

# CONVERGENCE ANALYSIS OF ALTERNATING DIRECTION METHOD OF MULTIPLIERS FOR A FAMILY OF NONCONVEX PROBLEMS

MINGYI HONG, ZHI-QUAN LUO AND MEISAM RAZAVIYAYN\*

**Abstract.** The alternating direction method of multipliers (ADMM) is widely used to solve large-scale linearly constrained optimization problems, convex or nonconvex, in many engineering fields. However there is a general lack of theoretical understanding of the algorithm when the objective function is nonconvex. In this paper we analyze the convergence of the ADMM for solving certain nonconvex *consensus* and *sharing* problems. We show that the classical ADMM converges to the set of stationary solutions, provided that the penalty parameter in the augmented Lagrangian is chosen to be sufficiently large. For the sharing problems, we show that the ADMM is convergent regardless of the number of variable blocks. Our analysis does not impose any assumptions on the iterates generated by the algorithm, and is broadly applicable to many ADMM variants involving proximal update rules and various flexible block selection rules.

**AMS(MOS) Subject Classifications:** 49, 90.

**1. Introduction.** Consider the following linearly constrained (possibly nonsmooth or/and nonconvex) problem with  $K$  blocks of variables  $\{x_k\}_{k=1}^K$ :

$$\begin{aligned} \min \quad & f(x) := \sum_{k=1}^K g_k(x_k) + \ell(x_1, \dots, x_K) \\ \text{s.t.} \quad & \sum_{k=1}^K A_k x_k = q, \quad x_k \in X_k, \quad \forall k = 1, \dots, K \end{aligned} \quad (1.1)$$

where  $A_k \in \mathbb{R}^{M \times N_k}$  and  $q \in \mathbb{R}^M$ ;  $X_k \subset \mathbb{R}^{N_k}$  is a closed convex set;  $\ell(\cdot)$  is a smooth (possibly nonconvex) function; each  $g_k(\cdot)$  can be either a smooth function, or a convex nonsmooth function. Let us define  $A := [A_1, \dots, A_K]$ . The augmented Lagrangian for problem (1.1) is given by

$$L(x; y) = \sum_{k=1}^K g_k(x_k) + \ell(x_1, \dots, x_K) + \langle y, q - Ax \rangle + \frac{\rho}{2} \|q - Ax\|^2, \quad (1.2)$$

where  $\rho > 0$  is a constant representing the primal penalty parameter.

To solve problem (1.1), let us consider a popular algorithm called the alternating direction method of multipliers (ADMM), whose steps are given below:

---

\*A conference version of the paper has been presented in 2015 International Conference on Acoustics Speech and Signal Processing (ICASSP) [1].

M. Hong is with the Department of Industrial and Manufacturing Systems Engineering, Iowa State University, Ames, IA 50011, USA. Email: [mingyi@iastate.edu](mailto:mingyi@iastate.edu).

Z.-Q. Luo is with the Chinese University of Hong Kong, Shenzhen, China and Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN 55455, USA. Email: [luozq@cuhk.edu.cn](mailto:luozq@cuhk.edu.cn) and [luozq@umn.edu](mailto:luozq@umn.edu).

M. Razaviyayn is with the Department of Electrical Engineering, Stanford University. Email: [meisamr@stanford.edu](mailto:meisamr@stanford.edu)

M. Hong is supported by the National Science Foundation (NSF), grant CCF-1526078, and by the Air Force Office of Scientific Research (AFOSR), grant 15RT0767. Z.-Q. Luo is supported by the National Science Foundation (NSF), grant CCF-1526434.

**Algorithm 0. ADMM for Problem (1.1)**

At each iteration  $t + 1$ , update the primal variables:

$$x_k^{t+1} = \arg \min_{x_k \in X_k} L(x_1^{t+1}, \dots, x_{k-1}^{t+1}, x_k, x_{k+1}^t, \dots, x_K^t; y^t), \quad \forall k = 1, \dots, K. \quad (1.3)$$

Update the dual variable:

$$y^{t+1} = y^t + \rho(q - Ax^{t+1}). \quad (1.4)$$

The ADMM algorithm was originally introduced in early 1970s [2,3], and has since been studied extensively [4–7]. Recently it has become widely popular in modern big data related problems arising in machine learning, computer vision, signal processing, networking and so on; see [8–15] and the references therein. In practice, the algorithm often exhibits faster convergence than traditional primal-dual type algorithms such as the dual ascent algorithm [16–18] or the method of multipliers [19]. It is also particularly suitable for parallel implementation [8].

There is a vast literature that applies the ADMM to various problems in the form of (1.1). Unfortunately, theoretical understanding of the algorithm is still fairly limited. For example, most of its convergence analysis is done for certain special form of problem (1.1) — the *two-block convex separable* problems, where  $K = 2$ ,  $\ell = 0$  and  $g_1, g_2$  are both convex. In this case, ADMM is known to converge under very mild conditions; see [7] and [8]. Under the same conditions, several recent works [20–22] have shown that the ADMM converges with the sublinear rate of  $\mathcal{O}(\frac{1}{t})$  or  $o(\frac{1}{t})$ , and it converges with a rate  $\mathcal{O}(\frac{1}{t^2})$  when properly accelerated [23,24]. Reference [25] has shown that the ADMM converges linearly when the objective function as well as the constraints satisfy certain additional assumptions. For the *multi-block* separable convex problems where  $K \geq 3$ , it is known that the original ADMM can diverge for certain pathological problems [26]. Therefore, most research effort in this direction has been focused on either analyzing problems with additional conditions, or showing convergence for variants of the ADMM; see for example [26–34]. It is worth mentioning that when the objective function is not separable across the variables (e.g., the coupling function  $\ell(\cdot)$  appears in the objective), the convergence of the ADMM is still open, even in the case where  $K = 2$  and  $f(\cdot)$  is convex. Recent works of [29,35] have shown that when problem (1.1) is convex but not necessarily separable, and when certain error bound condition is satisfied, then the ADMM iteration converges to the set of primal-dual optimal solutions, provided that the dual stepsize decreases in time. Another recent work in this direction can be found in [36].

Unlike the convex case, for which the behavior of ADMM has been investigated quite extensively, when the objective becomes nonconvex, the convergence issue of ADMM remains largely open. Nevertheless, it has been observed by many researchers that the ADMM works extremely well for various applications involving nonconvex objectives, such as the nonnegative matrix factorization [37,38], phase retrieval [39], distributed matrix factorization [40], distributed clustering [41], sparse zero variance discriminant analysis [42], polynomial optimization [43], tensor decomposition [44], matrix separation [45], matrix completion [46], asset allocation [47], sparse feedback control [48] and so on. However, to the best of our knowledge, existing convergence analysis of ADMM for nonconvex problems is very limited — all known global convergence analysis needs to impose uncheckable conditions on the sequence generated by the algorithm. For example, references [43,45–47] show global convergence of the ADMM to the set of stationary solutions for their respective nonconvex problems, by

making the key assumptions that the limit points do exist, and that the successive differences of the iterates (both primal and dual) converge to zero. However such assumption is nonstandard and overly restrictive. It is not clear whether the same convergence result can be claimed without making assumptions on the iterates. Reference [49] analyzes a family of splitting algorithms (which includes the ADMM as a special case) for certain nonconvex quadratic optimization problem, and shows that they converge to the stationary solution when certain condition on the dual stepsize is met. We note that there has been many recent works proposing new algorithms to solve nonconvex and nonsmooth problems, for example [50–54]. However, these works do not deal with nonconvex problems with linearly coupling constraints, and their analysis does not directly apply to the ADMM-type methods.

The aim of this paper is to provide some theoretical justification on the good performance of the ADMM for nonconvex problems. Specifically, we establish the convergence of ADMM for certain types of nonconvex problems including the consensus and sharing problems without making any assumptions on the iterates. Our analysis shows that, as long as the objective functions  $g_k$ 's and  $\ell$  satisfy certain regularity conditions, and the penalty parameter  $\rho$  is chosen large enough (with computable bounds), then the iterates generated by the ADMM is guaranteed to converge to the set of stationary solutions. It should be noted that our analysis covers many variants of the ADMM including per-block proximal update and flexible block selection. An interesting consequence of our analysis is that for a particular reformulation of the sharing problem, the *multi-block* ADMM algorithm converges, regardless of the convexity of the objective function. Finally, to facilitate possible applications to other nonconvex problems, we highlight the main proof steps in our analysis framework that can guarantee the global convergence of the ADMM iterates (1.3)–(1.4) to the set of stationary solutions.

## 2. The Nonconvex Consensus Problem.

**2.1. The Basic Problem.** Consider the following nonconvex global consensus problem with regularization

$$\begin{aligned} \min \quad & f(x) := \sum_{k=1}^K g_k(x) + h(x) \\ \text{s.t.} \quad & x \in X \end{aligned} \tag{2.1}$$

where  $g_k$ 's are a set of smooth, possibly nonconvex functions, while  $h(x)$  is a convex nonsmooth regularization term. This problem is related to the convex global consensus problem discussed heavily in [8, Section 7], but with the important difference that  $g_k$ 's can be nonconvex.

In many practical applications,  $g_k$ 's need to be handled by a single agent, such as a thread or a processor. This motivates the following consensus formulation. Let us introduce a set of new variables  $\{x_k\}_{k=0}^K$ , and transform problem (2.1) equivalently to the following linearly constrained problem

$$\begin{aligned} \min \quad & \sum_{k=1}^K g_k(x_k) + h(x_0) \\ \text{s.t.} \quad & x_k = x_0, \forall k = 1, \dots, K, \quad x_0 \in X. \end{aligned} \tag{2.2}$$

We note that after reformulation, the problem dimension is increased by  $K$  due to the introduction of auxiliary variables  $\{x_1, \dots, x_K\}$ . Consequently, solving the refor-

mulated problem (2.2) distributedly may not be as efficient (in terms of total number of iterations required) as applying the centralized algorithms [50–54] directly to the original problem (2.1). Nonetheless, a major benefit of solving the reformulated problem (2.2) is the flexibility of allowing each distributed agent to handle a single *local* variable  $x_k$  and a *local* function  $g_k$ .

The augmented Lagrangian function is given by

$$L(\{x_k\}, x_0; y) = \sum_{k=1}^K g_k(x_k) + h(x_0) + \sum_{k=1}^K \langle y_k, x_k - x_0 \rangle + \sum_{k=1}^K \frac{\rho_k}{2} \|x_k - x_0\|^2. \quad (2.3)$$

Note that this augmented Lagrangian is slightly different from the one expressed in (1.2), as we have used a set of different penalization parameters  $\{\rho_k\}$ , one for each equality constraint  $x_k = x_0$ . We note that there can be many other variants of the basic consensus problem, such as the *general form consensus optimization*, the *sharing* problem and so on. We will discuss some of those variants in the later sections.

**2.2. The ADMM Algorithm for Nonconvex Consensus.** The problem (2.2) can be solved distributedly by applying the classical ADMM. The details are given in the table below.

