coupled but their function  $\phi(\mathbf{x}, y)$  is Lipschitz differentiable.

The model in the first scenario is

$$\begin{aligned} & \underset{x_0, x_1, \dots, x_p, y}{\text{minimize}} && f(x_0, x_1, \dots, x_p) + h(y) \\ & \text{subject to} && A_0 x_0 + A_1 x_1 + \dots + A_p x_p + B y = b, \end{aligned} \tag{7}$$

where the function  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  ( $n = \sum_{i=0}^p n_i$ ) is proper, continuous, and possibly nonsmooth, and the function  $h : \mathbb{R}^q \rightarrow \mathbb{R}$  is proper and differentiable. Both  $f$  and  $h$  can be nonconvex.

**Theorem 1** *Suppose that A1-A3 and the following assumptions hold.*

A4 (**objective- $f$  regularity**)  $f$  has the form

$$f(\mathbf{x}) := g(\mathbf{x}) + \sum_{i=0}^p f_i(x_i)$$

where

(i)  $g(\mathbf{x})$  is Lipschitz differentiable with constant  $L_g$ ,

(ii) Either

a.  $f_0$  is lower semi-continuous,  $f_i(x_i)$  is restricted prox-regular (definition 2) for  $i = 1, \dots, p$ ; Or,

b.  $\sup\{\|d\| : x_0 \in S, d \in \partial f_0(x_0)\}$  is bounded for any bounded set  $S$ ,  $f_i(x_i)$  is continuous and piecewise linear (definition 1) for  $i = 1, \dots, p$ ;

A5 (**objective- $h$  regularity**)  $h(y)$  is Lipschitz differentiable with constant  $L_h$ ;

Then, Algorithm 1 converges subsequently for any sufficiently large  $\beta$  (the lower bound is given in Lemma 9), that is, starting from any  $x_2^0, \dots, x_p^0, y^0, w^0$ , it generates a sequence that is bounded, has at least one limit point, and that each limit point  $(\mathbf{x}^*, y^*, w^*)$  is a stationary point of  $L_\beta$ , namely,  $0 \in \partial L_\beta(\mathbf{x}^*, y^*, w^*)$ .

In addition, if  $L_\beta$  is a Kurdyka-Lojasiewicz (KL) function [32, 5, 1], then  $(\mathbf{x}^k, y^k, w^k)$  converges globally to the unique limit point  $(\mathbf{x}^*, y^*, w^*)$ .

Assumptions A4 and A5 regulate the objective functions. None of the functions needs to be convex.  $f_0$  can be any lower semi-continuous function, and the non-Lipschitz differentiable parts  $f_1, \dots, f_n$  of  $f$  shall satisfy either Definition 1 or Definition 2. Under Assumptions A4 and A5, the augmented Lagrangian function  $L_\beta$  is continuous.

It will be easy to see, from our proof in Section 3.3, that the Lipschitz differentiable assumption on  $g$  can be relaxed to hold just in any bounded set, since the boundedness of  $\{\mathbf{x}^k\}$  is established before that property is used in our proof. Consequently,  $g$  can be functions like  $e^x$ , whose derivative is not globally Lipschitz.

Functions satisfying the KL inequality include real analytic functions, semi-algebraic functions and locally strongly convex functions (more information can be referred to Sec. 2.2 in [53] and references therein).

In the second scenario,  $\mathbf{x}$  and  $y$  can be coupled in the objective as shown in (1), but the objective needs to be smooth.

**Theorem 2** *Suppose that A1-A3 hold and  $\phi$  in (1) is Lipschitz differentiable with constant  $L_\phi$ . Then, Algorithm 1 has the same subsequential and global convergence results as stated in Theorem 1.*

Although Theorems 1 and 2 impose different conditions on the objective functions, their proofs are similar. Hence, we will focus on proving Theorem 1 first and leave the proof of Theorem 2 to the Appendix.



### 3 Proof

#### 3.1 Keystones

The following properties hold for Algorithm 1 under our assumptions. Here, we first list them and present Proposition 2, which establishes convergence assuming these properties. Then in the next two subsections, we prove these properties.

P1 (**boundedness**)  $\{\mathbf{x}^k, y^k, w^k\}$  is bounded, and  $L_\beta(\mathbf{x}^k, y^k, w^k)$  is lower bounded.

P2 (**sufficient descent**) There is  $C_1(\beta) > 0$  such that for all sufficiently large  $k$ , we have

$$L_\beta(\mathbf{x}^k, y^k, w^k) - L_\beta(\mathbf{x}^{k+1}, y^{k+1}, w^{k+1}) \geq C_1(\beta) \left( \|B(y^{k+1} - y^k)\|^2 + \sum_{i=1}^p \|A_i(x_i^k - x_i^{k+1})\|^2 \right). \quad (8)$$

P3 (**subgradient bound**) There exists  $C_2(\beta) > 0$  and  $d^{k+1} \in \partial L_\beta(\mathbf{x}^{k+1}, y^{k+1}, w^{k+1})$  such that

$$\|d^{k+1}\| \leq C_2(\beta) \left( \|B(y^{k+1} - y^k)\| + \sum_{i=1}^p \|A_i(x_i^{k+1} - x_i^k)\| \right). \quad (9)$$

It is our intention to start  $i$  at 1, thus skipping the  $x_0$ -block, in (8) and (9).

The proposition below is standard and not new though it does not appear exactly in the literature.

**Proposition 2** *Suppose that when an algorithm is applied to the problem (7), its sequence  $(\mathbf{x}^k, y^k, w^k)$  satisfies P1–P3. Then, the sequence has at least a limit point  $(\mathbf{x}^*, y^*, w^*)$ , and any limit point  $(\mathbf{x}^*, y^*, w^*)$  is a stationary point. That is,  $0 \in \partial L_\beta(\mathbf{x}^*, y^*, w^*)$ , or equivalently,*

$$0 = \mathbf{A}\mathbf{x}^* + B y^*, \quad (10a)$$

$$0 \in \partial f(\mathbf{x}^*) + \mathbf{A}^T w^*, \quad (10b)$$

$$0 \in \partial h(y^*) + B^T w^*. \quad (10c)$$

Furthermore, the running best rates<sup>2</sup> of the sequences  $\{\|B(y^{k+1} - y^k)\|^2 + \sum_{i=1}^p \|A_i(x_i^k - x_i^{k+1})\|^2\}$  and  $\{\|d^{k+1}\|\}$  are  $o(\frac{1}{k})$  and  $o(\frac{1}{\sqrt{k}})$ , respectively. Moreover, if  $L_\beta$  is a KL function, then  $(\mathbf{x}^k, y^k, w^k)$  converges globally to the unique point  $(\mathbf{x}^*, y^*, w^*)$ .

*Proof* The proof is standard. Similar steps are found in, for example, [1, 53].

By P1, the sequence  $(\mathbf{x}^k, y^k, w^k)$  is bounded, so there exist a convergent subsequence and a limit point, denoted by  $(\mathbf{x}^{k_s}, y^{k_s}, w^{k_s})_{s \in \mathbb{N}} \rightarrow (\mathbf{x}^*, y^*, w^*)$  as  $s \rightarrow +\infty$ . By P1 and P2,  $L_\beta(\mathbf{x}^k, y^k, w^k)$  is monotonically nonincreasing and lower bounded, and therefore  $\|A_i x_i^k - A_i x_i^{k+1}\| \rightarrow 0$  and  $\|B y^k - B y^{k+1}\| \rightarrow 0$  as  $k \rightarrow \infty$ . Based on P3, there exists  $d^k \in \partial L_\beta(\mathbf{x}^k, y^k, w^k)$  such that  $\|d^k\| \rightarrow 0$ . In particular,  $\|d^{k_s}\| \rightarrow 0$  as  $s \rightarrow \infty$ . By definition of general subgradient [41, Definition 8.3], we have  $0 \in \partial L_\beta(\mathbf{x}^*, y^*, w^*)$ .

The running best rate of the sequence  $\{\|B(y^{k+1} - y^k)\|^2 + \sum_{i=1}^p \|A_i(x_i^k - x_i^{k+1})\|^2\}$  can be easily obtained via taking advantage of [17, Lemma 1.2] or [27, Theorem 3.3.1]). By (9), it is obvious that the running best rate of the sequence  $\{\|d^{k+1}\|\}$  is  $o(\frac{1}{\sqrt{k}})$ .

Similar to the proof of Theorem 2.9 in [1], we can claim the global convergence of the considered sequence  $(\mathbf{x}^k, y^k, w^k)_{k \in \mathbb{N}}$  under the KL assumption of  $L_\beta$ .  $\square$

<sup>2</sup> A nonnegative sequence  $a_k$  induces its running best sequence  $b_k = \min\{a_i : i \leq k\}$ ; therefore,  $a_k$  has running best rate of  $o(1/k)$  if  $b_k = o(1/k)$ .

In P2, the sufficient descent inequality (8) needs hold only for all large  $k$ , not all  $k$ . In our analysis, P1 gives subsequence convergence, P2 measures the augmented Lagrangian descent, and P3 bounds the subgradient by total point changes. The reader should consider P1–P3 when generalizing Algorithm 1, for example, by replacing the direct minimization subproblems to prox-gradient or inexact subproblems.

### 3.2 Preliminaries

In this section, we give some useful lemmas that will be used in the main proof. To save space, throughout this section we assume Assumptions A1–A5 hold, and let

$$(\mathbf{x}^+, y^+, w^+) := (\mathbf{x}^{k+1}, y^{k+1}, w^{k+1}). \quad (11)$$

In addition, we let  $A_{<s}x_{<s} := \sum_{i<s} A_i x_i$  and, in a similar fashion,  $A_{>s}x_{>s} := \sum_{i>s} A_i x_i$ .

**Lemma 2** *If  $\beta > \bar{M}^2 L_h$  ( $\bar{M}$  is defined in A3), all the subproblems in Algorithm 1 are well defined.*

This lemma is on its own, so we leave its proof to the appendix.

**Lemma 3 (bound dual by primal)** *Let  $\lambda_{++}(B^T B)$  be the smallest strictly-positive eigenvalue of  $B^T B$ ,  $C \triangleq L_h \bar{M} \lambda_{++}^{-1/2}(B^T B)$ . For all  $k \in \mathbb{N}$ , it holds that*

- (a)  $B^T w^k = -\nabla h(y^k)$ .
- (b)  $\|w^+ - w^k\| \leq C \|By^+ - By^k\|$ .

*Proof* Part (a) follows directly from the optimality condition of  $y^k$ :  $0 = \nabla h(y^k) + B^T w^{k-1} + B^T \beta(\mathbf{A}\mathbf{x}^k + By^k)$ , and  $w^k = w^{k-1} + \beta(\mathbf{A}\mathbf{x}^k + By^k)$ .