<b>Algorithm 1. The Classical ADMM for Problem (2.2)</b>	
At each iteration $t + 1$ , compute:	
$x_0^{t+1} = \arg \min_{x_0 \in X} L(\{x_k^t\}, x_0; y^t).$	(2.4)
Each node $k$ computes $x_k$ by solving:	
$x_k^{t+1} = \arg \min_{x_k} g_k(x_k) + \langle y_k^t, x_k - x_0^{t+1} \rangle + \frac{\rho_k}{2} \ x_k - x_0^{t+1}\ ^2.$	(2.5)
Each node $k$ updates the dual variable:	
$y_k^{t+1} = y_k^t + \rho_k (x_k^{t+1} - x_0^{t+1}).$	(2.6)

In the  $x_0$  update step, if the nonsmooth penalization  $h(\cdot)$  does not appear in the objective, then this step can be written as

$$x_0^{t+1} = \arg \min_{x_0 \in X} L(\{x_k^t\}, x_0; y^t) = \text{proj}_X \left[ \frac{\sum_{k=1}^K \rho_k x_k^t + \sum_{k=1}^K y_k^t}{\sum_{k=1}^K \rho_k} \right]. \quad (2.7)$$

Note that the above algorithm has the exact form as the classical ADMM described in [8], where the variable  $x_0$  is taken as the first block of primal variable, and the collection  $\{x_k\}_{k=1}^K$  as the second block. The two primal blocks are updated in a sequential (i.e., Gauss-Seidel) manner, followed by an inexact dual ascent step.

In what follows, we consider a more general version of ADMM which includes Algorithm 1 as a special case. In particular, we propose a *flexible* ADMM algorithm in which there is a greater flexibility in choosing the order of the update of both the primal and the dual variables. Specifically, we consider the following two types of variable block update order rules: let  $k = 0, 2, \dots, K$  be the indices for the primal variable blocks  $x_0, x_1, x_2, \dots, x_K$ , and let  $\mathcal{C}^t \subseteq \{0, 1, \dots, K\}$  denote the set of variables updated in iteration  $t$ , then

1. *Randomized update rule*: At each iteration  $t + 1$ , a variable block  $k$  is chosen randomly with probability  $p_k^{t+1}$ ,

$$\Pr(k \in \mathcal{C}^{t+1} \mid x_0^t, y^t, \{x_k^t\}) = p_k^{t+1} \geq p_{\min} > 0. \quad (2.8)$$

2. *Essentially cyclic update rule*: There exists a given period  $T \geq 1$  during which each index is updated at least once. More specifically, at iteration  $t$ , update all the variables in an index set  $\mathcal{C}^t$  whereby

$$\bigcup_{i=1}^T \mathcal{C}^{t+i} = \{0, 1, \dots, K\}, \forall t. \quad (2.9)$$

We call this update rule a *period- $T$  essentially cyclic update rule*.

**Algorithm 2. The Flexible ADMM for Problem (2.2)**

Let  $\mathcal{C}^1 = \{0, \dots, K\}$ ,  $t = 0, 1, \dots$ .

At each iteration  $t + 1$ , do:

**If**  $t + 1 \geq 2$ , pick an index set  $\mathcal{C}^{t+1} \subseteq \{0, \dots, K\}$ .

**If**  $0 \in \mathcal{C}^{t+1}$ , compute:

$$x_0^{t+1} = \operatorname{argmin}_{x \in X} L(\{x_k^t\}, x_0; y^t). \quad (2.10)$$

**Else**  $x_0^{t+1} = x_0^t$ .

**If**  $k \neq 0$  and  $k \in \mathcal{C}^{t+1}$ , node  $k$  computes  $x_k$  by solving:

$$x_k^{t+1} = \arg \min_{x_k} g_k(x_k) + \langle y_k^t, x_k - x_0^{t+1} \rangle + \frac{\rho_k}{2} \|x_k - x_0^{t+1}\|^2. \quad (2.11)$$

Update the dual variable:

$$y_k^{t+1} = y_k^t + \rho_k (x_k^{t+1} - x_0^{t+1}). \quad (2.12)$$

**Else**  $x_k^{t+1} = x_k^t$ ,  $y_k^{t+1} = y_k^t$ .

We note that the randomized version of Algorithm 2 is similar to that of the convex consensus algorithms studied in [55] and [56]. It is also related to the *randomized BSUM-M* algorithm studied in [29]. The difference with the latter is that in the randomized BSUM-M, the dual variable is viewed as an additional block that can be randomly picked (independent of the way that the primal blocks are picked), whereas in Algorithm 2, the dual variable  $y_k$  is always updated whenever the corresponding primal variable  $x_k$  is updated. To the best of our knowledge, the period- $T$  essentially cyclic update rule is a new variant of the ADMM.

Notice that Algorithm 1 is simply the period-1 essentially cyclic rule, which is a special case of Algorithm 2. Therefore we will focus on analyzing Algorithm 2. To this end, we make the following assumption.

**Assumption A.**

- A1. There exists a positive constant  $L_k > 0$  such that

$$\|\nabla_k g_k(x_k) - \nabla_k g_k(z_k)\| \leq L_k \|x_k - z_k\|, \forall x_k, z_k, k = 1, \dots, K.$$

Moreover,  $h$  is convex (possibly nonsmooth);  $X$  is a closed convex set.

- A2. For all  $k$ , the penalty parameter  $\rho_k$  is chosen large enough such that:
1. For all  $k$ , the  $x_k$  subproblem (2.11) is strongly convex with modulus  $\gamma_k(\rho_k)$ ;
  2. For all  $k$ ,  $\rho_k \gamma_k(\rho_k) > 2L_k^2$  and  $\rho_k \geq L_k$ .
- A3.  $f(x)$  is bounded from below over  $X$ , that is,

$$\underline{f} := \min_{x \in X} f(x) > -\infty.$$

We have the following remarks regarding to the assumptions made above.

- As  $\rho_k$  inceases, the subproblem (2.11) will be eventually strongly convex with respect to  $x_k$ . The corresponding strong convexity modulus  $\gamma_k(\rho_k)$  is a monotonic increasing function of  $\rho_k$ .
- Whenever  $g_k(\cdot)$  is nonconvex (therefore  $\rho_k > \gamma_k(\rho_k)$ ), the condition  $\rho_k \gamma_k(\rho_k) \geq 2L_k^2$  implies  $\rho_k \geq L_k$ .
- By construction,  $L(\{x_k\}, x_0; y)$  is also strongly convex with respect to  $x_0$ , with a modulus  $\gamma := \sum_{k=1}^K \rho_k$ .
- Assumption A makes no assumption on the *iterates* generated by the algorithm. This is in contrast to the existing analysis of the nonconvex ADMM algorithms [37, 43, 46].

Now we begin to analyze Algorithm 2. We first make several definitions. Let  $t(k)$  (resp.  $t(0)$ ) denote the latest iteration index that  $x_k$  (resp.  $x_0$ ) is updated before iteration  $t + 1$ , i.e.,

$$\begin{aligned} t(k) &= \max \{r \mid r \leq t, k \in \mathcal{C}^r\}, \quad k = 1, \dots, K, \\ t(0) &= \max \{r \mid r \leq t, 0 \in \mathcal{C}^r\}. \end{aligned} \quad (2.13)$$

This definition implies that  $x_k^t = x_k^{t(k)}$  for all  $k = 0, \dots, K$ .

Also define new vectors  $\hat{x}_0^{t+1}$ ,  $\{\hat{x}_k^{t+1}\}$ ,  $\hat{y}^{t+1}$  and  $\{\tilde{x}_k^{t+1}\}$ ,  $\tilde{y}^{t+1}$  by

$$\hat{x}_0^{t+1} = \arg \min_{x_0 \in X} L(\{x_k^t\}, x_0; y^t), \quad (2.14a)$$

$$\hat{x}_k^{t+1} = \arg \min_{x_k} g_k(x_k) + \langle y_k^t, x_k - \hat{x}_0^{t+1} \rangle + \frac{\rho_k}{2} \|x_k - \hat{x}_0^{t+1}\|^2, \quad \forall k \quad (2.14b)$$

$$\hat{y}_k^{t+1} = y_k^t + \rho_k (\hat{x}_k^{t+1} - \hat{x}_0^{t+1}). \quad (2.14c)$$

$$\tilde{x}_k^{t+1} = \arg \min_{x_k} g_k(x_k) + \langle y_k^t, x_k - x_0^t \rangle + \frac{\rho_k}{2} \|x_k - x_0^t\|^2, \quad \forall k \quad (2.14d)$$

$$\tilde{y}_k^{t+1} = y_k^t + \rho_k (\tilde{x}_k^{t+1} - x_0^t). \quad (2.14e)$$

In words,  $(\hat{x}_0^{t+1}, \{\hat{x}_k^{t+1}\}, \hat{y}^{t+1})$  is a “virtual” iterate assuming that all variables are updated at iteration  $t + 1$ .  $(\tilde{x}_k^{t+1}, \tilde{y}^{t+1})$  is a “virtual” iterate for the case where  $x_0$  is not updated but the rest of variables are updated.

We first show that the size of the successive difference of the dual variables can be bounded above by that of the primal variables.

LEMMA 2.1. *Suppose Assumption A holds. Then for Algorithm 2 with either randomized or essentially cyclic update rule, the following are true*

$$L_k^2 \|x_k^{t+1} - x_k^t\|^2 \geq \|y_k^{t+1} - y_k^t\|^2, \quad \forall k = 1, \dots, K, \quad (2.15a)$$

$$L_k^2 \|\hat{x}_k^{t+1} - x_k^t\|^2 \geq \|\hat{y}_k^{t+1} - y_k^t\|^2, \quad \forall k = 1, \dots, K. \quad (2.15b)$$

$$L_k^2 \|\tilde{x}_k^{t+1} - x_k^t\|^2 \geq \|\tilde{y}_k^{t+1} - y_k^t\|^2, \quad \forall k = 1, \dots, K. \quad (2.15c)$$

**Proof.** We will show the first inequality. The second inequality follows a similar line of argument.

To prove (2.15a), first note that the case for  $k \notin \mathcal{C}^{t+1}$  is trivial, as both sides of (2.15a) evaluate to zero. Suppose  $k \in \mathcal{C}^{t+1}$ . From the  $x_k$  update step (2.11) we have the following optimality condition

$$\nabla g_k(x_k^{t+1}) + y_k^t + \rho_k(x_k^{t+1} - x_0^{t+1}) = 0, \forall k \in \mathcal{C}^{t+1}/\{0\}. \quad (2.16)$$

Combined with the dual variable update step (2.12) we obtain

$$\nabla g_k(x_k^{t+1}) = -y_k^{t+1}, \forall k \in \mathcal{C}^{t+1}/\{0\}. \quad (2.17)$$

Combining this with Assumption A1, and noting that for any given  $k$ ,  $y_k$  and  $x_k$  are always updated in the same iteration, we obtain for all  $k \in \mathcal{C}^{t+1}/\{0\}$

$$\begin{aligned} \|y_k^{t+1} - y_k^t\| &= \|y_k^{t+1} - y_k^{t(k)}\| \\ &= \|\nabla g_k(x_k^{t+1}) - \nabla g_k(x_k^{t(k)})\| \leq L_k \|x_k^{t+1} - x_k^{t(k)}\| = L_k \|x_k^{t+1} - x_k^t\|. \end{aligned}$$

The desired result follows.  $\square$

Next, we use (2.15a) to bound the difference of the augmented Lagrangian.

**LEMMA 2.2.** *For Algorithm 2 with either randomized or period- $T$  essentially cyclic update rule, we have the following*

$$\begin{aligned} &L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t) \\ &\leq \sum_{k \neq 0, k \in \mathcal{C}^{t+1}} \left( \frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2} \right) \|x_k^{t+1} - x_k^t\|^2 - \frac{\gamma}{2} \|x_0^{t+1} - x_0^t\|^2. \end{aligned} \quad (2.18)$$

**Proof.** We first split the successive difference of the augmented Lagrangian by

$$\begin{aligned} &L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t) \\ &= (L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^{t+1}\}, x_0^{t+1}; y^t)) \\ &\quad + (L(\{x_k^{t+1}\}, x_0^{t+1}; y^t) - L(\{x_k^t\}, x_0^t; y^t)). \end{aligned} \quad (2.19)$$

The first term in (2.19) can be bounded by

$$\begin{aligned} &L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^{t+1}\}, x_0^{t+1}; y^t) \\ &= \sum_{k=1}^K \langle y_k^{t+1} - y_k^t, x_k^{t+1} - x_0^{t+1} \rangle \\ &\stackrel{(a)}{=} \sum_{k \neq 0, k \in \mathcal{C}^{t+1}} \frac{1}{\rho_k} \|y_k^{t+1} - y_k^t\|^2 \end{aligned} \quad (2.20)$$

where in (a) we have use (2.12), and the fact that  $y_k^{t+1} - y_k^t = 0$  for all variable block  $x_k$  that has not been updated (i.e.,  $k \neq 0, k \notin \mathcal{C}^{t+1}$ ). The second term in (2.19) can

be bounded by

$$\begin{aligned}
& L(\{x_k^{t+1}\}, x_0^{t+1}; y^t) - L(\{x_k^t\}, x_0^t; y^t) \\
&= L(\{x_k^{t+1}\}, x_0^{t+1}; y^t) - L(\{x_k^t\}, x_0^{t+1}; y^t) + L(\{x_k^t\}, x_0^{t+1}; y^t) - L(\{x_k^t\}, x_0^t; y^t) \\
&\stackrel{(a)}{\leq} \sum_{k=1}^K \left( \langle \nabla_{x_k} L(\{x_k^{t+1}\}, x_0^{t+1}; y^t), x_k^{t+1} - x_k^t \rangle - \frac{\gamma_k(\rho_k)}{2} \|x_k^{t+1} - x_k^t\|^2 \right) \\
&\quad + \langle \zeta_{x_0}^{t+1}, x_0^{t+1} - x_0^t \rangle - \frac{\gamma}{2} \|x_0^{t+1} - x_0^t\|^2 \\
&\stackrel{(b)}{=} \sum_{k \neq 0, k \in \mathcal{C}^{t+1}} \left( \langle \nabla_{x_k} L(\{x_k^{t+1}\}, x_0^{t+1}; y^t), x_k^{t+1} - x_k^t \rangle - \frac{\gamma_k(\rho_k)}{2} \|x_k^{t+1} - x_k^t\|^2 \right) \\
&\quad + \iota\{0 \in \mathcal{C}^{t+1}\} \left( \langle \zeta_{x_0}^{t+1}, x_0^{t+1} - x_0^t \rangle - \frac{\gamma}{2} \|x_0^{t+1} - x_0^t\|^2 \right) \\
&\stackrel{(c)}{\leq} - \sum_{k \neq 0, k \in \mathcal{C}^{t+1}} \frac{\gamma_k(\rho_k)}{2} \|x_k^{t+1} - x_k^t\|^2 - \iota\{0 \in \mathcal{C}^{t+1}\} \frac{\gamma}{2} \|x_0^{t+1} - x_0^t\|^2, \tag{2.21}
\end{aligned}$$

where in (a) we have used the fact that  $L(\{x_k\}, x_0; y)$  is strongly convex w.r.t. each  $x_k$  and  $x_0$ , with modulus  $\gamma_k(\rho_k)$  and  $\gamma$ , respectively, and that

$$\zeta_{x_0}^{t+1} \in \partial_{x_0} L(\{x_k^t\}, x_0^{t+1}; y^t)$$

is some subgradient vector; in (b) we have used the fact that when  $k \notin \mathcal{C}^{t+1}$  (resp.  $0 \notin \mathcal{C}^{t+1}$ ),  $x_k^{t+1} = x_k^t$  (resp.  $x_0^{t+1} = x_0^t$ ), and we have defined  $\iota\{0 \in \mathcal{C}^{t+1}\}$  as the indicator function that takes the value 1 if  $0 \in \mathcal{C}^{t+1}$  is true, and takes value 0 otherwise; in (c) we have used the optimality of each subproblem (2.11) and (2.10) (where  $\zeta_{x_0}^{t+1}$  is specialized to the subgradient vector that satisfies the optimality condition for problem (2.10)).

Combining the above two inequalities (2.20) and (2.21), we obtain

$$\begin{aligned}
& L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t) \\
&\leq - \sum_{k \neq 0, k \in \mathcal{C}^{t+1}} \frac{\gamma_k(\rho_k)}{2} \|x_k^{t+1} - x_k^t\|^2 + \sum_{k \neq 0, k \in \mathcal{C}^{t+1}} \frac{1}{\rho_k} \|y_k^{t+1} - y_k^t\|^2 - \iota\{0 \in \mathcal{C}^{t+1}\} \frac{\gamma}{2} \|x_0^{t+1} - x_0^t\|^2 \\
&\leq \sum_{k \neq 0, k \in \mathcal{C}^{t+1}} \left( \frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2} \right) \|x_k^{t+1} - x_k^t\|^2 - \iota\{0 \in \mathcal{C}^{t+1}\} \frac{\gamma}{2} \|x_0^{t+1} - x_0^t\|^2
\end{aligned}$$

where the last inequality is due to (2.15a). The desired result is obtained by noticing the fact that when  $0 \notin \mathcal{C}^{t+1}$ , we have  $x_0^{t+1} - x_0^t = 0$ .  $\square$

The above result implies that if the following condition is satisfied:

$$\rho_k \gamma_k(\rho_k) \geq 2L_k^2, \quad \forall k = 1, \dots, K, \tag{2.22}$$

then the value of the augmented Lagrangian function will always decrease. Note that as long as  $\gamma_k(\rho_k) \neq 0$ , one can always find a  $\rho_k$  large enough such that the above condition is satisfied, as the left hand side (lhs) of (2.22) is monotonically increasing w.r.t.  $\rho_k$ , while the right hand side (rhs) is a constant.

Next we show that  $L(\{x_k^t\}, x_0^t; y^t)$  is in fact convergent.



LEMMA 2.3. Suppose Assumption A is true. Let  $\{\{x_k^t\}, x_0^t, y^t\}$  be generated by Algorithm 2 with either the essentially cyclic rule or the randomized rule. Then the following limit exists and is lower bounded by  $\underline{f}$  defined in Assumption A3:

$$\lim_{t \rightarrow \infty} L(\{x_k^t\}, x_0^t, y^t) \geq \underline{f}. \quad (2.23)$$

**Proof.** Notice that the augmented Lagrangian function can be expressed as

$$\begin{aligned} & L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) \\ &= h(x_0^{t+1}) + \sum_{k=1}^K \left( g_k(x_k^{t+1}) + \langle y_k^{t+1}, x_k^{t+1} - x_0^{t+1} \rangle + \frac{\rho_k}{2} \|x_k^{t+1} - x_0^{t+1}\|^2 \right) \\ &\stackrel{(a)}{=} h(x_0^{t+1}) + \sum_{k=1}^K \left( g_k(x_k^{t+1}) + \langle \nabla g_k(x_k^{t+1}), x_0^{t+1} - x_k^{t+1} \rangle + \frac{\rho_k}{2} \|x_k^{t+1} - x_0^{t+1}\|^2 \right) \\ &\stackrel{(b)}{\geq} h(x_0^{t+1}) + \sum_{k=1}^K g_k(x_0^{t+1}) = f(x_0^{t+1}) \end{aligned} \quad (2.24)$$

where (b) comes from the Lipschitz continuity of the gradient of  $g_k$ 's (Assumption A1), and the fact that  $\rho_k \geq L_k$  for all  $k = 1, \dots, K$  (Assumption A2). To see why (a) is true, we first observe that due to (2.17), we have for all  $k \neq 0$  and  $k \in \mathcal{C}^{t+1}$

$$\langle y_k^{t+1}, x_k^{t+1} - x_0^{t+1} \rangle = \langle \nabla g_k(x_k^{t+1}), x_0^{t+1} - x_k^{t+1} \rangle.$$

For all  $k \neq 0$  and  $k \notin \mathcal{C}^{t+1}$ , it follows from  $x_k^{t+1} = x_k^t = x_k^{t(k)} = x_k^{t(k)+1}$  and  $y_k^{t+1} = y_k^t = y_k^{t(k)} = y_k^{t(k)+1}$  that

$$\begin{aligned} \langle y_k^{t+1}, x_k^{t+1} - x_0^{t+1} \rangle &= \langle y_k^{t(k)+1}, x_k^{t(k)+1} - x_0^{t+1} \rangle \\ &= \langle \nabla g_k(x_k^{t(k)+1}), x_0^{t+1} - x_k^{t(k)+1} \rangle = \langle \nabla g_k(x_k^{t+1}), x_0^{t+1} - x_k^{t+1} \rangle. \end{aligned}$$

Combining these two cases shows that (a) is true.

Clearly, (2.24) and Assumption A3 together imply that  $L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1})$  is lower bounded. This combined with (2.18) says that whenever the penalty parameter  $\rho_k$ 's are chosen sufficiently large (as per Assumption A2),  $L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1})$  is monotonically decreasing and is convergent. This completes the proof.  $\square$

We are now ready to prove our first main result, which asserts that the sequence of iterates generated by Algorithm 2 converges to the set of stationary solution of problem (2.2).

THEOREM 2.4. Assume that Assumption A is satisfied. Then we have the following

1. We have  $\lim_{t \rightarrow \infty} \|x_k^{t+1} - x_0^{t+1}\| = 0$ ,  $k = 1, \dots, K$ , deterministically for the essentially cyclic update rule and almost surely for the randomized update rule.
2. Let  $(\{x_k^*\}, x_0^*, y^*)$  denote any limit point of the sequence  $\{\{x_k^{t+1}\}, x_0^{t+1}, y^{t+1}\}$  generated by Algorithm 2. Then the following statement is true (deterministically for the essentially cyclic update rule and almost surely for the randomized

update rule)

$$\begin{aligned} 0 &= \nabla g_k(x_k^*) + y_k^*, \quad k = 1, \dots, K. \\ x_0^* &\in \arg \min_{x \in X} h(x) + \sum_{k=1}^K \langle y_k^*, x_k^* - x \rangle \\ x_k^* &= x_0^*, \quad k = 1, \dots, K. \end{aligned}$$

That is, any limit point of Algorithm 2 is a stationary solution of problem (2.2).