Then let us prove Part (b). Since  $w^+ - w^k = \beta(\mathbf{A}\mathbf{x}^+ + By^+) \in \text{Im}(B)$ , we get

$$\|w^+ - w^k\| \leq \lambda_{++}^{-1/2}(B^T B) \|B^T(w^+ - w^k)\| = \lambda_{++}^{-1/2}(B^T B) \|\nabla h(y^+) - \nabla h(y^k)\| \leq C \|By^+ - By^k\|.$$

The last inequality follows from the Lipschitz property of  $\nabla h$  and Lemma 1.  $\square$

### 3.3 Main proof

This subsection proves Theorem 1 for Algorithm 1 under Assumptions A1–A5. For all  $k \in \mathbb{N}$  and  $i = 0, \dots, p$ , because of the optimality of  $x_i^k$ , we can introduce the following *general subgradients*  $d_i^k$  and  $\bar{d}_i^k$ ,

$$\bar{d}_i^k := -(A_i^T w^+ + \beta \rho_i^k) \in \partial_i f(x_{<i}^+, x_i^+, x_{>i}^k), \quad (12)$$

$$d_i^k := -\nabla_i g(x_{<i}^+, x_i^+, x_{>i}^k) + \bar{d}_i^k \in \partial f_i(x_i^+), \quad (13)$$

where

$$\rho_i^k := A_i^T (A_{>i} x_{>i}^k - A_{>i} x_{>i}^+) + A_i^T (By^k - By^+).$$

The next two lemmas estimate the descent of  $L_\beta(\mathbf{x}, y, w)$  at each iteration.

**Lemma 4 (descent of  $L_\beta$  during  $x_i$  update)** *The iterates in Algorithm 1 satisfy*

1.  $L_\beta(x_{<i}^+, \mathbf{x}_i^k, x_{>i}^k, y^k, w^k) \geq L_\beta(x_{<i}^+, \mathbf{x}_i^+, x_{>i}^k, y^k, w^k)$ ,  $i = 0, \dots, p$ ;

2.  $L_\beta(\mathbf{x}^k, y^k, w^k) \geq L_\beta(\mathbf{x}^+, y^k, w^k)$ ;

3.  $L_\beta(\mathbf{x}^k, y^k, w^k) - L_\beta(\mathbf{x}^+, y^k, w^k) = \sum_{i=0}^p r_i$ , where

$$r_i := f(x_{<i}^+, x_i^k, x_{>i}^k) - f(x_{<i}^+, x_i^+, x_{>i}^k) - \langle \bar{d}_i^k, x_i^k - x_i^+ \rangle + \frac{\beta}{2} \|A_i x_i^k - A_i x_i^+\|^2 \geq 0, \quad (14)$$

where  $\bar{d}_i^k$  is defined in (12).

4. For  $i = 1, \dots, p$  (without the block  $i = 0$ ), if

$$f_i(x_i^k) + \frac{\gamma_i}{2} \|x_i^k - x_i^+\|^2 \geq f_i(x_i^+) + \langle d_i^k, x_i^k - x_i^+ \rangle, \quad (15)$$

holds with constant  $\gamma_i \geq 0$  (later, this condition will be shown to hold), then we have

$$r_i \geq \frac{\beta - \gamma_i \bar{M}^2 - L_g \bar{M}^2}{2} \|A_i x_i^k - A_i x_i^+\|^2, \quad (16)$$

where the constants  $L_g$  and  $\bar{M}$  are defined in Assumptions A4 and A3, respectively.

*Proof Part 1* follows directly from the minimization subproblems, which give  $x_i^+$ . **Part 2** is a result of

$$L_\beta(\mathbf{x}^k, y^k, w^k) - L_\beta(\mathbf{x}^+, y^k, w^k) = \sum_{i=1}^p (L_\beta(x_{<i}^+, x_i^k, x_{>i}^k, y^k, w^k) - L_\beta(x_{<i}^+, x_i^+, x_{>i}^k, y^k, w^k)),$$

and part 1. **Part 3:** Each term in the sum equals  $f(x_{<i}^+, x_i^k, x_{>i}^k) - f(x_{<i}^+, x_i^+, x_{>i}^k)$  plus

$$\begin{aligned} & \langle w^k, A_i x_i^k - A_i x_i^+ \rangle + \frac{\beta}{2} \|A_{<i} x_{<i}^+ + A_i x_i^k + A_{>i} x_{>i}^k + B y^k\|^2 - \frac{\beta}{2} \|A_{<i} x_{<i}^+ + A_i x_i^+ + A_{>i} x_{>i}^k + B y^k\|^2 \\ &= \langle w^k, A_i x_i^k - A_i x_i^+ \rangle + \langle \beta (A_{<i} x_{<i}^+ + A_i x_i^+ + A_{>i} x_{>i}^k + B y^k), A_i x_i^k - A_i x_i^+ \rangle + \frac{\beta}{2} \|A_i x_i^k - A_i x_i^+\|^2 \\ &= \langle A_i^T w^+ + \beta \rho_i^k, x_i^k - x_i^+ \rangle + \frac{\beta}{2} \|A_i x_i^k - A_i x_i^+\|^2 \end{aligned}$$

where the first equality follows from the cosine rule:  $\|b + c\|^2 - \|a + c\|^2 = \|b - a\|^2 + 2\langle a + c, b - a \rangle$  with  $b = A_i x_i^k$ ,  $a = A_i x_i^+$ , and  $c = A_{<i} x_{<i}^+ + A_{>i} x_{>i}^k + B y^k$ .

**Part 4.** Let  $d_i^k$  be defined in (13). From the inequalities (6) and (15), we get

$$f_i(x_i^k) - f_i(x_i^+) - \langle d_i^k, x_i^k - x_i^+ \rangle \geq -\frac{\gamma_i}{2} \|x_i^k - x_i^+\|^2 \geq -\frac{\gamma_i \bar{M}^2}{2} \|A_i x_i^k - A_i x_i^+\|^2. \quad (17)$$

By Assumption A4 part (i) and inequality (6), we also get

$$g(x_{<i}^+, x_i^k, x_i^k) - g(x_{<i}^+, x_i^+, x_{>i}^k) - \langle \nabla_i g(x_{<i}^+, x_i^+, x_{>i}^k), x_i^k - x_i^+ \rangle \geq -\frac{L_g}{2} \|x_i^k - x_i^+\|^2 \geq -\frac{L_g \bar{M}^2}{2} \|A_i x_i^k - A_i x_i^+\|^2. \quad (18)$$

Finally, rewriting the expression of  $r_i$  and applying (17) and (18) we obtain

$$\begin{aligned} r_i &= (g(x_{<i}^+, x_i^k, x_i^k) - g(x_{<i}^+, x_i^+, x_{>i}^k) - \langle \nabla_i g(x_{<i}^+, x_i^+, x_{>i}^k), x_i^k - x_i^+ \rangle) \\ &\quad + (f_i(x_i^k) - f_i(x_i^+) - \langle d_i^k, x_i^k - x_i^+ \rangle) + \frac{\beta}{2} \|A_i x_i^k - A_i x_i^+\|^2 \\ &\geq \frac{\beta - \gamma_i \bar{M}^2 - L_g \bar{M}^2}{2} \|A_i x_i^k - A_i x_i^+\|^2. \end{aligned}$$

□

The assumption (15) in part 4 is the same as (4) in Definition 2 except the latter holds for more functions due to the exclusion set  $S_M$ . In order to relax (15) to (4), we must find  $M$  and specify the exclusion set  $S_M$ . (This complicates our analysis but is necessary for nonconvex functions such as the  $\ell_q$  quasi-norm.) We will finally achieve this relaxation in Lemma 9.

**Lemma 5 (descent of  $L_\beta$  due to  $y$  and  $w$  updates)** *If  $\beta > 2(L_h\bar{M}^2 + 1 + C)$ , where  $C$  is the constant in Lemma 3 and  $L_h$  is the Lipschitz constant in Assumption A5, then for any  $k \in \mathbb{N}$*

$$L_\beta(\mathbf{x}^+, y^k, w^k) - L_\beta(\mathbf{x}^+, y^+, w^+) \geq \|By^+ - By^k\|^2. \quad (19)$$

*Proof* Because  $\beta/2 > L_h\bar{M}^2 + 1 + C$  and  $\beta^{-1} < 1/C$ , we know

$$\frac{\beta}{2} - \frac{C^2}{\beta} - \frac{L_h\bar{M}^2}{2} > L_h\bar{M}^2 + 1 + C - C - \frac{L_h\bar{M}^2}{2} > 1. \quad (20)$$

From Assumption A5 and Lemma 3 part 2, it follows

$$\begin{aligned} & L(\mathbf{x}^+, y^k, w^k) - L(\mathbf{x}^+, y^+, w^+) \\ &= h(y^k) - h(y^+) + \langle w^+, By^k - By^+ \rangle + \frac{\beta}{2} \|By^+ - By^k\|^2 - \frac{1}{\beta} \|w^+ - w^k\|^2 \end{aligned} \quad (21)$$

$$\begin{aligned} & \geq -\frac{L_h\bar{M}^2}{2} \|By^+ - By^k\|^2 + \frac{\beta}{2} \|By^+ - By^k\|^2 - \frac{C^2}{\beta} \|By^+ - By^k\|^2 \\ &= \|By^+ - By^k\|^2, \end{aligned} \quad (22)$$

The last inequality holds because of (20).  $\square$

Based on Lemma 4 and Lemma 5, we now establish the following results:

**Lemma 6 (Monotone, lower-bounded  $L_\beta$  and (P1) bounded sequence)** *If  $\beta > 2(L_h\bar{M}^2 + 1 + C)$  as in Lemma 5, then the sequence  $(\mathbf{x}^k, y^k, w^k)$  of Algorithm 1 satisfies*

1.  $L_\beta(\mathbf{x}^k, y^k, w^k) \geq L_\beta(\mathbf{x}^+, y^+, w^+)$ .
2.  $L_\beta(\mathbf{x}^k, y^k, w^k)$  is lower bounded for all  $k \in \mathbb{N}$  and converges as  $k \rightarrow \infty$ .
3.  $\{\mathbf{x}^k, y^k, w^k\}$  is bounded.

*Proof* Part 1. It is a direct result of Lemma 4 part 2, and Lemma 5.

Part 2. By Assumption A2, there exists  $y'$  such that  $\mathbf{A}\mathbf{x}^k + By' = 0$  and  $y' = H(By')$ . By A1–A2, we have

$$f(\mathbf{x}^k) + h(y') \geq \min_{\mathbf{x}, y} \{f(\mathbf{x}) + h(y) : \mathbf{A}\mathbf{x} + By = 0\} > -\infty.$$

Then we have

$$\begin{aligned} L_\beta(\mathbf{x}^k, y^k, w^k) &= f(\mathbf{x}^k) + h(y^k) + \langle B^T w^k, y^k - y' \rangle + \frac{\beta}{2} \|\mathbf{A}\mathbf{x}^k + By^k\|^2 \\ &= f(\mathbf{x}^k) + h(y^k) + \langle \nabla h(y^k), y' - y^k \rangle + \frac{\beta}{2} \|\mathbf{A}\mathbf{x}^k + By^k\|^2 \\ (\text{Lemma 1, } \nabla h \text{ is Lipschitz}) &\geq f(\mathbf{x}^k) + h(y') + \frac{\beta - L_h\bar{M}^2}{2} \|\mathbf{A}\mathbf{x}^k + By^k\|^2 \\ &> -\infty. \end{aligned}$$

Part 3. From parts 1 and 2,  $L_\beta(\mathbf{x}^k, y^k, w^k)$  is upper bounded by  $L_\beta(\mathbf{x}^0, y^0, w^0)$  and so are  $f(\mathbf{x}^k) + h(y')$  and  $\|\mathbf{A}\mathbf{x}^k + By^k\|^2$ . By Assumption A1,  $\{\mathbf{x}^k\}$  is bounded and, therefore,  $\{By^k\}$  is also bounded. By Lemma 1, we know that  $\{y^k\}$  is bounded. By Lemma 3,  $\{B^T w^k\}$  is also bounded. Similar to the proof in Lemma 3 b,  $w^k - w^0 \in \text{Im}(B)$ . Therefore, the boundedness of  $B^T w^k$  implies the boundedness of  $w^k$ .  $\square$

It is important to remark that, once  $\beta$  is larger than the threshold, the constants and bounds in Lemmas 5 and 6 only rely on the objective  $f(x) + h(y)$ , constraint  $\mathbf{A}$ ,  $B$ , and the initial point  $\mathbf{x}^0, y^0, w^0$ . They are *independent of  $\beta$* , which is essential to the proof of Lemma 9 below.