3. If  $X$  is a compact set, then the sequence of iterates generated by Algorithm 2 converges to the set of stationary solutions of problem (2.2). That is,

$$\lim_{t \rightarrow \infty} \text{dist}((\{x_k^t\}, x_0^t, y^t); Z^*) = 0, \quad (2.25)$$

where  $Z^*$  is the set of primal-dual stationary solutions of problem (2.2);  $\text{dist}(x; Z^*)$  denotes the distance between a vector  $x$  and the set  $Z^*$ , i.e.,

$$\text{dist}(x; Z^*) = \min_{\hat{x} \in Z^*} \|x - \hat{x}\|.$$

**Proof.** We first show part (1) of the theorem. For the essentially cyclic update rule, Lemma 2.2 implies that

$$\begin{aligned} &L(\{x_k^{t+T}\}, x_0^{t+T}; y^{t+T}) - L(\{x_k^t\}, x_0^t; y^t) \\ &\leq \sum_{i=1}^T \sum_{k \neq 0, k \in \mathcal{C}^{t+i}} \left( \frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2} \right) \|x_k^{t+i} - x_k^{t+i-1}\|^2 - \frac{\gamma}{2} \|x_0^{t+i-1} - x_0^{t+i}\|^2 \\ &= \sum_{i=1}^T \sum_{k=1}^K \left( \frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2} \right) \|x_k^{t+i} - x_k^{t+i-1}\|^2 - \frac{\gamma}{2} \|x_0^{t+i-1} - x_0^{t+i}\|^2, \end{aligned}$$

where the last equality follows from the fact  $x_k^{t+i} = x_k^{t+i-1}$  if  $k \notin \mathcal{C}^{t+i}$  and  $k \neq 0$ . Using the fact that each index in  $\{0, \dots, K\}$  will be updated at least once during  $[t, t+T]$ , as well as Lemma 2.3 and the bounds for  $\rho_k$ 's in Assumption A2, we have

$$\|x_0^{t+1} - x_0^{t(0)}\| \rightarrow 0, \quad \|x_k^{t+1} - x_k^{t(k)}\| \rightarrow 0, \quad \forall k = 1, \dots, K. \quad (2.26)$$

By Lemma 2.1, we further obtain  $\|y_k^{t+1} - y_k^{t(k)}\| \rightarrow 0$  for all  $k = 1, 2, \dots, K$ . In light of the dual update step of Algorithm 2, the fact that  $\|y_k^{t+1} - y_k^{t(k)}\| \rightarrow 0$  implies that  $\|x_k^{t+1} - x_0^{t+1}\| \rightarrow 0$ .

For the randomized update rule, we can take the conditional expectation (over the choice of the blocks) on both sides of (2.18) and obtain

$$\begin{aligned}
& \mathbb{E} [L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t) \mid \{x_k^t\}, x_0^t; y^t] \\
& \leq \mathbb{E} \left[ \sum_{k \neq 0, k \in \mathcal{C}^{t+1}} \left( \frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2} \right) \|x_k^{t+1} - x_k^t\|^2 - \frac{\gamma}{2} \|x_0^{t+1} - x_0^t\|^2 \mid \{x_k^t\}, x_0^t; y^t \right] \\
& \leq \sum_{k=1}^K p_k p_0 \left( \frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2} \right) \|\hat{x}_k^{t+1} - x_k^t\|^2 - p_0 \frac{\gamma}{2} \|\hat{x}_0^{t+1} - x_0^t\|^2 \\
& \quad + \sum_{k=1}^K p_k (1 - p_0) \left( \frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2} \right) \|\tilde{x}_k^{t+1} - x_k^t\|^2 \\
& \leq p_{\min}^2 \sum_{k=1}^K \left( \frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2} \right) \|\hat{x}_k^{t+1} - x_k^t\|^2 - p_{\min} \frac{\gamma}{2} \|\hat{x}_0^{t+1} - x_0^t\|^2
\end{aligned}$$

where in the last two inequalities, we have used the fact that  $\rho_k$ 's satisfy Assumption A2, hence  $\frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2} < 0$  for all  $k$ ; the last inequality follows from the fact that  $p_k \geq p_{\min}$  for all  $k = 0, \dots, K$ . Note that by Lemma 2.3,  $L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - \underline{f} \geq 0$  for all  $t$ , where  $\underline{f}$  is defined in Assumption A3. Then let us subtract both sides of the above inequality by  $\underline{f}$ , and invoke the Supermartingale Convergence Theorem [57, Proposition 4.2]. We conclude that  $L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1})$  is convergent almost surely (a.s.), and that

$$\|\hat{x}_0^{t+1} - x_0^t\| \rightarrow 0, \quad \|\hat{x}_k^{t+1} - x_k^t\| \rightarrow 0, \quad \forall k = 1, \dots, K, \quad \text{a.s.} \quad (2.27)$$

By Lemma 2.1, we further obtain  $\|\hat{y}_k^{t+1} - y_k^t\| \rightarrow 0$ , a.s. and for all  $k = 1, 2, \dots, K$ . Finally, from the definition of  $\hat{y}^{t+1}$ , we see that  $\|\hat{y}_k^{t+1} - y_k^t\| \rightarrow 0$  a.s. implies that  $\|\hat{x}_0^{t+1} - \hat{x}_k^{t+1}\| \rightarrow 0$  a.s. for all  $k = 1, 2, \dots, K$ .

Next we show part (2) of the theorem. For simplicity, we consider only the essentially cyclic rule as the proof for the randomized rule is similar. We begin by examining the optimality condition for the  $x_k$  and  $x_0$  subproblems at iteration  $t+1$ . Suppose  $k \neq 0$ ,  $k \in \mathcal{C}^{t+1}$ , then we have

$$\nabla g_k(x_k^{t+1}) + y_k^t + \rho_k(x_k^{t+1} - x_0^{t+1}) = 0. \quad (2.28)$$

Similarly, suppose  $0 \in \mathcal{C}^{t+1}$ , then there exists an  $\eta^{t+1} \in \partial h(x_0^{t+1})$  such that

$$\left\langle x - x_0^{t+1}, \eta^{t+1} - \sum_{k=1}^K (y_k^t - \rho_k(x_0^{t+1} - x_0^t)) \right\rangle \geq 0, \quad \forall x \in X.$$

These expressions imply that

$$\begin{aligned}
& \nabla g_k(x_k^{t+1}) + y_k^t + \rho_k(x_k^{t+1} - x_0^{t+1}) = 0, \quad k \neq 0, \quad k \in \mathcal{C}^{t+1} \\
& h(x) - h(x_0^{t+1}) + \left\langle x - x_0^{t+1}, \sum_{k=1}^K (-y_k^t + \rho_k(x_0^{t+1} - x_0^t)) \right\rangle \geq 0, \quad \forall x \in X, \quad \text{if } 0 \in \mathcal{C}^{t+1}.
\end{aligned} \quad (2.29)$$

Using the definition of the essentially cyclic update rule, we have that for all  $t$

$$\begin{aligned} \nabla g_k(x_k^{r(k)}) + y_k^{r(k)} &= 0, \forall k \neq 0, \text{ for some } r(k) \in [t, t+T], \\ h(x) - h(x_0^{r(0)}) + \left\langle x - x_0^{r(0)}, \sum_{k=1}^K \left( -y_k^{r(0)-1} + \rho_k(x_0^{r(0)} - x_k^{r(0)-1}) \right) \right\rangle &\geq 0, \quad (2.30) \\ \forall x \in X, \text{ for some } r(0) \in [t, t+T]. \end{aligned}$$

Note that  $T$  is finite, and that  $\|x_k^{t+1} - x_k^t\| \rightarrow 0$ ,  $\|x_0^{t+1} - x_0^t\| \rightarrow 0$  and  $\|y_k^{t+1} - y_k^t\| \rightarrow 0$ , we have

$$\begin{aligned} \|x_k^{r(k)} - x_k^{t+1}\| &\rightarrow 0, \forall k, \quad \|x_0^{r(0)} - x_0^{t+1}\| \rightarrow 0, \\ \|y_k^{t+1} - y_k^{r(k)}\| &\rightarrow 0, \quad \|y_k^{t+1} - y_k^{r(0)-1}\| \rightarrow 0, \forall k. \end{aligned} \quad (2.31)$$

Using this result, taking limit for (2.30), and using the fact that  $\|x_k^{t+1} - x_k^t\| \rightarrow 0$ ,  $x_0^{t+1} \rightarrow x_0^*$ ,  $x_k^{t+1} \rightarrow x_k^*$ ,  $y_k^{t+1} \rightarrow y_k^*$  for all  $k$ , we have

$$\begin{aligned} \nabla g_k(x_k^*) + y_k^* &= 0, \quad k = 1, \dots, K \\ h(x) - h(x_0^*) + \sum_{k=1}^K \langle x - x_0^*, -y_k^* \rangle &\geq 0, \quad \forall x \in X. \end{aligned} \quad (2.32)$$

Due to the fact that  $\|y_k^{t+1} - y_k^t\| \rightarrow 0$  for all  $k$ , we have that the primal feasibility is achieved in the limit, i.e.,

$$x_k^* = x_0^*, \quad \forall k = 1, \dots, K. \quad (2.33)$$

This set of equalities together with (2.32) imply

$$h(x) + \sum_{k=1}^K \langle x_k^* - x, y_k^* \rangle - \left( h(x_0^*) + \sum_{k=1}^K \langle x_k^* - x_0^*, y_k^* \rangle \right) \geq 0, \quad \forall x \in X. \quad (2.34)$$

This concludes the proof of part (2).

To prove part 3, we first show that there exists a limit point for each of the sequences  $\{x_k^t\}$ ,  $\{x_0^t\}$  and  $\{y^t\}$ . Let us consider only the essentially cyclic rule. Due to the compactness assumption of  $X$ , it is obvious that  $\{x_0^t\}$  must have a limit point. Also by a similar argument leading to (2.26), we see that  $\|x_k^t - x_0^t\| \rightarrow 0$ , thus for each  $k$ ,  $x_k^t$  must also lie in a compact set thus have a limit point. Note that the Lipschitz continuity of  $\nabla g_k$  combined with the compactness of the set  $X$  implies that the set  $\{\nabla g_k(x) \mid x \in X\}$  is bounded, therefore  $\{\nabla g_k(x_k^t)\}$  is a bounded sequence. Using (2.17), we conclude that  $\{y_k^t\}$  is also a bounded sequence, therefore must have at least one limit point.

We prove part 3 by contradiction. Because the feasible set is compact, then  $\{x_k^t\}$  lies in a compact set. From the argument in the previous part it is easy to see that  $\{x_0^t\}$ ,  $\{y^t\}$  also lie in some compact sets. Then every subsequence will have a limit point. Suppose that there exists a subsequence  $\{x_k^{t_j}\}$ ,  $x_0^{t_j}$  and  $\{y^{t_j}\}$  such that

$$(\{x_k^{t_j}\}, x_0^{t_j}, y^{t_j}) \rightarrow (\{\hat{x}_k\}, \hat{x}_0, \hat{y}) \quad (2.35)$$

where  $(\{\hat{x}_k\}, \hat{x}_0, \hat{y})$  is some limit point, and by part 2, we have  $(\hat{x}_k, \hat{x}_0, \hat{y}) \in Z^*$ . By further restricting to a subsequence if necessary, we can assume that  $(\hat{x}_k, \hat{x}_0, \hat{y})$  is the unique limit point.

Suppose that this sequence does not converge to the set of stationary solutions, i.e.,

$$\lim_{j \rightarrow \infty} \text{dist} \left( (\{x_k^{t_j}\}, x_0^{t_j}, y^{t_j}); Z^* \right) = \gamma > 0. \quad (2.36)$$

Then it follows that there exists some  $J(\gamma) > 0$  such that

$$\|(\{x_k^{t_j}\}, x_0^{t_j}, y^{t_j}) - (\{\hat{x}_k\}, \hat{x}_0, \hat{y})\| \leq \gamma/2, \quad \forall j \geq J(\gamma).$$

By the definition of the distance function we have

$$\text{dist} \left( (\{x_k^{t_j}\}, x_0^{t_j}, y^{t_j}); Z^* \right) \leq \text{dist} \left( (\{x_k^{t_j}\}, x_0^{t_j}, y^{t_j}), (\{\hat{x}_k\}, \hat{x}_0, \hat{y}) \right).$$

Combining the above two inequalities we must have

$$\text{dist} \left( (\{x_k^{t_j}\}, x_0^{t_j}, y^{t_j}); Z^* \right) \leq \gamma/2, \quad \forall t_j \geq T_j(\gamma).$$

This contradicts to (2.36). The desired result is proven.  $\square$

The analysis presented above is different from the conventional analysis of the ADMM algorithm where the main effort is to bound the distance between the current iterate and the optimal solution set. The above analysis is partly motivated by our previous analysis of the convergence of ADMM for multi-block convex problems, where the progress of the algorithm is measured by the combined decrease of certain primal and dual gaps; see [27, Theorem 3.1]. Nevertheless, the nonconvexity of the problem makes it difficult to estimate either the primal or the dual optimality gaps. Therefore we choose to use the decrease of the augmented Lagrangian as a measure of the progress of the algorithm.