**Lemma 7 (Asymptotic regularity)**  $\lim_{k \rightarrow \infty} \|By^k - By^+\| = 0$  and  $\lim_{k \rightarrow \infty} \|w^k - w^+\| = 0$ .

*Proof* The first result follows directly from Lemmas 4, 5, and 6 (part 2), and the second result from Lemma 3 part (a) and that  $\nabla h$  is Lipschitz.  $\square$

The lemma below corresponds to Assumption A4, part ii-b.

**Lemma 8 (Boundedness for piecewise linear  $f_i$ 's)** Consider the case that  $f_i$ ,  $i = 1, \dots, p$ , are piecewise linear. There exist constants  $M^* > 0$  (independent of  $\beta$ ),  $\bar{M}$  and  $L_g$  defined in A3 and A4, respectively, for any  $\epsilon_0 > 0$ , when  $\beta > \max\{2(M^* + 1)/\epsilon_0^2, L_h \bar{M}^2 + 1 + C\}$ , there exists  $k_{\text{pl}} \in \mathbb{N}$  such that the followings hold for all  $k > k_{\text{pl}}$ :

1.  $\|A_i x_i^+ - A_i x_i^k\| < \epsilon_0$  and  $\|x_i^+ - x_i^k\| < \bar{M} \epsilon_0$ ,  $i = 1, \dots, p$ ;
2.  $\|\nabla g(\mathbf{x}^k) - \nabla g(\mathbf{x}^+)\| < p \bar{M} L_g \epsilon_0$ .

*Proof* Part 1. Since the number  $K$  of the linear pieces of  $f_i$  is finite for  $i = 1, \dots, p$ ,  $\partial f_0$  is bounded for  $x$  in any bounded set  $S$ , and  $\{\mathbf{x}^k, y^k, w^k\}$  is bounded (see Lemma 6),  $\partial_i f(x_{<i}^+, x_i^+, x_{>i}^k)$  are uniformly bounded for all  $k$  and  $i$ . Since  $\bar{d}_i^k \in \partial_i f(x_{<i}^+, x_i^+, x_{>i}^k)$  (see (12)), the first three terms of  $r_i$  (see (14)) are bounded by a universal constant  $M^*$  independent of  $\beta$ :

$$f(x_{<i}^+, x_i^k, x_{>i}^k) - f(x_{<i}^+, x_i^+, x_{>i}^k) - \langle \bar{d}_i^k, x_i^k - x_i^+ \rangle \in [-M^*, M^*].$$

Hence, as long as  $\beta > 2(M^* + 1)/\epsilon_0^2$ ,

$$\|A_i x_i^+ - A_i x_i^k\| \geq \epsilon_0 \Rightarrow r_i \geq \frac{\beta}{2} \epsilon_0^2 - M^* > 1 \quad (23)$$

$$\Rightarrow L_\beta(x_{<i}^+, \mathbf{x}_i^k, x_{>i}^k, y^k, w^k) - 1 > L_\beta(x_{<i}^+, \mathbf{x}_i^+, x_{>i}^k, y^k, w^k). \quad (24)$$

By Lemmas 4, 5, and 6, this means  $L_\beta(\mathbf{x}^k, y^k, w^k) - 1 > L_\beta(\mathbf{x}^+, y^+, w^+)$ . Since  $\{L_\beta(\mathbf{x}^k, y^k, w^k)\}$  is lower bounded,  $\|A_i x_i^+ - A_i x_i^k\| \geq \epsilon_0$  can only hold for finitely many  $k$ . Then, we get part 1, along with Lemma 1. Part 2 follows from  $\|\nabla g(\mathbf{x}^k) - \nabla g(\mathbf{x}^+)\| \leq L_g \|\mathbf{x}^k - \mathbf{x}^+\|$ , part 1 above, and Lemma 1.  $\square$

**Lemma 9 (Sufficient descent property P2)** Suppose

$$\beta > \max \left\{ 2(M + 1)/\epsilon_0^2, L_h \bar{M}^2 + 1 + C, \sum_{i=1}^p \gamma_i \bar{M}^2 + L_g \bar{M}^2 \right\},$$

where  $\gamma_i$  ( $i = 1, \dots, p$ ) and  $\epsilon_0$  are constants only depending on  $f$ ,  $M > M^*$  is a constant independent of  $\beta$ . Then, Algorithm 1 satisfies the sufficient descent property P2.

It is worth noting that the proof below would be much simpler if there are only two blocks, instead of  $p + 2$ , or if assume *prox-regular* functions  $f_i$  instead of the less restrictive *restricted prox-regular* functions.

*Proof* We will show the lower bound (16) for  $i = 1, \dots, p$ , which, along with Lemma 4 part 3 and Lemma 5, establishes the sufficient descent property P2.

We shall obtain the lower bound (16) in the backward order  $i = p, (p-1), \dots, 1$ . In light of Lemmas 4, 5, and 6, each lower bound (16) for  $r_i$  gives us  $\|A_i x_i^k - A_i x_i^+\| \rightarrow 0$  as  $k \rightarrow \infty$ . We will first show (16) for  $r_p$ . Then, after we do the same for  $r_{p-1}, \dots, r_{i+1}$ , we will get  $\|A_j x_j^k - A_j x_j^+\| \rightarrow 0$  for  $j = p, p-1, \dots, i+1$ , using which we will get the lower bound (16) for the next  $r_i$ . We must take this backward order since  $\rho_i^k$  (see (13)) includes the terms  $A_j x_j^k - A_j x_j^+$  for  $j = p, p-1, \dots, i+1$ .

Our proof for each  $i$  is divided into two cases. In Case 1,  $f_i$ 's are restricted prox-regular (cf. Definition 2), we will get (16) for  $r_i$  by validating the condition (15) in Lemma 4 part 4 for  $f_i$ . In Case 2,  $f_i$ 's are piecewise linear (cf. Definition 1), we will show that (15) holds for  $\gamma_i = 0$  for  $k \geq k_{p1}$ , and following the proof of Lemma 4 part 4, we directly get (16) with  $\gamma_i = 0$ .

**Base step**, take  $i = p$ .

*Case 1)*  $f_p$  is restricted prox-regular. At  $i = p$ , the inclusion (13) simplifies to

$$d_p^k := -(\nabla_p g(\mathbf{x}^+) + A_p^T w^+) - \beta A_p^T (By^k - By^+) \in \partial f_p(x_p^+). \quad (25)$$

By Lemma 6 part 3 and the continuity of  $\nabla g$ , there exists a constant  $M > M^*$  (independent of  $\beta$ ) such that

$$\|\nabla_p g(\mathbf{x}^+) + A_p^T w^+\| \leq M - 1.$$

By Lemma 7, there exists  $k_p \in \mathbb{N}$  such that, for  $k > k_p$ ,

$$\beta \|A_p^T (By^k - By^+)\| \leq 1.$$

Then, we apply the triangle inequality to (25) to obtain

$$\|d_p^k\| \leq \|\nabla_p g(\mathbf{x}^+) + A_p^T w^+\| + \beta \|A_p^T (By^k - By^+)\| \leq M.$$

Use this  $M$  to define  $S_M$  in Definition 2, which qualifies  $f_p$  for (4) and thus validates the assumption in Lemma 4 part 4, proving the lower bound (16) for  $r_p$ . As already argued, we get  $\lim_k \|A_p x_p^k - A_p x_p^+\| = 0$ .

*Case 2):*  $f_i$ 's are piecewise linear (cf. Definition 1). From  $\|By^k - By^+\| \rightarrow 0$  and  $\|w^k - w^+\| \rightarrow 0$  (Lemma 7) and  $\|\nabla g(\mathbf{x}^k) - \nabla g(\mathbf{x}^+)\| < p\bar{M}L_g\epsilon_0$  (Lemma 8). In light of (25),  $d_p^k \in \partial f_p(x_p^+)$ ,  $d_p^+ \in \partial f_p(x_p^{k+2})$  such that  $\|d_p^+ - d_p^k\| < 2p\bar{M}L_g\epsilon_0$  for all sufficiently large  $k$ .

Note that  $\epsilon_0 > 0$  can be *arbitrarily* small. Given  $d_p^k \in \partial f_p(x_p^+)$  and  $d_p^+ \in \partial f_p(x_p^{k+2})$ , when the following two properties both hold: (i)  $\|d_p^+ - d_p^k\| < 2p\bar{M}L_g\epsilon_0$  and (ii)  $\|x_p^+ - x_p^k\| < \bar{M}\epsilon_0$  (Lemma 8 part 1), we can conclude that  $x_p^+$  and  $x_p^k$  belongs to the same  $\bar{U}_j$ . Suppose  $x_p^+ \in \bar{U}_{j_1}$  and  $x_p^k \in \bar{U}_{j_2}$ . Because of (ii), the polyhedron  $U_{j_1}$  is adjacent to the polyhedron  $U_{j_2}$  or  $j_1 = j_2$ . If  $\bar{U}_{j_1}$  and  $\bar{U}_{j_2}$  are adjacent ( $j_1 \neq j_2$ ) and  $a_{j_1} = a_{j_2}$ , then we can concatenate  $\bar{U}_{j_1}$  and  $\bar{U}_{j_2}$  together and all the following analysis carries through. If  $\bar{U}_{j_1}$  and  $\bar{U}_{j_2}$  are adjacent ( $j_1 \neq j_2$ ) and  $a_{j_1} \neq a_{j_2}$ , then property (i) is only possible if at least one of  $x_p^+, x_p^k$  belongs to their intersection  $\bar{U}_{j_1} \cap \bar{U}_{j_2}$  so we can include both points in either  $\bar{U}_{j_1}$  or  $\bar{U}_{j_2}$ , again giving us  $j_1 = j_2$ . Since  $x_p^+, x_p^k \in \bar{U}_{j_1}$  and  $d_p^k \in \partial f_p(x_p^+)$ , from the convexity of the linear function, we have

$$f_p(x_p^k) - f_p(x_p^+) - \langle d_p^k, x_p^k - x_p^+ \rangle \geq 0,$$

which strengthens the inequality (15) for  $i = p$  with  $\gamma_p = 0$ . By following the proof for Lemma 4 part 4, we get the lower bound (16) for  $r_p$  with  $\gamma_p = 0$ . As already argued, we get  $\lim_{k \rightarrow \infty} \|A_p x_p^k - A_p x_p^+\| = 0$ .

**Inductive step**, let  $i \in \{p-1, \dots, 1\}$  and make the inductive assumption:  $\lim_{k \rightarrow \infty} \|A_j x_j^k - A_j x_j^+\| = 0$ ,  $j = p, \dots, i+1$ , which together with  $\lim_{k \rightarrow \infty} \|By^k - By^+\| = 0$  (Lemma 7) gives  $\lim_{k \rightarrow \infty} \rho_i^k = 0$  (defined in (13)).

*Case 1)*  $f_i$  is restricted prox-regular. From (13), we have

$$d_i^k = -(\nabla_i g(x_{<i}^+, x_i^+, x_{>i}^k) + A_p^T w^+) - \beta \rho_i^k \in \partial f_i(x_i^+). \quad (26)$$

Following a similar argument in the case  $i = p$  above, there exists  $k_i \in \mathbb{N}$  such that, for  $k > \max\{k_p, k_{p-1}, \dots, k_i\}$ , we have

$$\|d_i^k\| \leq \|\nabla_i g(x_{<i}^+, x_i^+, x_{>i}^k) + A_p^T w^+\| + \beta \|\rho_i^k\| \leq M.$$

Use this  $M$  to define  $S_M$  in Definition 2 for  $f_i$  and thus validates the assumption in Lemma 4 part 4 for  $f_i$ . Therefore, we get the lower bound (16) for  $r_i$  and thus  $\lim_k \|A_i x_i^k - A_i x_i^+\| = 0$ .

*Case 2):*  $f_i$ 's are piecewise linear (cf. Definition 1). The argument is the same as in the base step for case 2, except at its beginning we must use  $d_i^k$  in (26) instead of  $d_p^k$  in (25). Therefore, we skip this part.

*Finally*, by combining  $r_i \geq C_1 \|A_i x_i^k - A_i x_i^+\|^2$ , for  $i = 1, \dots, p$ , with Lemmas 4 and 5, we establish the sufficient descent property P2.

□

**Lemma 10 (Subgradient bound property P3)** *Algorithm 1 satisfies Property P3.*

*Proof* Because  $f(\mathbf{x}) = g(\mathbf{x}) + \sum_{i=1}^p f_i(x_i)$  and  $g$  is  $C^1$ , we know

$$\partial L_\beta(\mathbf{x}^+, y^+, w^+) = \left( \left\{ \frac{\partial L_\beta}{\partial x_i} \right\}_{i=1}^p, \nabla_y L_\beta, \nabla_w L_\beta \right).$$

In order to prove the lemma, we only need to show that each block of  $\partial L_\beta$  can be controlled by some constant depending on  $\beta$ . So it suffices to prove

$$\|\nabla_w L_\beta\| \leq \frac{C}{\beta} \|By^+ - By^k\|, \quad (27)$$

$$\|\nabla_y L_\beta\| \leq L_h \bar{M} \|By^+ - By^k\|, \quad (28)$$

and, for  $s = 0, \dots, p$ , there exists  $d_s \in \frac{\partial L}{\partial x_s}$  such that

$$\|d_s\| \leq (\lambda_{\max}(A_s)\beta + L_h \bar{M} + \lambda_{\max}(A_s)C) \left( \sum_{i=1}^p \|A_i x_i^+ - A_i x_i^k\| + \|By^+ - By^k\| \right). \quad (29)$$

In order to prove (27), we have  $\nabla_w L_\beta = \mathbf{A}\mathbf{x}^+ + By^+ = \frac{1}{\beta}(w^+ - w^k)$ . By Lemma 3,  $\|\nabla_w L_\beta\| \leq \frac{C}{\beta} \|By^+ - By^k\|$ . In order to prove (28), notice that  $\nabla_y L_\beta = B^T(w^+ - w^k)$  and apply Lemma 3. In order to prove (29), observe that

$$\begin{aligned} \frac{\partial L_\beta}{\partial x_s} &= \nabla_s g(x^+) + \partial f_s(x_s^+) + A_s^T w^+ + \beta A_s^T (Ax^+ + By^+) \\ &= \nabla_s g(x_{\leq s}^+, x_{> s}^k) + \partial f_s(x_s^+) + A_s^T w^k + \beta A_s^T (A_{\leq s} x_{\leq s}^+ + A_{> s} x_{> s}^k + By^k) \end{aligned} \quad (30)$$

$$+ A_s^T (w^+ - w^k) + \beta A_s^T (A_{> s} x_{> s}^+ - A_{> s} x_{> s}^k + By^+ - By^k) + \nabla_s g(x^+) - \nabla_s g(x_{\leq s}^+, x_{> s}^k). \quad (31)$$

In (30), the first order optimal condition for  $x_s^+$  yields

$$0 \in \nabla_s g(x_{\leq s}^+, x_{> s}^k) + \partial f_s(x_s^+) + A_s^T w^k + \beta A_s^T (A_{\leq s} x_{\leq s}^+ + A_{> s} x_{> s}^k + B y^k).$$

Thus,  $d_s := -\left(A_s^T(w^+ - w^k) + \beta A_s^T(A_{> s} x_{> s}^+ - A_{> s} x_{> s}^k + B y^+ - B y^k) + \nabla_s g(x^+) - \nabla_s g(x_{\leq s}^+, x_{> s}^k)\right) \in \frac{\partial L_s}{\partial x_s}$ . Denote the biggest singular value of  $A_s$  to be  $\lambda_{\max}(A_s)$ , we have

$$\|d_s\| \leq (\lambda_{\max}(A_s)\beta + L_h \bar{M} + \lambda_{\max}(A_s)C) \sum_{i=1}^p \|A_i x_i^+ - A_i x_i^k\| + \|B y^+ - B y^k\|.$$

That completes the proof.  $\square$

*Proof (of Theorem 1).*

Lemmas 5, 9, and 10 establish the properties P1–P3. Theorem 1 follows from Proposition 2.  $\square$

## 4 Discussion

### 4.1 Tightness of assumptions

In this section, we demonstrate the tightness of the assumptions in Theorem 1 and compare them with related recent works. We only focus on results that do *not* make assumptions on the iterates themselves.

Hong et al. [23] uses  $\nabla h(y^k)$  to bound  $w^k$ . This inspired our analysis. They studied ADMM for nonconvex consensus and sharing problem. Their assumptions for the sharing problem are

- (i)  $f = \sum_i f_i$ ,  $f_i$  is Lipschitz differentiable or convex.  $\text{dom}(f)$  is a closed bounded set.
- (ii)  $h$  is Lipschitz differentiable.
- (iii)  $A_i$  has full column rank,  $B$  is the identity matrix.

The boundedness of  $\text{dom}(f)$  in part (i) implies Assumption A1, (iii) implies A2 and A3, (i) implies A4, and (ii) implies A5. Our assumptions on  $f$  and the matrices  $A, B$  are much weaker.

Wang et al. [46] studies the so-called Bregman ADMM and includes the standard ADMM as an special case. By setting all the auxiliary functions in their algorithm to zero, their assumptions for the standard ADMM reduce to

- (a)  $B$  is invertible.
- (b)  $h$  is Lipschitz differentiable and lower bounded. There exists  $\beta_0 > 0$  such that  $h - \beta_0 \nabla h$  is lower bounded.
- (c)  $f = \sum_{i=0}^p f_i(x_i)$  where  $f_i$ ,  $i = 0, \dots, p$  is strongly convex.

It is easy to see that (a), (b) and (c) imply Assumptions A1 and A3, (a) implies A2, (c) implies A4 and (b) implies A5. Therefore, their assumptions are stronger than ours. We have much more relaxed conditions on  $f$ , which can have a coupled Lipschitz differentiable term with separable restricted prox-regular or piecewise linear parts. We also have a simpler assumption on the boundedness without using  $h - \nabla h$ .

Li and Pong [30] studies ADMM and its proximal version for nonconvex objectives. Their assumptions for ADMM are

- (1)  $p = 1$  and  $f$  is lower semi-continuous.
- (2)  $h \in C^2$  with bounded Hessian matrix  $c_2 I \succeq \nabla^2 h \succeq c_1 I$  where  $c_2 > c_1 > 0$ .



- (3)  $A$  is the identity matrix,  $B$  is full row rank.  
(4)  $h$  is coercive and  $f$  is lower bounded.

The assumptions (3) and (4) imply our assumption A1 and A4, (3) implies A2 and A3, and (2) implies A5. Our assumptions on  $h$  and the matrices  $A, B$  are more general.

In summary, our convergence conditions for ADMM on nonconvex problems are the most general to the best of our knowledge. It is natural to ask whether our assumptions can be further weakened. We will provide some examples to demonstrate that, while A1, A4 and A3 can probably be further weakened, A5 and A2 are essential in the convergence of nonconvex ADMM and cannot be completely dropped in general. In [10], their divergence example is

$$\underset{x_1, x_2, y}{\text{minimize}} \quad 0 \tag{32a}$$

$$\text{subject to} \quad \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} x_1 + \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} x_2 + \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix} y = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}. \tag{32b}$$

Another related example is shown in [30, Example 7].

$$\underset{x_1, x_2, y}{\text{minimize}} \quad \iota_{S_1}(x_1) + \iota_{S_2}(x_2) \tag{33a}$$

$$\text{subject to} \quad x_1 = y \tag{33b}$$

$$x_2 = y, \tag{33c}$$

where  $S_1 = \{x = (x_1, x_2) \mid x_2 = 0\}$ ,  $S_2 = \{(0, 0), (2, 1), (2, -1)\}$ . These two examples satisfy A1 and A4-A5 but fail to satisfy A2. Without A2, ADMM is generally incapable to find a feasible point at all, let alone a stationary point. Therefore, A2 is indispensable.

To see the necessity of A5 (the smoothness of  $h$ ), consider another divergence example

$$\underset{x, y}{\text{minimize}} \quad -|x| + |y| \tag{34a}$$

$$\text{subject to} \quad x = y, \quad x \in [-1, 1]. \tag{34b}$$

For any  $\beta > 0$ , with the initial point  $(x^0, y^0, w^0) = (-\frac{2}{\beta}, 0, -1)$ , we get the sequence  $(x^{2k+1}, y^{2k+1}, w^{2k+1}) = (\frac{2}{\beta}, 0, 1)$  and  $(x^{2k}, y^{2k}, w^{2k}) = (-\frac{2}{\beta}, 0, -1)$  for  $k \in \mathbb{N}$ , which diverges. This problem satisfies all the assumptions except A5, without which  $w^k$  cannot be controlled by  $y^k$  anymore. Therefore, A5 is also indispensable.

#### 4.2 Primal variables' update order in ADMM

We discuss about the update order of  $\{x_i\}_{i=0}^p$  and  $y$  in this subsection. Theorem 1 and Theorem 2 apply to the ADMM in which the primal variables  $x_0, \dots, x_p$  are sequentially updated in a fixed order. With minor changes to the proof, both theorems still hold for free update orders of  $x_1, \dots, x_p$ , possibly different between iterations, as long as  $x_0$  is always the first and  $y$  is always the last primal variable to update, just before  $w$ . For example,  $x_1, \dots, x_p$  can be randomly permuted before each iteration, which may help avoid low-quality local solutions.