Next we analyze the iteration complexity of the vanilla ADMM (i.e., Algorithm 1). To state our result, let us define the *proximal gradient* of the augmented Lagrangian function as

$$\tilde{\nabla} L(\{x_k\}, x_0, y) = \begin{bmatrix} x_0 - \text{prox}_h[x_0 - \nabla_{x_0}(L(\{x_k\}, x_0, y) - h(x_0))] \\ \nabla_{x_1} L(\{x_k\}, x_0, y) \\ \vdots \\ \nabla_{x_K} L(\{x_k\}, x_0, y) \end{bmatrix} \quad (2.37)$$

where  $\text{prox}_h[z] := \arg \min_x h(x) + \frac{1}{2}\|x - z\|^2$  is the proximity operator. We will use the following quantity to measure the progress of the algorithm

$$P(\{x_k^t\}, x^t, y^t) := \|\tilde{\nabla} L(\{x_k^t\}, x_0^t, y^t)\|^2 + \sum_{k=1}^K \|x_k^t - x_0^t\|^2.$$

It can be verified that if  $P(\{x_k^t\}, x^t, y^t) \rightarrow 0$ , then a stationary solution of the problem (2.2) is obtained. We have the following iteration complexity result.

**THEOREM 2.5.** *Suppose Assumption A is satisfied. Let  $T(\epsilon)$  denote an iteration index in which the following inequality is achieved*

$$T(\epsilon) := \min \{t \mid P(\{x_k^t\}, x^t, y^t) \leq \epsilon, t \geq 0\}$$

for some  $\epsilon > 0$ . Then there exists some constant  $C > 0$  such that

$$\epsilon \leq \frac{C(L(\{x_k^1\}, x_0^1, y^1) - f)}{T(\epsilon)}. \quad (2.38)$$

where  $\underline{f}$  is defined in Assumption A3.

**Proof.** We first show that there exists a constant  $\sigma_1 > 0$  such that

$$\|\tilde{\nabla}L(\{x_k^t\}, x_0^t, y^t)\| \leq \sigma_1 \left( \|x_0^{t+1} - x_0^t\| + \sum_{k=1}^K \|x_k^{t+1} - x_k^t\| \right), \quad \forall r \geq 1. \quad (2.39)$$

This proof follows similar steps of [27, Lemma 2.5]. From the optimality condition of the  $x_0$  update step (2.10) we have

$$x_0^{t+1} = \text{prox}_h \left[ x_0^{t+1} - \sum_{k=1}^K \rho_k \left( x_0^{t+1} - x_k^t - \frac{y_k^t}{\rho_k} \right) \right].$$

This implies that

$$\begin{aligned} & \|x_0^t - \text{prox}_h [x_0^t - \nabla_{x_0}(L(\{x_k^t\}, x_0^t, y^t) - h(x_0^t))] \| \\ &= \left\| x_0^t - x_0^{t+1} + x_0^{t+1} - \text{prox}_h \left[ x_0^t - \sum_{k=1}^K \rho_k \left( x_0^t - x_k^t - \frac{y_k^t}{\rho_k} \right) \right] \right\| \\ &\leq \|x_0^t - x_0^{t+1}\| + \left\| \text{prox}_h \left[ x_0^{t+1} - \sum_{k=1}^K \rho_k \left( x_0^{t+1} - x_k^t - \frac{y_k^t}{\rho_k} \right) \right] \right. \\ &\quad \left. - \text{prox}_h \left[ x_0^t - \sum_{k=1}^K \rho_k \left( x_0^t - x_k^t - \frac{y_k^t}{\rho_k} \right) \right] \right\| \\ &\leq 2\|x_0^{t+1} - x_0^t\| + \sum_{k=1}^K \rho_k \|x_0^t - x_0^{t+1}\| \end{aligned} \quad (2.40)$$

where in the last inequality we have used the nonexpansiveness of the proximity operator.

Similarly, the optimality condition of the  $x_k$  subproblem is given by

$$\nabla g_k(x_k^{t+1}) + \rho_k \left( x_k^{t+1} - x_0^{t+1} + \frac{y_k^t}{\rho_k} \right) = 0.$$

Therefore we have

$$\begin{aligned} & \|\nabla_{x_k} L(\{x_k^t\}, x_0^t, y^t)\| \\ &= \|\nabla g_k(x_k^t) + \rho_k(x_k^t - x_0^t + \frac{y_k^t}{\rho_k})\| \\ &= \left\| \left( \nabla g_k(x_k^t) + \rho_k(x_k^t - x_0^t + \frac{y_k^t}{\rho_k}) \right) - \left( \nabla g_k(x_k^{t+1}) + \rho_k(x_k^{t+1} - x_0^{t+1} + \frac{y_k^t}{\rho_k}) \right) \right\| \\ &\leq (L_k + \rho_k) \|x_k^t - x_k^{t+1}\| + \rho_k \|x_0^t - x_0^{t+1}\|. \end{aligned} \quad (2.41)$$

Therefore, combining (2.40) and (2.41), we have

$$\|\tilde{\nabla}L(\{x_k^t\}, x_0^t, y^t)\| \leq \left( 2 + \sum_{k=1}^K 2\rho_k \right) \|x_0^t - x_0^{t+1}\| + \sum_{k=1}^K (L_k + \rho_k) \|x_k^t - x_k^{t+1}\|. \quad (2.42)$$

By taking  $\sigma_1 = \max \left\{ (2 + \sum_{k=1}^K 2\rho_k), L_1 + \rho_1, \dots, L_K + \rho_K \right\}$ , (2.39) is proved.

According to Lemma 2.1, we have

$$\sum_{k=1}^K \|x_k^t - x_0^t\| = \sum_{k=1}^K \frac{1}{\rho_k} \|y_k^{t+1} - y_k^t\| \leq \sum_{k=1}^K \frac{L_k}{\rho_k} \|x_k^{t+1} - x_k^t\|. \quad (2.43)$$

The inequalities (2.42) – (2.43) implies that for some  $\sigma_3 > 0$

$$\begin{aligned} & \sum_{k=1}^K \|x_k^t - x_0^t\|^2 + \|\tilde{\nabla} L(\{x_k^t\}, x_0^t, y^t)\|^2 \\ & \leq \sigma_3 \left( \|x_0^t - x_0^{t+1}\|^2 + \sum_{k=1}^K \|x_k^t - x_k^{t+1}\|^2 \right). \end{aligned} \quad (2.44)$$

According to Lemma 2.2, there exists a constant  $\sigma_2 = \min \left\{ \left\{ \frac{\gamma_k(\rho_k)}{2} - \frac{L_k^2}{\rho_k} \right\}_{k=1}^K, \frac{\gamma}{2} \right\}$  such that

$$\begin{aligned} & L(\{x_k^t\}, x_0^t; y^t) - L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) \\ & \geq \sigma_2 \left( \sum_{k=1}^K \|x_k^{t+1} - x_k^t\|^2 + \|x_0^{t+1} - x_0^t\|^2 \right). \end{aligned} \quad (2.45)$$

Combining (2.44) and (2.45) we have

$$\begin{aligned} & \sum_{k=1}^K \|x_k^t - x_0^t\|^2 + \|\tilde{\nabla} L(\{x_k^t\}, x_0^t, y^t)\|^2 \\ & \leq \frac{\sigma_3}{\sigma_2} (L(\{x_k^t\}, x_0^t; y^t) - L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1})). \end{aligned}$$

Summing both sides of the above inequality over  $t = 1, \dots, r$ , we have

$$\begin{aligned} & \sum_{t=1}^r \sum_{k=1}^K \|x_k^t - x_0^t\|^2 + \|\tilde{\nabla} L(\{x_k^t\}, x_0^t, y^t)\|^2 \\ & \leq \frac{\sigma_3}{\sigma_2} (L(\{x_k^1\}, x_0^1; y^1) - L(\{x_k^{r+1}\}, x_0^{r+1}; y^{r+1})) \\ & \leq \frac{\sigma_3}{\sigma_2} (L(\{x_k^1\}, x_0^1; y^1) - \underline{f}) \end{aligned}$$

where in the last inequality we have used the fact that  $L(\{x_k^{r+1}\}, x_0^{r+1}; y^{r+1})$  is decreasing and lower bounded by  $\underline{f}$  (cf. Lemmas 2.2–2.3).

By utilizing the definition of  $T(\epsilon)$  and  $P(\{x_k^t\}, x^t, y^t)$ , the above inequality becomes

$$T(\epsilon)\epsilon \leq \frac{\sigma_3}{\sigma_2} (L(\{x_k^1\}, x_0^1; y^1) - \underline{f}) \quad (2.46)$$

Dividing both sides by  $T(\epsilon)$ , and by setting  $C = \sigma_3/\sigma_2$ , the desired result is obtained.

□

**2.3. The Proximal ADMM.** One potential limitation of Algorithms 1 and 2 is the requirement that each subproblem (2.11) needs to be solved *exactly*, while in certain practical applications cheap iterations are preferred. In this section, we consider an important extension of Algorithm 1–2 in which the above restriction is removed. The main idea is to take a proximal step instead of minimizing the augmented Lagrangian function exactly with respect to each variable block. Like in the previous section, we will analyze a generalized version, termed the *flexible* proximal ADMM, where there is more freedom in choosing the update schedules.

**Algorithm 3. A Flexible Proximal ADMM for Problem (2.2)**

At each iteration  $t + 1$ , compute:

$$x_0^{t+1} = \operatorname{argmin}_{x \in X} L(\{x_k^t\}, x_0; y^t). \quad (2.47)$$

Pick a set  $\mathcal{C}^{t+1} \subseteq \{1, \dots, K\}$ .

If  $k \in \mathcal{C}^{t+1}$ , update  $x_k$  by solving:

$$x_k^{t+1} = \operatorname{argmin}_{x_k} \langle \nabla g_k(x_0^{t+1}), x_k - x_0^{t+1} \rangle + \langle y_k^t, x_k - x_0^{t+1} \rangle + \frac{\rho_k + L_k}{2} \|x_k - x_0^{t+1}\|^2. \quad (2.48)$$

Update the dual variable:

$$y_k^{t+1} = y_k^t + \rho_k (x_k^{t+1} - x_0^{t+1}). \quad (2.49)$$

Else let  $x_k^{t+1} = x_k^t, y_k^{t+1} = y_k^t$ .

Notice that the  $x_k$  update step is different from the conventional proximal update (e.g., [8]). In particular, the linearization is done with respect to  $x_0^{t+1}$  instead of  $x_k^t$  computed in the previous iteration. This modification is instrumental in the convergence analysis of Algorithm 3.

Here we use the *period- $T$  essentially cyclic rule* to decide the set  $\mathcal{C}^{t+1}$  at each iteration. We note that there is a slight difference of the update schedule used in Algorithm 3 and Algorithm 2. In Algorithm 3, the block variable  $x_0$  is updated *in every iteration* while in Algorithm 2 the update of  $x_0$  is also governed by block selection rules.

Now we begin analyzing Algorithm 3. We make the following assumptions in this section (in addition to Assumption A1 and A3).