In general, including the last block  $y$  in the permutation causes ADMM to diverge. A simple example is

$$\begin{aligned} & \underset{x, y \in \mathbb{R}}{\text{minimize}} && x(1 + y) \\ & \text{subject to} && x - y = 0. \end{aligned}$$

It is easy to check that, if we fix the update order to either  $x, y, w$  or  $y, x, w$  for all iterations, Algorithm 1 converges. However, if we alternate between the two update orders, we obtain (with  $\alpha := 1/\beta$ ) the diverging sequence  $(x^{2k+1}, y^{2k+1}, w^{2k+1}) = (2\alpha(\alpha - 1), -\alpha, \alpha - 1)$  and  $(x^{2k}, y^{2k}, w^{2k}) = (-\alpha, 2\alpha(\alpha - 1), -\alpha)$ . Another divergent example when primal variables' update order alternates is the following convex and nonsmooth problem:

$$\underset{x, y}{\text{minimize}} \quad 2|x - 1| + |y| \tag{35a}$$

$$\text{subject to } x = y. \tag{35b}$$

### 4.3 Inexact optimization of subproblems

Note that all subproblems in Algorithm 1 should be solved exactly. This might restrict the wide use of the algorithm in real applications. Thus, the convergence of the inexact version of Algorithm 1 is discussed here. We extend the developed convergence results to the following inexact version of Algorithm 1 under some additional assumptions. More specifically, we assume that the sequence  $\{\mathbf{x}^k, y^k, w^k\}$  generated by the inexact version of Algorithm 1 satisfies

P1' (**boundedness**)  $\{\mathbf{x}^k, y^k, w^k\}$  is bounded, and  $L_\beta(\mathbf{x}^k, y^k, w^k)$  is lower bounded;

P2' (**sufficient descent**) there is  $C_1 > 0$  such that for all sufficiently large  $k$ , we have

$$L_\beta(\mathbf{x}^k, y^k, w^k) - L_\beta(\mathbf{x}^{k+1}, y^{k+1}, w^{k+1}) \geq C_1 (\|B(y^{k+1} - y^k)\|^2 + \sum_{i=1}^p \|A_i(x_i^k - x_i^{k+1})\|^2) - \eta_k, \tag{36}$$

P3' (**subgradient bound**) and there exists  $d^{k+1} \in \partial L_\beta(\mathbf{x}^{k+1}, y^{k+1}, w^{k+1})$  such that

$$\|d^{k+1}\| \leq C_2 (\|B(y^{k+1} - y^k)\| + \sum_{i=1}^p \|A_i(x_i^{k+1} - x_i^k)\|) + \eta_k. \tag{37}$$

When  $\sum_k \eta_k < \infty$ , the convergence results in Theorem 1 still hold for this sequence. This is because Proposition 2 still holds when the error is summable. However, when a specific algorithm is applied to solve these subproblems inexactly, it might require some additional conditions, and we leave this in the future work.

## 5 Applications

In this section, we apply the developed convergence results to several well-known applications.

---

A) Statistical learning

Statistical learning models often involve two terms in the objective function. The first term is used to measure the fitting error. The second term is a regularizer to control the model complexity. Generally speaking, it can be written as

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \sum_{i=1}^p l_i(A_i x - b_i) + r(x), \quad (38)$$

where  $A_i \in \mathbb{R}^{m_i \times n}$ ,  $b_i \in \mathbb{R}^{m_i}$  and  $x \in \mathbb{R}^n$ . Examples of the fitting measure  $l_i$  include least squares, logistic functions, and other smooth functions. The regularizers can be some sparsity-inducing functions [2, 14, 56, 58, 59, 60] such as MCP, SCAD,  $\ell_q$ . Take LASSO as an example,

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|y - Ax\|^2 + \lambda \|x\|_1.$$

The first term  $\|y - Ax\|^2$  measures the difference between the linear model  $Ax$  and outcome  $y$ . The second term  $\|x\|_1$  measures the sparsity of  $x$ .

In order to solve (38) using ADMM, we reformulate it as

$$\begin{aligned} & \underset{x, \{z_i\}_{i=1}^p}{\text{minimize}} \quad r(x) + \sum_{i=1}^p l_i(A_i z_i - b_i), \\ & \text{subject to} \quad x = z_i, \quad \forall i = 1, \dots, p. \end{aligned} \quad (39)$$

Algorithm 2 gives the standard ADMM algorithm for this problem.

---

**Algorithm 2** ADMM for (39)

---

Denote  $\mathbf{z} = [z_1; z_2; \dots; z_p]$ ,  $\mathbf{w} = [w_1; w_2; \dots; w_p]$ .

**Initialize**  $x^0, \mathbf{z}^0, \mathbf{w}^0$  arbitrarily;

**while** stopping criterion are not satisfied **do**

$$x^{k+1} \leftarrow \underset{x}{\text{argmin}} \quad r(x) + \frac{\beta}{2} \sum_{i=1}^p (z_i^k + \frac{w_i^k}{\beta} - x)^2;$$

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

$$z_s^{k+1} \leftarrow \underset{z_s}{\text{argmin}} \quad l_s(A_s z_s - b_s) + \frac{\beta}{2} (z_s + \frac{w_s^k}{\beta} - x^{k+1})^2;$$

$$w_s^{k+1} = w_s^k + \beta(z_s^{k+1} - x^{k+1});$$

**end for**

$k \leftarrow k + 1;$

**end while**

return  $x^k$ .

---

Based on Theorem 1, we have the corollary.

**Corollary 1** Let  $r(x) = \|x\|_q^q = \sum_i |x_i|^q$ ,  $0 < q \leq 1$  or any piecewise linear function, if

- i) (Coercivity)  $r(x) + \sum_i l_i(A_i x + b_i)$  is coercive;
- ii) (Smoothness) For each  $i = 1, \dots, p$ ,  $l_i$  is Lipschitz differentiable.

then for sufficiently large  $\beta$ , the sequence  $(x^k, \mathbf{z}^k, \mathbf{w}^k)$  generated by Algorithm 2 has limit points and all of its limit points are stationary points of the augmented Lagrangian  $L_\beta$ .

*Proof* Rewrite the optimization to a standard form, we have

$$\begin{aligned} & \underset{x, \{z_i\}_{i=1}^p}{\text{minimize}} && r(x) + \sum_{i=1}^p l_i(A_i z_i - b_i), \end{aligned} \quad (40a)$$

$$\text{subject to} \quad E x + \mathbf{I}_{np} z = 0. \quad (40b)$$

where  $E = -[\mathbf{I}_n; \dots; \mathbf{I}_n] \in \mathbb{R}^{np \times n}$ ,  $\mathbf{I}_{np} \in \mathbb{R}^{np \times np}$  is the identity matrix, and  $z = [z_1; \dots; z_p] \in \mathbb{R}^{np}$ . Fitting (40) to the standard form (7), there are two blocks  $(x, z)$  and  $B = \mathbf{I}_{np}$ .  $f(x) = r(x)$  and  $h(z) = \sum_{i=1}^p l_i(A_i z_i - b_i)$ .

Now let us check A1–A5. A1 holds because of i). A2 holds because  $B = \mathbf{I}_{np}$ . A5 holds because of ii). A3 holds because  $E$  and  $\mathbf{I}_{np}$  both have full column ranks. Hence, it remains to verify A4 that  $r(x) = \sum_i |x_i|^q$  is restricted prox-regular. When  $q = 1$ , this is trivial so we only consider the nonconvex case  $0 < q < 1$ . The set of general subgradient of  $r(\cdot)$  is

$$\partial r(x) = \{d = [d_1; \dots; d_n] \mid d_i = q \cdot \text{sign}(x_i) |x_i|^{q-1} \text{ if } x_i \neq 0; d_i \in \mathbb{R} \text{ if } x_i = 0\}.$$

For any two positive constants  $C > 0$  and  $M > 1$ , take  $\gamma = \max(\frac{4(pC^q + MC)}{c^2}, q(1-q)c^{q-2})$ , where  $c \triangleq \frac{1}{3} M^{\frac{1}{1-q}}$ . The exclusion set  $S_M$  contains the set  $\{x \mid \min_{x_i \neq 0} |x_i| \leq 3c\}$ . For any point  $z \in \mathbb{B}(0, C)/S_M$  and  $y \in \mathbb{B}(0, C)$ , if  $\|z - y\| \leq c$ , then  $\text{supp}(z) \subset \text{supp}(y)$  and  $\|z\|_0 \leq \|y\|_0$ , where  $\mathbb{B}(0, C) \triangleq \{x \mid \|x\| < C\}$ ,  $\text{supp}(z)$  denotes the index set of all non-zero elements of  $z$  and  $\|z\|_0$  denotes the cardinality of  $\text{supp}(z)$ . Define

$$y'_i = \begin{cases} y_i & i \in \text{supp}(z) \\ 0 & i \notin \text{supp}(z) \end{cases}, \quad i = 1, \dots, p.$$

Then the following line of proof holds,

$$\begin{aligned} \|y\|_q^q - \|z\|_q^q - \langle d, y - z \rangle &\stackrel{(a)}{\geq} \|y'\|_q^q - \|z\|_q^q - \langle d, y' - z \rangle \\ &\stackrel{(b)}{\geq} -\frac{q(1-q)}{2} c^{q-2} \|z - y'\|^2 \\ &\stackrel{(c)}{\geq} -\frac{q(1-q)}{2} c^{q-2} \|z - y\|^2. \end{aligned} \quad (41)$$

(a) holds because for any  $i \notin \text{supp}(z)$ ,  $|y_i| < c$ , which means  $|y_i|^q \geq M y_i$ . (b) holds because  $r(x)$  is twice differentiable along the line segment connecting  $z$  and  $y'$ , and the second order derivative is no bigger than  $q(1-q)c^{q-2}$ . (c) holds because  $\|z - y\| \geq \|z - y'\|$ .

If  $\|z_1 - z_2\| > c$ , then we have

$$\|z_1\|_q^q - \|z_2\|_q^q - \langle d, z_1 - z_2 \rangle \geq -(2pC^q + 2MC) \geq -\frac{2pC^q + 2MC}{c^2} \|z_1 - z_2\|^2. \quad (42)$$

Combining (41) and (42) yields the result. This verifies A4 and completes the proof.  $\square$

---

B) Minimization on compact manifold

Compact manifolds and their projection operators such as spherical manifolds  $S^{n-1}$ , Stiefel manifolds (the set of  $p$  orthonormal vectors  $x_1, \dots, x_p \in \mathbb{R}^n$ ,  $p \leq n$ ) and Grassmann manifolds (the set of subspaces in  $\mathbb{R}^n$  of dimension  $p$ ) often arise in optimization. Some recent studies and algorithms can be found in [51, 29, 33]. A simple example is:

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && J(x), \\ & \text{subject to} && \|x\|^2 = 1, \end{aligned} \tag{43}$$

More generally, let  $S$  be any compact set. We consider the problem

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && J(x), \\ & \text{subject to} && x \in S, \end{aligned} \tag{44}$$

which can be rewritten to the following form:

$$\begin{aligned} & \underset{x, y}{\text{minimize}} && \iota_S(x) + J(y), \\ & \text{subject to} && x - y = 0, \end{aligned} \tag{45a}$$

$$\tag{45b}$$

where  $\iota_S(\cdot)$  is the indicator function:  $\iota_S(x) = 0$  if  $x \in S$  or  $\infty$  if  $x \notin S$ . Applying ADMM to solve this problem, we get Algorithm 3.

---

**Algorithm 3** ADMM for minimization on a compact set (45)

---

```

Initialize  $x^0, y^0, w^0$  arbitrarily;
while stopping criterion are not satisfied do
   $x^{k+1} \leftarrow \text{Proj}_S(y^k - \frac{w^k}{\beta});$ 
   $y^{k+1} \leftarrow \text{argmin}_y J(y) + \frac{\beta}{2} \|y - \frac{w^k}{\beta} - x^{k+1}\|^2;$ 
   $w^{k+1} \leftarrow w^k + \beta(y^{k+1} - x^{k+1});$ 
   $k \leftarrow k + 1.$ 
end while
return  $x^k.$ 

```

---

Based on Theorem 1, we have the following corollary.

**Corollary 2** *If  $J$  is Lipschitz differentiable, then for any sufficiently large  $\beta$ , the sequence  $(x^k, y^k, w^k)$  generated by Algorithm 3 has at least one limit point, and each limit point is a stationary point of the augmented Lagrangian  $L_\beta$ .*

*Proof* To show this corollary, we shall verify Assumptions A1–A5.

Assumption A1 holds because the feasible set is a bounded set and  $J$  is lower bounded on the feasible set. A2 and A3 hold because both  $A$  and  $B$  are identity matrices. A5 holds because  $J$  is Lipschitz differentiable. A4 holds because  $\iota_S$  is lower semi-continuous.

---

C) Smooth optimization over complementarity constraints

We consider the following optimization problem over complementarity constraints.

$$\begin{aligned} & \underset{\{x,y\}}{\text{minimize}} && h(x,y) \\ & \text{subject to} && x^T y = 0, x \geq 0, y \geq 0, \end{aligned} \tag{46}$$

where  $h(x, y)$  is a smooth function with Lipschitz differentiable gradient. The considered problem is a special case of the mathematical programming with equilibrium constraints (MPEC) [11], and includes the linear complementarity problem (LCP) [13] as a special case. In order to apply the ADMM algorithm to solve this problem, we introduce two auxiliary variables  $x', y' \in \mathbb{R}^n$  and define the complementarity set  $S \triangleq \{(x, y) : x^T y = 0, x \geq 0, y \geq 0\}$ . With these notations, problem (46) can be reformulated as follows

$$\begin{aligned} & \underset{\{x',y',x,y\}}{\text{minimize}} && \iota_S(x', y') + h(x, y) \\ & \text{subject to} && x' - x = 0, \quad y' - y = 0, \end{aligned} \tag{47}$$

where  $\iota_S(x', y')$  denotes the indicator function of the set  $S$ . Furthermore, let  $\mathbf{x}_0 = \begin{pmatrix} x' \\ y' \end{pmatrix}$  and the second block  $\mathbf{y} = \begin{pmatrix} x \\ y \end{pmatrix}$ . Then (47) becomes

$$\begin{aligned} & \underset{\mathbf{x}_0, \mathbf{y}}{\text{minimize}} && \iota_S(\mathbf{x}_0) + h(\mathbf{y}) \\ & \text{subject to} && \mathbf{x}_0 - \mathbf{y} = 0. \end{aligned} \tag{48}$$

**Corollary 3** *Assume that  $h$  is Lipschitz differentiable and coercive over the complementarity set, then for sufficiently large  $\beta$ , the sequence  $(\mathbf{x}_0^k, \mathbf{y}^k, w^k)$  generated by Algorithm ADMM applied to (48) has limit points and all of its limit points are stationary points of the augmented Lagrangian  $L_\beta$ .*

*Proof* In order to prove this corollary, we only need to verify Assumptions **A1-A5**. **A1** holds for the coercivity of  $h$  over  $S$  and the specific form of  $\iota_S$ . **A2** is obvious due to in this case  $A = I$  and  $B = -I$ . **A3** holds for both  $I$  and  $-I$  being full column rank. **A4** can be satisfied by setting  $f_0 = \iota_S$  and  $g \equiv h$ . **A5** holds due to the Lipschitz differentiability of  $h$ . Thus, according to Theorem 1, we complete the proof.  $\square$

## D) Matrix decomposition

ADMM has also been applied to solve matrix related problems, such as sparse principle component analysis (PCA) [24], matrix decomposition [44, 48], matrix completion [7], matrix recovery [38], non-negative matrix factorization [54, 45] and background/foreground extraction [8, 56].

In the following, we take the video surveillance image-flow problem as an example. A video can be formulated as a matrix  $V$  where each column is a vectorized image of a video frame. It can be generally decomposed into three parts, background, foreground, and noise. The background has low rank since it does not move. The derivative of the foreground is small because foreground (such as human beings, other moving

objectives) moves relatively slowly. The noise is generally assumed to be Gaussian and thus can be modeled via Frobenius norm.

More specifically, consider the following matrix decomposition model:

$$\underset{X, Y, Z}{\text{minimize}} \quad p(X) + \sum_{i=1}^{m-1} \|Y_i - Y_{i+1}\| + \|Z\|_F^2, \quad (49)$$

$$\text{subject to} \quad V = X + Y + Z, \quad (50)$$

where  $X, Y, Z, V \in \mathbb{R}^{n \times m}$ ,  $Y_i$  is the  $i$ th column of  $Y$ ,  $\|\cdot\|_F$  is the Frobenius norm, and  $p(X)$  is any lower bounded lower semi-continuous penalty function, for example, the Schatten- $q$  quasi-norm  $\|X\|_q$  ( $0 < q \leq 1$ ):

$$\|A\|_q = \sum_{i=1}^n \sigma_i^q,$$

where  $\sigma_i$  is the  $i$ th largest singular value of  $A$ .

The corresponding ADMM algorithm is given in Algorithm 4.

---

#### Algorithm 4 ADMM for (49)

---

**Initialize**  $Y^0, Z^0, W^0$  arbitrarily;

**while** stopping criterion are not satisfied **do**

$$X^{k+1} \leftarrow \underset{X}{\text{argmin}} p(X) + \frac{\beta}{2} \|X + Y^k + Z^k - V + W^k / \beta\|_F^2;$$

$$Y^{k+1} \leftarrow \underset{Y}{\text{argmin}} \sum_{i=1}^m \|Y_i - Y_j\| + \frac{\beta}{2} \|X^{k+1} + Y + Z^k - V + W^k / \beta\|_F^2;$$

$$Z^{k+1} \leftarrow \underset{Z}{\text{argmin}} \|Z\|_F^2 + \frac{\beta}{2} \|X^{k+1} + Y^{k+1} + Z - V + W^k / \beta\|_F^2;$$

$$W^{k+1} \leftarrow W^k + \beta(X^{k+1} + Y^{k+1} + Z^{k+1} - V);$$

$$k \leftarrow k + 1;$$

**end while**

return  $X^k, Y^k, Z^k$ .

---

**Corollary 4** For a sufficiently large  $\beta$ , the sequence  $(X^k, Y^k, Z^k, W^k)$  generated by Algorithm 4 has at least one limit point, and each limit point is a stationary point of the augmented Lagrangian function  $L_\beta$ .

*Proof* Let us verify Assumptions A1–A5. Assumption A1 holds because of the coercivity of  $\|\cdot\|_F$  and  $\|\cdot\|_q$ . A2 and A3 hold because all the coefficient matrices are identity matrices. A5 holds because  $\|\cdot\|_F^2$  is Lipschitz differentiable. A4 holds because  $p$  is lower semi-continuous.

## 6 Conclusion

This paper studied the convergence of ADMM, in its multi-block and original cyclic update form, for nonconvex and nonsmooth optimization. The objective can be certain nonconvex and nonsmooth functions while the constraints are coupled linear equalities. Our results theoretically demonstrate that ADMM, as a variant of ALM, may converge under weaker conditions than ALM. While ALM generally requires the objective function to be smooth, ADMM only requires it to have a smooth part  $h(y)$  while the remaining part  $f(\mathbf{x})$  can be coupled, nonconvex, and include separable nonsmooth functions and indicator functions of constraint sets.

Our results relax the previous assumptions (e.g., semi-convexity) and allow the nonconvex functions such as  $\ell_q$  quasi-norm ( $0 < q < 1$ ), Schatten- $q$  quasi-norm, SCAD, and others that often appear in sparse optimization. They also allow nonconvex constraint sets such as unit spheres, matrix manifolds, and complementarity constraints.

The underlying proof technique identifies an exclusion set where the sequence does not enter after finitely many iterations. We also manage to have a very general first block  $x_0$ . We show that while the middle  $p$  blocks  $x_1, \dots, x_p$  can be updated in an arbitrary order for different iterations, the first block  $x_0$  should be updated at first and the last block  $y$  at last; otherwise, the concerned iterates may diverge according to the existing example.

Our results can be applied to problems in matrix decomposition, sparse recovery, machine learning, and optimization on compact smooth manifolds and lead to novel convergence guarantees.

## Acknowledgements

We would like to thank Drs. Wei Shi, Ting Kei Pong, and Qing Ling for their insightful comments, and Drs. Xin Liu and Yangyang Xu for helpful discussions.

## Appendix

*Proof (Lemma 1)* By the definitions of  $H$  in A3(a) and  $y^k$ , we have  $y^k = H(By^k)$ . Therefore,  $\|y^{k_1} - y^{k_2}\| = \|H(By^{k_1}) - H(By^{k_2})\| \leq \bar{M}\|By^{k_1} - By^{k_2}\|$ . Similarly, by the optimality of  $x_i^k$ , we have  $x_i^k = F_i(A_i x_i^k)$ . Therefore,  $\|x_i^{k_1} - x_i^{k_2}\| = \|F_i(A_i x_i^{k_1}) - F_i(A_i x_i^{k_2})\| \leq \bar{M}\|A_i x_i^{k_1} - A_i x_i^{k_2}\|$ .  $\square$

*Proof (Lemma 2)* Let us first show that the  $y$ -subproblem is well defined. To begin with, we will show that  $h(y)$  is lower bounded by a quadratic function of  $By$ :

$$h(y) \geq h(H(0)) - (\bar{M}\|\nabla h(H(0))\|) \cdot \|By\| - \frac{L_h \bar{M}^2}{2} \|By\|^2.$$

By A3, we know  $h(y)$  is lower bounded by  $h(H(By))$ :

$$h(y) \geq h(H(By)).$$

Because of A5 and A3,  $h(H(By))$  is lower bounded by a quadratic function of  $By$ :

$$h(H(By)) \geq h(H(0)) + \langle \nabla h(H(0)), H(By) - H(0) \rangle - \frac{L_h}{2} \|H(By) - H(0)\|^2 \quad (51)$$

$$\geq h(H(0)) - \|\nabla h(H(0))\| \cdot \bar{M} \cdot \|By\| - \frac{L_h \bar{M}^2}{2} \|By\|^2 \quad (52)$$

Therefore  $h(y)$  is also bounded by the quadratic function:

$$h(y) \geq h(H(0)) - \|\nabla h(H(0))\| \cdot \bar{M} \cdot \|By\| - \frac{L_h \bar{M}^2}{2} \|By\|^2.$$



Recall that  $y$ -subproblem is to minimize the Lagrangian function w.r.t.  $y$ , by neglecting other constants, it is equivalent to minimize:

$$\operatorname{argmin} P(y) := h(y) + \langle w^k + \beta \mathbf{A}x^+, By \rangle + \frac{\beta}{2} \|By\|^2. \quad (53)$$

Because  $h(y)$  is lower bounded by  $-\frac{L_h \bar{M}^2}{2} \|By\|^2$ , when  $\beta > L_h \bar{M}$ ,  $P(y) \rightarrow \infty$  as  $\|By\| \rightarrow \infty$ . This shows that  $y$ -subproblem is coercive with respect to  $By$ . Because  $P(y)$  is lower semi-continuous and  $\operatorname{argmin} h(y)$  s.t.  $By = u$  has a unique solution for each  $u$ , the minimal point of  $P(y)$  must exist and the  $y$ -subproblem is well defined.