**Assumption B.** For all  $k$ , the penalty parameter  $\rho_k$  is chosen large enough such that:

$$\alpha_k := \frac{\rho_k - 7L_k}{2} - \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) 2L_k^2 > 0 \quad (2.50)$$

$$\beta_k := \frac{\rho_k}{2} - T^2 \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) 8L_k^2 > 0 \quad (2.51)$$

$$\rho_k \geq 5L_k, \quad k = 1, \dots, K. \quad (2.52)$$

Again let  $t(k)$  denote the last iteration that  $x_k$  is updated before  $t + 1$ , i.e.,

$$t(k) = \max \{r \mid r \leq t, k \in \mathcal{C}^r\}, \quad k = 1, \dots, K. \quad (2.53)$$

Note that we do not need  $t(0)$  anymore since  $x_0$  is updated in every iteration. Clearly, we have  $x_k^t = x_k^{t(k)}$  and as a result,  $y_k^t = y_k^{t(k)}$ . We have the following result.



LEMMA 2.6. *Suppose Assumption B and Assumptions A1, A3 are satisfied. Then for Algorithm 3, the following is true for the essentially cyclic block selection rule*

$$2L_k^2(4\|x_0^{t+1} - x_0^{t(k)}\|^2 + \|x_k^{t+1} - x_k^t\|^2) \geq \|y_k^{t+1} - y_k^t\|^2, \quad k = 1, \dots, K. \quad (2.54)$$

**Proof.** Suppose  $k \notin \mathcal{C}^{t+1}$ , then the inequality is trivially true, as  $y_k^{t+1} = y_k^t$ .

For any  $k \in \mathcal{C}^{t+1}$ , we observe from the update of  $x_k$  step (2.48) that the following is true

$$\nabla g_k(x^{t+1}) + y_k^t + (\rho_k + L_k)(x_k^{t+1} - x_0^{t+1}) = 0, \quad k \in \mathcal{C}^{t+1}, \quad (2.55)$$

or equivalently

$$\nabla g_k(x^{t+1}) + L_k(x_k^{t+1} - x_0^{t+1}) = -y_k^{t+1}, \quad k \in \mathcal{C}^{t+1}. \quad (2.56)$$

Therefore we have, for all  $k \in \mathcal{C}^{t+1}$

$$\begin{aligned} \|y_k^{t+1} - y_k^t\| &= \|y_k^{t+1} - y_k^{t(k)}\| \\ &= \|\nabla g_k(x_0^{t+1}) - \nabla g_k(x_0^{t(k)}) + L_k(x_k^{t+1} - x_0^{t+1}) - L_k(x_k^{t(k)} - x_0^{t(k)})\| \\ &= \|\nabla g_k(x_0^{t+1}) - \nabla g_k(x_0^{t(k)}) + L_k(x_k^{t+1} - x_0^{t+1}) - L_k(x_k^t - x_0^{t(k)})\| \\ &\leq L_k(2\|x_0^{t+1} - x_0^{t(k)}\| + \|x_k^{t+1} - x_k^t\|) \end{aligned}$$

where the last step follows from triangular inequality and the fact  $x_k^t = x_k^{t(k)}$  (cf. the definition of  $t(k)$ ). The above result further implies that

$$2L_k^2(4\|x_0^{t+1} - x_0^{t(k)}\|^2 + \|x_k^{t+1} - x_k^t\|^2) \geq \|y_k^{t+1} - y_k^t\|^2, \quad k = 1, \dots, K \quad (2.57)$$

which is the desired result.  $\square$

Next, we upper bound the successive difference of the augmented Lagrangian. To this end, let us define the following functions

$$\begin{aligned} \ell_k(x_k; x_0^{t+1}, y^t) &= g_k(x_k) + \langle y_k^t, x_k - x_0^{t+1} \rangle + \frac{\rho_k}{2} \|x_k - x_0^{t+1}\|^2 \\ u_k(x_k; x_0^{t+1}, y^t) &= g_k(x_0^{t+1}) + \langle \nabla g_k(x_0^{t+1}), x_k - x_0^{t+1} \rangle \\ &\quad + \langle y_k^t, x_k - x_0^{t+1} \rangle + \frac{\rho_k + L_k}{2} \|x_k - x_0^{t+1}\|^2. \end{aligned}$$

Using these short-hand definitions, we have

$$L(\{x_k^{t+1}\}, x_0^{t+1}, y^t) = \sum_{k=1}^K \ell_k(x_k^{t+1}; x_0^{t+1}, y^t) \quad (2.58)$$

$$x_k^{t+1} = \arg \min_{x_k} u_k(x_k; x_0^{t+1}, y^t), \quad \forall k \in \mathcal{C}^{t+1}. \quad (2.59)$$

The lemma below bounds the difference between  $\ell_k(x_k^{t+1}; x_0^{t+1}, y^t)$  and  $\ell_k(x_k^t; x_0^{t+1}, y^t)$ .

LEMMA 2.7. *Suppose Assumption A1 is satisfied. Let  $\{x_k^t, x_0^t, y^t\}$  be generated by Algorithm 3 with essential cyclic block update rule. Then we have the following*

$$\begin{aligned} &\ell_k(x_k^{t+1}; x_0^{t+1}, y^t) - \ell_k(x_k^t; x_0^{t+1}, y^t) \\ &\leq -\frac{\rho_k - 7L_k}{2} \|x_k^{t+1} - x_k^t\|^2 + \frac{4L_k}{\rho_k^2} \|y_k^{t+1} - y_k^t\|^2, \quad k = 1, \dots, K. \end{aligned} \quad (2.60)$$

**Proof.** When  $k \notin \mathcal{C}^{t+1}$ , the inequality is trivially true. We focus on the case  $k \in \mathcal{C}^{t+1}$ . From the definition of  $\ell_k(\cdot)$  and  $u_k(\cdot)$  we have the following

$$\ell_k(x_k; x_0^{t+1}, y^t) \leq u_k(x_k; x_0^{t+1}, y^t), \forall x_k, k = 1, \dots, K. \quad (2.61)$$

Observe that when  $k \in \mathcal{C}^{t+1}$ ,  $x_k^{t+1}$  is generated according to (2.59). Due to the strong convexity of  $u_k(x_k; x_0^{t+1}, y^t)$  with respect to  $x_k$ , we have

$$u_k(x_k^{t+1}; x_0^{t+1}, y^t) - u_k(x_k^t; x_0^{t+1}, y^t) \leq -\frac{\rho_k + L_k}{2} \|x_k^t - x_k^{t+1}\|^2, \forall k \in \mathcal{C}^{t+1}. \quad (2.62)$$

Further, we have the following series of inequalities

$$\begin{aligned} & u_k(x_k^t; x_0^{t+1}, y^t) - \ell_k(x_k^t; x_0^{t+1}, y^t) \\ &= g_k(x_0^{t+1}) + \langle \nabla g_k(x_0^{t+1}), x_k^t - x_0^{t+1} \rangle + \langle y_k^t, x_k^t - x_0^{t+1} \rangle + \frac{\rho_k + L_k}{2} \|x_k^t - x_0^{t+1}\|^2 \\ &\quad - \left( g_k(x_k^t) + \langle y_k^t, x_k^t - x_0^{t+1} \rangle + \frac{\rho_k}{2} \|x_k^t - x_0^{t+1}\|^2 \right) \\ &= g_k(x_0^{t+1}) - g_k(x_k^t) + \langle \nabla g_k(x_0^{t+1}), x_k^t - x_0^{t+1} \rangle + \frac{L_k}{2} \|x_k^t - x_0^{t+1}\|^2 \\ &\leq \langle \nabla g_k(x_0^{t+1}) - \nabla g_k(x_k^t), x_k^t - x_0^{t+1} \rangle + L_k \|x_k^t - x_0^{t+1}\|^2 \\ &\leq 2L_k \|x_k^t - x_0^{t+1}\|^2 \leq 4L_k (\|x_k^t - x_k^{t+1}\|^2 + \|x_k^{t+1} - x_0^{t+1}\|^2), \end{aligned} \quad (2.63)$$

where the first two inequalities follow from Assumption A1. Combining (2.61) – (2.63) we obtain

$$\begin{aligned} & \ell_k(x_k^{t+1}; x_0^{t+1}, y^t) - \ell_k(x_k^t; x_0^{t+1}, y^t) \\ &\leq u_k(x_k^{t+1}; x_0^{t+1}, y^t) - u_k(x_k^t; x_0^{t+1}, y^t) + u_k(x_k^t; x_0^{t+1}, y^t) - \ell_k(x_k^t; x_0^{t+1}, y^t) \\ &\leq -\frac{\rho_k - 7L_k}{2} \|x_k^t - x_k^{t+1}\|^2 + 4L_k \|x_k^{t+1} - x_0^{t+1}\|^2 \\ &= -\frac{\rho_k - 7L_k}{2} \|x_k^t - x_k^{t+1}\|^2 + \frac{4L_k}{\rho_k^2} \|y_k^{t+1} - y_k^t\|^2, \forall k \in \mathcal{C}^{t+1}. \end{aligned}$$

The desired result then follows.  $\square$

Next, we bound the difference of the augmented Lagrangian function values.

LEMMA 2.8. Assume the same set up as in Lemma 2.7. Then we have

$$\begin{aligned} & L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^1\}, x_0^1; y^1) \\ &\leq -\sum_{i=1}^t \sum_{k=1}^K \alpha_k \|x_k^{i+1} - x_k^i\|^2 - \sum_{i=1}^t \sum_{k=1}^K \beta_k \|x_0^{i+1} - x_0^i\|^2 \end{aligned} \quad (2.64)$$

where we  $\beta_k$  and  $\alpha_k$  are the positive constants defined in (2.50) and (2.51).

**Proof.** We first bound the successive difference  $L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t)$ .

Again we decompose it as in (2.19), and bound the resulting two differences separately.

The first term in (2.19) can be again expressed as

$$L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^{t+1}\}, x_0^{t+1}; y^t) = \sum_{k=1}^K \frac{1}{\rho_k} \|y_k^{t+1} - y_k^t\|^2.$$

To bound the second term in (2.19), we use Lemma 2.7. We use an argument similar to the proof of (2.21) to obtain

$$\begin{aligned}
& L(\{x_k^{t+1}\}, x_0^{t+1}; y^t) - L(\{x_k^t\}, x_0^t; y^t) \\
&= L(\{x_k^{t+1}\}, x_0^{t+1}; y^t) - L(\{x_k^t\}, x_0^{t+1}; y^t) + L(\{x_k^t\}, x_0^{t+1}; y^t) - L(\{x_k^t\}, x_0^t; y^t) \\
&= \sum_{k=1}^K (\ell_k(x_k^{t+1}; x_0^{t+1}, y^t) - \ell_k(x_k^t; x_0^{t+1}, y^t)) + L(\{x_k^t\}, x_0^{t+1}; y^t) - L(\{x_k^t\}, x_0^t; y^t) \\
&\leq -\sum_{k=1}^K \left( \frac{\rho_k - 7L_k}{2} \|x_k^{t+1} - x_k^t\|^2 - \frac{4L_k}{\rho_k^2} \|y_k^{t+1} - y_k^t\|^2 \right) - \frac{\gamma}{2} \|x_0^{t+1} - x_0^t\|^2 \quad (2.65)
\end{aligned}$$

where the last inequality follows from Lemma 2.7 and the strong convexity of  $L(\{x_k^t\}, x_0; y^t)$  with respect to the variable  $x$  (with modulus  $\gamma = \sum_{k=1}^K \rho_k$ ) at  $x_0 = x_0^{t+1}$ .

Combining the above two inequalities, we obtain

$$\begin{aligned}
& L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t) \\
&\leq \sum_{k=1}^K \left( -\frac{\rho_k - 7L_k}{2} \|x_k^{t+1} - x_k^t\|^2 + \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) \|y_k^{t+1} - y_k^t\|^2 \right) - \frac{\gamma}{2} \|x_0^{t+1} - x_0^t\|^2 \\
&\stackrel{(a)}{\leq} \sum_{k=1}^K \left( -\frac{\rho_k - 7L_k}{2} \|x_k^{t+1} - x_k^t\|^2 + \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) 2L_k^2 (4\|x_0^{t+1} - x_0^{t(k)}\|^2 + \|x_k^{t+1} - x_k^t\|^2) \right) \\
&\quad - \frac{\gamma}{2} \|x_0^{t+1} - x_0^t\|^2 \\
&\stackrel{(b)}{=} -\sum_{k=1}^K \left( \frac{\rho_k - 7L_k}{2} - \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) 2L_k^2 \right) \|x_k^{t+1} - x_k^t\|^2 - \sum_{k=1}^K \left( \frac{\rho_k}{2} \right) \|x_0^{t+1} - x_0^t\|^2 \\
&\quad + \sum_{k=1}^K \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) 8L_k^2 \|x_0^{t(k)} - x_0^{t+1}\|^2 \\
&\leq -\sum_{k=1}^K \left( \frac{\rho_k - 7L_k}{2} - \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) 2L_k^2 \right) \|x_k^{t+1} - x_k^t\|^2 - \sum_{k=1}^K \left( \frac{\rho_k}{2} \right) \|x_0^{t+1} - x_0^t\|^2 \\
&\quad + \sum_{k=1}^K T \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) 8L_k^2 \sum_{i=0}^{\min\{T-1, t-1\}} \|x_0^{t-i+1} - x_0^{t-i}\|^2 \quad (2.66)
\end{aligned}$$

where in (a) we have used (2.54); in (b) we have used the fact that  $\gamma = \sum_{k=1}^K \rho_k$ ; in the last inequality we have used the definition of the period- $T$  essentially cyclic update rule which implies that

$$\begin{aligned}
& \|x_0^{t+1} - x_0^{t(k)}\| \leq \sum_{i=0}^{\min\{T-1, t-1\}} \|x_0^{t-i+1} - x_0^{t-i}\| \\
&\implies \|x_0^{t+1} - x_0^{t(k)}\|^2 \leq T \sum_{i=0}^{\min\{T-1, t-1\}} \|x_0^{t-i+1} - x_0^{t-i}\|^2. \quad (2.67)
\end{aligned}$$

Then for any given  $t$ , the difference  $L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^1\}, x_0^1; y^1)$  is obtained by summing (2.66) over all iterations. Specifically, we obtain

$$\begin{aligned}
& L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^1\}, x_0^1; y^1) \\
& \leq - \sum_{i=1}^t \sum_{k=1}^K \left( \frac{\rho_k - 7L_k}{2} - \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) 2L_k^2 \right) \|x_k^{i+1} - x_k^i\|^2 \\
& \quad - \sum_{i=1}^t \sum_{k=1}^K \left( \frac{\rho_k}{2} - T^2 \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) 8L_k^2 \right) \|x_0^{i+1} - x_0^i\|^2 \\
& = - \sum_{i=1}^t \sum_{k=1}^K \alpha_k \|x_k^{i+1} - x_k^i\|^2 - \sum_{i=1}^t \sum_{k=1}^K \beta_k \|x_0^{i+1} - x_0^i\|^2.
\end{aligned}$$

This completes the proof.  $\square$

We conclude that to make the rhs of (2.64) negative at each iteration, it is sufficient to require that  $\alpha_k > 0$  and  $\beta_k > 0$  for all  $k$ , or more specifically:

$$\begin{aligned}
& \frac{\rho_k - 7L_k}{2} - \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) 2L_k^2 > 0, \quad k = 1, \dots, K, \\
& \frac{\rho_k}{2} - T^2 \left( \frac{4L_k}{\rho_k^2} + \frac{1}{\rho_k} \right) 8L_k^2 > 0, \quad k = 1, \dots, K.
\end{aligned} \tag{2.68}$$

Note that one can always find a set of  $\rho_k$ 's large enough such that the above condition is satisfied.

Next we show that  $L(\{x_k^t\}, x_0^t; y^t)$  is convergent.

**LEMMA 2.9.** *Suppose Assumption A1, A3 and Assumption B are satisfied. Then Algorithm 3 with period- $T$  essentially cyclic update rule generates a sequence of augmented Lagrangian, whose limit exists and is bounded below by  $\underline{f}$ .*

**Proof.** Observe that the augmented Lagrangian can be expressed as

$$\begin{aligned}
& L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) \\
& = h(x_0^{t+1}) + \sum_{k=1}^K \left( g_k(x_k^{t+1}) + \langle y_k^{t+1}, x_k^{t+1} - x_0^{t+1} \rangle + \frac{\rho_k}{2} \|x_k^{t+1} - x_0^{t+1}\|^2 \right) \\
& \stackrel{(a)}{=} h(x_0^{t+1}) + \sum_{k=1}^K \left( g_k(x_k^{t+1}) + \langle \nabla g_k(x_0^{t+1}) + L_k(x_k^{t+1} - x_0^{t+1}), x_0^{t+1} - x_k^{t+1} \rangle + \frac{\rho_k}{2} \|x_k^{t+1} - x_0^{t+1}\|^2 \right) \\
& = h(x_0^{t+1}) + \sum_{k=1}^K \left( g_k(x_k^{t+1}) + \langle \nabla g_k(x_0^{t+1}), x_0^{t+1} - x_k^{t+1} \rangle + \frac{\rho_k - 2L_k}{2} \|x_k^{t+1} - x_0^{t+1}\|^2 \right) \\
& \stackrel{(b)}{\geq} h(x_0^{t+1}) + \sum_{k=1}^K \left( g_k(x_0^{t+1}) + \frac{\rho_k - 5L_k}{2} \|x_k^{t+1} - x_0^{t+1}\|^2 \right) \\
& = f(x_0^{t+1}) + \sum_{k=1}^K \frac{\rho_k - 5L_k}{2} \|x_k^{t+1} - x_0^{t+1}\|^2
\end{aligned} \tag{2.69}$$

where (a) is from (2.56); (b) is due to the following inequalities

$$\begin{aligned}
g_k(x_0^{t+1}) &\leq g_k(x_k^{t+1}) + \langle \nabla g_k(x_k^{t+1}), x_0^{t+1} - x_k^{t+1} \rangle + \frac{L_k}{2} \|x_k^{t+1} - x_0^{t+1}\|^2 \\
&= g_k(x_k^{t+1}) + \langle \nabla g_k(x_k^{t+1}) - \nabla g_k(x_0^{t+1}), x_0^{t+1} - x_k^{t+1} \rangle \\
&\quad + \langle \nabla g_k(x_0^{t+1}), x_0^{t+1} - x_k^{t+1} \rangle + \frac{L_k}{2} \|x_k^{t+1} - x_0^{t+1}\|^2 \\
&\leq g_k(x_k^{t+1}) + \langle \nabla g_k(x_0^{t+1}), x_0^{t+1} - x_k^{t+1} \rangle + \frac{3L_k}{2} \|x_k^{t+1} - x_0^{t+1}\|^2.
\end{aligned}$$

Clearly, combining the inequality (2.69) with Assumptions B and A3 yields that  $L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1})$  is lower bounded. It follows from Lemma 2.8 that whenever the penalty parameter  $\rho_k$ 's are chosen sufficiently large (as per Assumption B),  $L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1})$  will monotonically decrease and is convergent. This completes the proof.  $\square$

Using Lemmas 2.6–2.9, we arrive at the following convergence result. The proof is similar to Theorem 2.4, and is thus omitted.

**THEOREM 2.10.** *Suppose that Assumptions A1, A3 and B hold. Then the following is true for Algorithm 3.*

1. We have  $\lim_{t \rightarrow \infty} \|x_0^{t+1} - x_k^{t+1}\| = 0$ ,  $k = 1, \dots, K$ .
2. Let  $(\{x_k^*\}, x_0^*, y^*)$  denote any limit point of the sequence  $\{\{x_k^{t+1}\}, x_0^{t+1}, y^{t+1}\}$  generated by Algorithm 3 with period- $T$  essentially cyclic block update rule. Then  $(\{x_k^*\}, x_0^*, y^*)$  is a stationary solution of problem (2.2).
3. If  $X$  is a compact set, then Algorithm 3 with period- $T$  essentially cyclic block update rule converges to the set of stationary solutions of problem (2.2). That is, the following is true

$$\lim_{t \rightarrow \infty} \text{dist}(\{x_k^t\}, x_0^t, y^t; Z^*) = 0. \quad (2.70)$$

where  $Z^*$  is the set of primal-dual stationary solutions of problem (2.2).

**3. The Nonconvex Sharing Problem.** Consider the following well-known sharing problem (see, e.g., [8, Section 7.3] for motivation)

$$\begin{aligned}
\min \quad & f(x_1, \dots, x_K) := \sum_{k=1}^K g_k(x_k) + \ell\left(\sum_{k=1}^K A_k x_k\right) \\
\text{s.t.} \quad & x_k \in X_k, \quad k = 1, \dots, K,
\end{aligned} \quad (3.1)$$

where  $x_k \in \mathbb{R}^{N_k}$  is the variable associated with a given agent  $k$ , and  $A_k \in \mathbb{R}^{M \times N_k}$  is some data matrix. The variables are coupled through the function  $\ell(\cdot)$ .

To facilitate distributed computation, this problem can be equivalently formulated into a linearly constrained problem by introducing an additional variable  $x_0 \in \mathbb{R}^M$ :

$$\begin{aligned}
\min \quad & \sum_{k=1}^K g_k(x_k) + \ell(x_0) \\
\text{s.t.} \quad & \sum_{k=1}^K A_k x_k = x_0, \quad x_k \in X_k, \quad k = 1, \dots, K.
\end{aligned} \quad (3.2)$$

The augmented Lagrangian for this problem is given by

$$L(\{x_k\}, x_0; y) = \sum_{k=1}^K g_k(x_k) + \ell(x_0) + \left\langle x_0 - \sum_{k=1}^K A_k x_k, y \right\rangle + \frac{\rho}{2} \left\| x_0 - \sum_{k=1}^K A_k x_k \right\|^2. \quad (3.3)$$

Note that we have chosen a special reformulation in (3.2): a *single* variable  $x_0$  is introduced which leads to a problem with a *single* linear constraint. Applying the classical ADMM to this reformulation leads to a *multi-block* ADMM algorithm in which  $K + 1$  block variables  $(\{x_k\}_{k=1}^K, x_0)$  are updated sequentially. As mentioned in the introduction, even in the case where the objective is convex, it is not known whether the multi-block ADMM converges in this case. Variants of the multi-block ADMM has been proposed in the literature to solve this type of multi-block problems; see recent developments in [26–30] and the references therein.

In this section, we show that the classical ADMM, together with several of its extensions using different block selection rules, converge even when the objective function is nonconvex. The main assumptions for convergence are that the penalty parameter  $\rho$  is large enough, and that the coupling function  $\ell(x_0)$  should be smooth (more detailed conditions will be given shortly). Similarly as in the previous sections, we consider a generalized version of ADMM with two types of block update rules: the period- $T$  essentially cyclic rule and the randomized rule. The detailed algorithm is given in the table below.