As for the  $x_i$ -subproblem,  $i = 0, \dots, p$ , ignoring the constants yields

$$\operatorname{argmin} L_\beta(x_{<i}^+, x_i, x_{>i}^k, y^k, w^k) = \operatorname{argmin} f(x_{<i}^+, x_i, x_{>i}^k) + \frac{\beta}{2} \left\| \frac{1}{\beta} w^k + A_{<i} x_{<i}^+ + A_{>i} x_{>i}^k + A_i x_i + By^k \right\|^2 \quad (54)$$

$$= \operatorname{argmin} f(x_{<i}^+, x_i, x_{>i}^k) + h(u) - h(u) + \frac{\beta}{2} \|Bu - By^k - \frac{1}{\beta} w^k\|^2. \quad (55)$$

where  $u = H(-A_{<i} x_{<i}^+ - A_{>i} x_{>i}^k - A_i x_i)$ . The first two terms are coercive bounded because  $A_{<i} x_{<i}^+ + A_{>i} x_{>i}^k + A_i x_i + Bu = 0$  and A1. The third and fourth terms are lower bounded because  $h$  is Lipschitz differentiable. Because the function is lower semi-continuous, all the subproblems are well defined.  $\square$

*Proof (Proposition 1)* Define the augmented Lagrangian function to be

$$L_\beta(x, y, w) = x^2 - y^2 + w(x - y) + \beta \|x - y\|^2.$$

It is clear that when  $\beta = 0$ ,  $L_\beta$  is not lower bounded for any  $w$ . We are going to show that for any  $\beta > 1$ , the duality gap is not zero.

$$\inf_{x \in [-1, 1], y \in \mathbb{R}} \sup_{w \in \mathbb{R}} L_\beta(x, y, w) > \sup_{w \in \mathbb{R}} \inf_{x \in [-1, 1], y \in \mathbb{R}} L_\beta(x, y, w).$$

On one hand, because  $\sup_{w \in \mathbb{R}} L_\beta(x, y, w) = +\infty$  when  $x \neq y$  and  $\sup_{w \in \mathbb{R}} L_\beta(x, y, w) = 0$  when  $x = y$ , we have

$$\inf_{x \in [-1, 1], y \in \mathbb{R}} \sup_{w \in \mathbb{R}} L_\beta(x, y, w) = 0.$$

On the other hand, let  $t = x - y$ ,

$$\sup_{w \in \mathbb{R}} \inf_{x \in [-1, 1], y \in \mathbb{R}} L_\beta(x, y, w) = \sup_{w \in \mathbb{R}} \inf_{x \in [-1, 1], t \in \mathbb{R}} t(2x - t) + wt + \beta t^2 = \sup_{w \in \mathbb{R}} \inf_{x \in [-1, 1], t \in \mathbb{R}} (w + 2x)t + (\beta - 1)t^2 \quad (56)$$

$$= \sup_{w \in \mathbb{R}} \inf_{x \in [-1, 1]} -\frac{(w + 2x)^2}{4(\beta - 1)} = \sup_{w \in \mathbb{R}} -\frac{\max((w - 2)^2, (w + 2)^2)}{4(\beta - 1)} = -\frac{1}{\beta - 1}. \quad (57)$$

This shows the duality gap is not zero (but it goes to 0 as  $\beta$  tends to  $\infty$ ).

Then let us show that ALM does not converge if  $\beta^k$  is bounded, i.e., there exists  $\beta > 0$  such that  $\beta^k \leq \beta$  for any  $k \in \mathbb{N}$ . Without loss of generality, we assume that  $\beta^k$  equals to the constant  $\beta$  for all  $k \in \mathbb{N}$ . This will not affect the proof. ALM consists of two steps

- 1)  $(x^{k+1}, y^{k+1}) = \operatorname{argmin}_{x, y} L_\beta(x, y, w^k)$ ,
- 2)  $w^{k+1} = w^k + \tau(x^{k+1} - y^{k+1})$ .

Since  $(x^{k+1} - y^{k+1}) \in \partial\psi(w^k)$  where  $\psi(w) = \inf_{x,y} L_\beta(x, y, w)$ , and we already know

$$\inf_{x,y} L_\beta(x, y, w) = -\frac{\max((w-2)^2, (w+2)^2)}{4(\beta-1)},$$

we have

$$w^{k+1} = \begin{cases} (1 - \frac{\tau}{2(\beta-1)})w^k - \frac{\tau}{\beta-1} & \text{if } w^k \geq 0 \\ (1 - \frac{\tau}{2(\beta-1)})w^k + \frac{\tau}{\beta-1} & \text{if } w^k \leq 0 \end{cases}.$$

Note that when  $w^k = 0$ , the optimization problem  $\inf_{x,y} L(x, y, 0)$  has two distinct minimal points which lead to two different values. This shows no matter how small  $\tau$  is,  $w^k$  will oscillate around 0 and never converge.

However, although the duality gap is not zero, ADMM still converges in this case. There are two ways to prove it. The first way is to check all the conditions in Theorem 1. Another way is to check the iterates directly. The ADMM iterates are

$$x^{k+1} = \max(-1, \min(1, \frac{\beta}{\beta+1}(y^k - \frac{w^k}{2\beta}))), \quad y^{k+1} = \frac{\beta}{\beta-1}(x^{k+1} + \frac{w^k}{2\beta}), \quad w^{k+1} = w^k + 2\beta(x^{k+1} - y^{k+1}). \quad (58)$$

The second equality shows that  $w^k = -2y^k$ , substituting it into the first and second equalities, we have

$$x^{k+1} = \max(-1, \min(1, y^k)), \quad y^{k+1} = \frac{1}{\beta-1}(\beta x^{k+1} - y^k). \quad (59)$$

Here  $|y^{k+1}| \leq \frac{\beta}{\beta-1} + \frac{1}{\beta-1}|y^k|$ . Thus after finite iterations,  $|y^k| \leq 2$  (assume  $\beta > 2$ ). If  $|y^k| \leq 1$ , the ADMM sequence converges obviously. If  $|y^k| > 1$ , without loss of generality, we could assume  $2 > y^k > 1$ . Then  $x^{k+1} = 1$ . It means  $0 < y^{k+1} < 1$ , so the ADMM sequence converges. Thus, we know for any initial point  $y^0$ , ADMM converges.

*Proof (Theorem 2)* Similar to the proof of Theorem 1, we only need to verify P1-P3 in Proposition 2. *Proof of P2:* Similar to Lemma 4 and Lemma 5, we have

$$\begin{aligned} & L_\beta(\mathbf{x}^k, y^k, w^k) - L_\beta(\mathbf{x}^{k+1}, y^{k+1}, w^{k+1}) \\ & \geq -\frac{1}{\beta} \|w^k - w^{k+1}\|^2 + \sum_{i=1}^p \frac{\beta - L_\phi \bar{M}}{2} \|A_i x_i^k - A_i x_i^{k+1}\|^2 + \frac{\beta - L_\phi \bar{M}}{2} \|B y^k - B y^{k+1}\|^2. \end{aligned} \quad (60)$$

Since  $B^T w^k = -\partial_y \phi(\mathbf{x}^k, y^k)$  for any  $k \in \mathbb{N}$ , we have

$$\|w^k - w^{k+1}\| \leq C_1 L_\phi \left( \sum_{i=1}^p \|x_i^k - x_i^{k+1}\| + \|y^k - y^{k+1}\| \right),$$

where  $C_1 = \min_{\lambda_i \neq 0} \lambda_i (B^T B)^{-1/2}$ ,  $\lambda_i (B^T B)$  is  $i$ th eigenvalue of  $B^T B$ , and  $L_\phi$  is the Lipschitz constant for  $\phi$ . Therefore, we have

$$\begin{aligned} & L_\beta(\mathbf{x}^k, y^k, w^k) - L_\beta(\mathbf{x}^{k+1}, y^{k+1}, w^{k+1}) \\ & \geq \left( \frac{\beta - L_\phi \bar{M}}{2} - \frac{C_1 L_\phi \bar{M}}{\beta} \right) \sum_{i=1}^p \|A_i x_i^k - A_i x_i^{k+1}\|^2 + \left( \frac{\beta - L_\phi \bar{M}}{2} - \frac{C_1 L_\phi \bar{M}}{\beta} \right) \|B y^k - B y^{k+1}\|^2. \end{aligned} \quad (61)$$

When  $\beta > \max(1, L_\phi \bar{M} + 2C_1 L_\phi \bar{M})$ , P2 holds.

*Proof of P1:* First of all, we have already shown  $L_\beta(\mathbf{x}^k, y^k, w^k) \geq L_\beta(\mathbf{x}^{k+1}, y^{k+1}, w^{k+1})$ , which means  $L_\beta(\mathbf{x}^k, y^k, w^k)$  decreases monotonically. There exists  $y'$  such that  $\mathbf{A}\mathbf{x}^k + By' = 0$  and  $y' = H(By')$ . In order to show  $L_\beta(\mathbf{x}^k, y^k, w^k)$  is lower bounded, we apply A1-A3 to get

$$\begin{aligned} L_\beta(\mathbf{x}^k, y^k, w^k) &= \phi(\mathbf{x}^k, y^k) + \langle w^k, \sum_{i=1}^p A_i x_i^k + By^k \rangle + \frac{\beta}{2} \left\| \sum_{i=1}^p A_i x_i^k + By^k \right\|^2 \\ &= \phi(\mathbf{x}^k, y^k) + \langle \partial_y \phi(\mathbf{x}^k, y^k), y' - y^k \rangle + \frac{\beta}{2} \|By^k - By'\|^2 \geq \phi(\mathbf{x}^k, y') + \frac{\beta}{4} \left\| \sum_{i=1}^p A_i x_i^k + By^k \right\|^2 > -\infty. \end{aligned} \quad (62)$$

This shows that  $L(\mathbf{x}^k, y^k, w^k)$  is lower bounded. If we view (62) from the opposite direction, it can be observed that

$$\phi(\mathbf{x}^k, y') + \frac{\beta}{4} \left\| \sum_{i=1}^p A_i x_i^k + By^k \right\|^2$$

is upper bounded by  $L_\beta(\mathbf{x}^0, y^0, w^0)$ . Then A1 ensures that  $\{\mathbf{x}^k, y^k\}$  is bounded. Therefore,  $w^k$  is bounded too.

*Proof of P3:* This part is trivial as  $\phi$  is Lipschitz differentiable. Hence we omit it.

## References

1. Attouch, H., Bolte, J., Svaiter, B.: Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized Gauss-Seidel methods. *Mathematical Programming* **137**(1-2), 91–129 (2013) 1.4, 1, 3.1
2. Bach, F., Jenatton, R., Mairal, J., Obozinski, G.: Optimization with sparsity-inducing penalties. *Foundations and Trends in Machine Learning* **4**(1), 1–106 (2012) 5
3. Bertsekas, D.P.: *Constrained optimization and Lagrange multiplier methods*. Academic press (2014) 1.2
4. Birgin, E.G., Martínez, J.M.: *Practical augmented Lagrangian methods for constrained optimization*, vol. 10. SIAM (2014) 1.2
5. Bolte, J., Daniilidis, A., Lewis, A.: The Lojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. *SIAM Journal on Optimization* **17**(4), 1205–1223 (2007) 1
6. Bouaziz, S., Tagliasacchi, A., Pauly, M.: Sparse iterative closest point. In: *Computer graphics forum*, vol. 32, pp. 113–123. Wiley Online Library (2013) 1
7. Cai, J.F., Candès, E.J., Shen, Z.: A singular value thresholding algorithm for matrix completion. *SIAM Journal on Optimization* **20**(4), 1956–1982 (2010) 5
8. Chartrand, R.: Nonconvex splitting for regularized low-rank+ sparse decomposition. *Signal Processing, IEEE Transactions on* **60**(11), 5810–5819 (2012) 5
9. Chartrand, R., Wohlberg, B.: A nonconvex admm algorithm for group sparsity with sparse groups. In: *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pp. 6009–6013. IEEE (2013) 1
10. Chen, C., He, B., Ye, Y., Yuan, X.: The direct extension of ADMM for multi-block convex minimization problems is not necessarily convergent. *Mathematical Programming* pp. 1–23 (2014) 1, 4.1
11. Chen C., Yuan, X., Zeng, S., Zhang, J., Penalty splitting methods for solving mathematical program with equilibrium constraints. Manuscript (private communication), (2016) 5
12. Conn, A.R., Gould, N.I., Toint, P.: A globally convergent augmented lagrangian algorithm for optimization with general constraints and simple bounds. *SIAM Journal on Numerical Analysis* **28**(2), 545–572 (1991) 1.2
13. Cottle, R., and Dantzig, G.: Complementary pivot theory of mathematical programming. *Linear Algebra and its Applications* **1**, 103–125 (1968) 5
14. Daubechies, I., DeVore, R., Fornasier, M., Güntürk, C.S.: Iteratively reweighted least squares minimization for sparse recovery. *Communications on Pure and Applied Mathematics* **63**(1), 1–38 (2010) 5

15. Davis, D., Yin, W.: Convergence rate analysis of several splitting schemes. In (Glowinski, R., Osher, S., Yin, W. ed.) *Splitting Methods in Communication, Imaging, Science and Engineering*. Springer, New York, (2016) 1.3
16. Davis, D., Yin, W.: Convergence rates of relaxed Peaceman-Rachford and ADMM under regularity assumptions. arXiv preprint arXiv:1407.5210 (2014). 1.3
17. Deng, W., Lai, M.J., Peng, Z., Yin, W.: Parallel multi-block ADMM with  $o(1/k)$  convergence. *Journal of Scientific Computing*, online first, (2016) 3.1
18. 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) 1.3
19. Glowinski, R.: Numerical methods for nonlinear variational problems. Springer series in computational physics. Springer-Verlag, New York (1984) 1.3
20. Glowinski, R., Marroco, A.: On the approximation by finite elements of order one, and resolution, penalisation-duality for a class of nonlinear dirichlet problems. *ESAIM: Mathematical Modelling and Numerical Analysis - Mathematical Modelling and Numerical Analysis* **9**(R2), 41–76 (1975) 1.3
21. He, B., Yuan, X.: On the  $o(1/n)$  convergence rate of the Douglas-Rachford alternating direction method. *SIAM Journal on Numerical Analysis* **50**(2), 700–709 (2012) 1.3
22. Hestenes, M.R.: Multiplier and gradient methods. *Journal of optimization theory and applications* **4**(5), 303–320 (1969) 1.2
23. Hong, M., Luo, Z.Q., Razaviyayn, M.: Convergence Analysis of Alternating Direction Method of Multipliers for a Family of Nonconvex Problems. *SIAM Journal on Optimization*, **26**(1), 337–364, (2016) 1.3, 1.4, 4.1
24. Hu, Y., Chi, E., Allen, G.I.: ADMM algorithmic regularization paths for sparse statistical machine learning. In (Glowinski, R., Osher, S., Yin, W. ed.) *Splitting Methods in Communication, Imaging, Science and Engineering*. Springer, New York, (2016) 5
25. Ivanov, M., Zlateva, N.: Abstract subdifferential calculus and semi-convex functions. *Serdica Mathematical Journal* **23**(1), 35p–58p (1997) 2.1
26. Jiang, B., Ma, S., Zhang, S.: Alternating direction method of multipliers for real and complex polynomial optimization models. *Optimization* **63**(6), 883–898 (2014) 1.3, 1.4
27. Knopp, K.: Infinite sequences and series. Courier Corporation (1956) 3.1
28. Kryštof, V., Zajíček, L.: Differences of two semiconvex functions on the real line. preprint (2015) 2.1
29. Lai, R., Osher, S.: A splitting method for orthogonality constrained problems. *Journal of Scientific Computing* **58**(2), 431–449 (2014) 1, 5
30. Li, G., Pong, T.K.: Global convergence of splitting methods for nonconvex composite optimization. *SIAM Journal on Optimization*, **25**(4), 2434–2460, (2015) 1.3, 1.4, 4.1, 4.1
31. Liavas, A.P., Sidiropoulos, N.D.: Parallel algorithms for constrained tensor factorization via the alternating direction method of multipliers. *IEEE Transactions on Signal Processing*, **63**(20), 5450 – 5463, (2015) 1
32. Łojasiewicz, S.: Sur la géométrie semi-et sous-analytique. *Ann. Inst. Fourier (Grenoble)* **43**(5), 1575–1595 (1993) 1
33. Lu, Z., Zhang, Y.: An augmented lagrangian approach for sparse principal component analysis. *Mathematical Programming* **135**(1-2), 149–193 (2012). 5
34. Magnússon, S., Weeraddana, P.C., Rabbat, M.G., Fischione, C.: On the convergence of alternating direction lagrangian methods for nonconvex structured optimization problems. *IEEE Transactions on Control of Network Systems*, **3**(3), 296 – 309 (2015) 1.3, 1.4
35. Mifflin, R.: Semismooth and semiconvex functions in constrained optimization. *SIAM Journal on Control and Optimization* **15**(6), 959–972 (1977) 2.1
36. Miksik, O., Vineet, V., Pérez, P., Torr, P.H., Cesson Sévigné, F.: Distributed non-convex ADMM-inference in large-scale random fields. In: *British Machine Vision Conference, BMVC* (2014) 1
37. Mo llenhoff, T., Strelakovsky, E., Moeller, M., Cremers, D.: The primal-dual hybrid gradient method for semiconvex splittings. *SIAM Journal on Imaging Sciences* **8**(2), 827–857 (2015) 2.1
38. Oymak, S., Mohan, K., Fazel, M., Hassibi, B.: A simplified approach to recovery conditions for low rank matrices. In: *Information Theory Proceedings (ISIT), 2011 IEEE International Symposium on*, pp. 2318–2322. IEEE (2011) 5
39. Poliquin, R., Rockafellar, R.: Prox-regular functions in variational analysis. *Transactions of the American Mathematical Society* **348**(5), 1805–1838 (1996) 2.1
40. Powell, M.J.: A method for non-linear constraints in minimization problems. UKAEA (1967) 1.2
41. Rockafellar, R.T., Wets, R.J.B.: *Variational Analysis* (2009) 2.1, 2.1, 3.1
42. Rosenberg, J., et al.: Applications of analysis on lipschitz manifolds. *Proc. Miniconferences on Harmonic Analysis and Operator Algebras* (Canberra, t987), *Proc. Centre for Math. Analysis* **16**, 269–283 (1988) 1
43. Shen, Y., Wen, Z., Zhang, Y.: Augmented lagrangian alternating direction method for matrix separation based on low-rank factorization. *Optimization Methods Software* **29**(2), 239–263 (2014) 1
44. Slavakis, K., Giannakis, G., Mateos, G.: Modeling and optimization for big data analytics:(statistical) learning tools for our era of data deluge. *Signal Processing Magazine, IEEE* **31**(5), 18–31 (2014) 5

- 
45. Sun, D.L., Fevotte, C.: Alternating direction method of multipliers for non-negative matrix factorization with the beta-divergence. In: Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference on, pp. 6201–6205. IEEE (2014) 1, 5
  46. Wang, F., Cao, W., Xu, Z.: Convergence of multi-block Bregman ADMM for nonconvex composite problems. arXiv preprint arXiv:1505.03063 (2015) 1.3, 1.4, 4.1
  47. Wang, F., Xu, Z., Xu, H.K.: Convergence of Bregman alternating direction method with multipliers for nonconvex composite problems. arXiv preprint arXiv:1410.8625 (2014) 1.3, 1.4
  48. Wang, X., Hong, M., Ma, S., Luo, Z.Q.: Solving multiple-block separable convex minimization problems using two-block alternating direction method of multipliers. arXiv preprint arXiv:1308.5294 (2013) 5
  49. Wen, Z., Peng, X., Liu, X., Sun, X., Bai, X.: Asset Allocation under the Basel Accord Risk Measures. ArXiv preprint ArXiv:1308.1321 (2013) 1
  50. Wen, Z., Yang, C., Liu, X., Marchesini, S.: Alternating direction methods for classical and ptychographic phase retrieval. *Inverse Problems* **28**(11), 115,010 (2012) 1
  51. Wen, Z., Yin, W.: A feasible method for optimization with orthogonality constraints. *Mathematical Programming* **142**(1-2), 397–434 (2013) 5
  52. Wikipedia: Schatten norm — Wikipedia, the free encyclopedia (2015). [Online; accessed 18-October-2015] 1
  53. Xu, Y., Yin, W.: A block coordinate descent method for regularized multiconvex optimization with applications to non-negative tensor factorization and completion. *SIAM Journal on Imaging Sciences* **6**(3), 1758–1789 (2013) 2.2, 3.1
  54. Xu, Y., Yin, W., Wen, Z., Zhang, Y.: An alternating direction algorithm for matrix completion with nonnegative factors. *Frontiers of Mathematics in China* **7**(2), 365–384 (2012) 1, 1.3, 1.4, 5
  55. Yan, M. and Yin, W.: Self equivalence of the alternating direction method of multipliers. In (Glowinski, R., Osher, S., Yin, W. ed.) *Splitting Methods in Communication, Imaging, Science and Engineering*. Springer, New York, (2016) 1
  56. Yang, L., Pong, T.K., Chen, X.: Alternating direction method of multipliers for nonconvex background/foreground extraction. arXiv preprint arXiv:1506.07029 (2015) 1, 5, 5
  57. You, S., Peng, Q.: A non-convex alternating direction method of multipliers heuristic for optimal power flow. In: *Smart Grid Communications (SmartGridComm), 2014 IEEE International Conference on*, pp. 788–793. IEEE (2014) 1
  58. Zeng J., Lin S., Wang Y., Xu Z.:  $L_{1/2}$  Regularization: convergence of iterative half thresholding algorithm, *IEEE Transactions on Signal Processing*, 62(9): 2317-2329, 2014. 5
  59. Zeng, J., Lin, S., Xu, Z.: Sparse regularization: Convergence of iterative jumping thresholding algorithm. *IEEE Transactions on Signal Processing*, 64(19): 5106-5117, 2016. 5
  60. Zeng, J., Peng, Z., Lin, S.: A Gauss-Seidel iterative thresholding algorithm for lq regularized least squares regression. arXiv preprint arXiv:1507.03173 (2015) 5