<p style="text-align: center;"><b>Algorithm 4. The Flexible ADMM for Problem (3.2)</b></p> <p>Let <math>\mathcal{C}^1 = \{0, \dots, K\}</math>, <math>t = 0, 1, \dots</math>.  At each iteration <math>t + 1</math>, do:</p> <p><b>If</b> <math>t + 1 \geq 2</math>, pick an index set <math>\mathcal{C}^{t+1} \subseteq \{0, \dots, K\}</math>.</p> <p><b>For</b> <math>k = 1, \dots, K</math>  <b>If</b> <math>k \in \mathcal{C}^{t+1}</math>, then agent <math>k</math> updates <math>x_k</math> by:</p> $x_k^{t+1} = \arg \min_{x_k \in X_k} g_k(x_k) - \langle y^t, A_k x_k \rangle + \frac{\rho}{2} \left\  x_0^t - \sum_{j < k} A_j x_j^{t+1} - \sum_{j > k} A_j x_j^t - A_k x_k \right\ ^2 \quad (3.4)$ <p><b>Else</b> <math>x_k^{t+1} = x_k^t</math>.  <b>If</b> <math>0 \in \mathcal{C}^{t+1}</math>, update the variable <math>x_0</math> by:</p> $x_0^{t+1} = \arg \min_x \ell(x_0) + \langle y^t, x_0 \rangle + \frac{\rho}{2} \left\  x_0 - \sum_{k=1}^K A_k x_k^{t+1} \right\ ^2. \quad (3.5)$ <p>Update the dual variable:</p> $y^{t+1} = y^t + \rho \left( x_0^{t+1} - \sum_{k=1}^K A_k x_k^{t+1} \right). \quad (3.6)$ <p><b>Else</b> <math>x_0^{t+1} = x_0^t</math>, <math>y^{t+1} = y^t</math>.</p>
---

The analysis of Algorithm 4 follows similar argument as that of Algorithm 3. Therefore we will only provide an outline for it.

First, we make the following assumptions in this section.

**Assumption C.**

C1. There exists a positive constant  $L > 0$  such that

$$\|\nabla\ell(x) - \nabla\ell(z)\| \leq L\|x - z\|, \forall x, z.$$

Moreover,  $X_k$ 's are closed convex sets; each  $A_k$  is full column rank so that  $\lambda_{\min}(A_k^T A_k) > 0$ , where  $\lambda_{\min}$  denotes the minimum eigenvalue of a matrix.

C2. The penalty parameter  $\rho$  is chosen large enough such that:

- (1) Each  $x_k$  subproblem (3.4) as well as the  $x_0$  subproblem (3.5) is strongly convex, with modulus  $\{\gamma_k(\rho)\}_{k=1}^K$  and  $\gamma(\rho)$ , respectively.
- (2)  $\rho\gamma(\rho) > 2L^2$ , and that  $\rho \geq L$ .

C3.  $f(x_1, \dots, x_K)$  is lower bounded over  $\prod_{k=1}^K X_k$ .

C4.  $g_k$  is either smooth nonconvex or convex (possibly nonsmooth). For the former case, there exists  $L_k > 0$  such that  $\|\nabla g_k(x_k) - \nabla g_k(z_k)\| \leq L_k\|x_k - z_k\|$ ,  $\forall x_k, z_k \in X_k$ .

Note that compared with Assumptions A and B, in this case we no longer require that each  $g_k$  to be smooth. Define an index set  $\mathcal{K} \subseteq \{1, \dots, K\}$ , such that  $g_k$  is convex if  $k \in \mathcal{K}$ , and nonconvex smooth otherwise. Further, the requirement that  $A_k$  is full column rank is needed to make the  $x_k$  subproblem (3.4) strongly convex.

Our convergence analysis consists of a series of lemmas whose proofs, for the most part, are omitted since they are similar to that of Lemma 2.1–Lemma 2.3.

LEMMA 3.1. *Suppose Assumption C is satisfied. Then for Algorithm 4 with either essentially cyclic rule or the randomized rule, the following is true*

$$\begin{aligned} \nabla\ell(x_0^{t+1}) &= -y^{t+1}, \text{ if } 0 \in \mathcal{C}^{t+1}, \quad L^2\|x_0^{t+1} - x_0^t\|^2 \geq \|y^{t+1} - y^t\|^2, \\ L^2\|x_0^{t+1} - x_0^{t(k)}\|^2 &\geq \|y^{t+1} - y^{t(k)}\|^2, \quad L^2\|\hat{x}_0^{t+1} - x_0^t\|^2 \geq \|\hat{y}^{t+1} - y^t\|^2. \end{aligned}$$

LEMMA 3.2. *Suppose Assumption C is satisfied. Then for Algorithm 4 with either essentially cyclic rule or the randomized rule, the following is true*

$$\begin{aligned} &L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t) \\ &\leq \sum_{k \neq 0, k \in \mathcal{C}^{t+1}} -\frac{\gamma_k(\rho)}{2} \|x_k^{t+1} - x_k^t\|^2 - \left( \frac{\gamma(\rho)}{2} - \frac{L^2}{\rho} \right) \|x_0^{t+1} - x_0^t\|^2. \end{aligned} \quad (3.7)$$

LEMMA 3.3. *Assume the same set up as in Lemma 3.2. Then the following limit exists and is bounded from below*

$$\lim_{t \rightarrow \infty} L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}). \quad (3.8)$$

**Proof.** We have the following series of inequalities

$$\begin{aligned}
& L(\{x_k^{r+1}\}, x_0^{r+1}; y^{r+1}) \\
&= \sum_{k=1}^K g_k(x_k^{t+1}) + \ell(x_0^{t+1}) + \left\langle x_0^{t+1} - \sum_{k=1}^K A_k x_k^{t+1}, y^{t+1} \right\rangle + \frac{\rho}{2} \left\| x_0^{t+1} - \sum_{k=1}^K A_k x_k^{t+1} \right\|^2 \\
&= \sum_{k=1}^K g_k(x_k^{t+1}) + \ell(x_0^{t+1}) + \left\langle \sum_{k=1}^K A_k x_k^{t+1} - x_0^{t+1}, \nabla \ell(x_0^{t+1}) \right\rangle + \frac{\rho}{2} \left\| x_0^{t+1} - \sum_{k=1}^K A_k x_k^{t+1} \right\|^2 \\
&\geq \sum_{k=1}^K g_k(x_k^{t+1}) + \ell \left( \sum_{k=1}^K A_k x_k^{t+1} \right) + \frac{\rho - L}{2} \left\| x_0^{t+1} - \sum_{k=1}^K A_k x_k^{t+1} \right\|^2.
\end{aligned}$$

The last inequality comes from the fact that

$$\ell \left( \sum_{k=1}^K A_k x_k^{t+1} \right) \leq \ell(x_0^{t+1}) + \left\langle \sum_{k=1}^K A_k x_k^{t+1} - x_0^{t+1}, \nabla \ell(x_0^{t+1}) \right\rangle + \frac{L}{2} \left\| x_0^{t+1} - \sum_{k=1}^K A_k x_k^{t+1} \right\|^2.$$

Using assumptions C2.– C3. leads to the desired result.  $\square$

We note that the above result holds true deterministically even if the randomized scheme is used. The reason is that at each iteration regardless of whether  $0 \in \mathcal{C}^{t+1}$ , we have  $y^{t+1} = -\nabla \ell(x^{t+1})$  because these two variables are always updated at the same iteration. The rest of the proof is not dependent on the algorithm.

We have the following main result for the nonconvex consensus problem.

**THEOREM 3.4.** *Suppose that Assumption C holds. Then the following is true for Algorithm 4, either deterministically for the essentially cyclic update rule or almost surely for the randomized update rule.*

1. We have  $\lim_{t \rightarrow \infty} \left\| \sum_{k=1}^K A_k x_k^{t+1} - x_0^{t+1} \right\| = 0$ ,  $k = 1, \dots, K$ .
2. Let  $(\{x_k^*\}, x_0^*, y^*)$  denote any limit point of the sequence  $\{\{x_k^{t+1}\}, x_0^{t+1}, y^{t+1}\}$  generated by Algorithm 4. Then  $(\{x_k^*\}, x_0^*, y^*)$  is a stationary solution of problem (3.2) in the sense that

$$\begin{aligned}
& x_k^* \in \arg \min_{x_k \in X_k} g_k(x_k) + \langle y^*, -A_k x_k \rangle, \quad k \in \mathcal{K}, \\
& \langle x_k - x_k^*, \nabla g_k(x_k^*) - A_k^T y^* \rangle \geq 0, \quad \forall x_k \in X_k, \quad k \notin \mathcal{K}, \\
& \nabla \ell(x_0^*) + y^* = 0, \\
& \sum_{k=1}^K A_k x_k^* = x_0^*.
\end{aligned}$$

3. If  $X_k$  is a compact set for all  $k$ , then Algorithm 4 converges to the set of stationary solutions of problem (3.2), i.e.,

$$\lim_{t \rightarrow \infty} \text{dist}(\{x_k^t\}, x_0^t, y^t; Z^*) = 0, \quad (3.9)$$

where  $Z^*$  is the set of primal-dual stationary solution for problem (3.2).

The following corollary specializes the previous convergence result to the case where all  $g_k$ 's as well as  $\ell$  are convex (but not necessarily strongly convex). We emphasize that this is still a nontrivial result, since unlike [27, 29, 31, 34], we do not require the dual stepsize to be small or the  $g_k$ 's and  $\ell$  to be strongly convex. Therefore it is not known whether the classical ADMM converges for the multi-block problem (3.2), even for the convex case.



COROLLARY 3.5. *Suppose that Assumptions C1 and C3 hold, and that  $g_k$  and  $\ell$  are convex. Further, suppose that Assumption C2 is weakened with the following assumption*

*C2' The penalty parameter  $\rho$  is chosen large enough such that  $\rho > \sqrt{2}L$ . Then the flexible ADMM algorithm (i.e., Algorithm 4), converges to the set of primal dual optimal solution  $(\{x_k^*\}, x^*, y^*)$  of problem (2.2), either deterministically for the essentially cyclic update rule or almost surely for the randomized update rule.*

Similar to the consensus problem, one can extend Algorithm 4 to its proximal version. Here the benefit offered by the proximal-type algorithms is twofold: *i*) one can remove the strong convexity requirement posed in Assumption C2-(1) ; *ii*) one can allow inexact and simple update for each block variable. However, the analysis is a bit more involved, as the penalty parameter  $\rho$  as well as the proximal coefficient for each subproblem needs to be carefully bounded. Due to the fact that the analysis follows almost identical steps as those in Section 2.3, we will not present them here.

**4. Extensions.** In this paper, we analyze the behavior of the ADMM method in the absence of convexity. We show that when the penalty parameter is chosen sufficiently large, the ADMM and several of its variants converge to the set of stationary solutions for certain consensus and sharing problems.

Our analysis is based on using the augmented Lagrangian as a potential function to guide the iterate convergence. This approach may be extended to other nonconvex problems. In particular, if the following set of sufficient conditions (see Assumption D below) are satisfied, then the convergence of the ADMM is guaranteed for the nonconvex problem (1.1). It is important to note that in practice these conditions should be verified case by case for different applications, just like what we have done for the consensus and sharing problems.

#### Assumption D

- D1. The iterations are well defined, meaning the function  $L(x^t; y^t)$  is uniformly lower bounded for all  $t$ .
- D2. There exists a constant  $\sigma > 0$  such that  $\|y^{t+1} - y^t\|^2 \leq \sigma \|x^{t+1} - x^t\|^2$ , for all  $t$ .
- D3.  $g_k(\cdot)$  is either smooth nonconvex or nonsmooth convex. The coupling function  $\ell(\cdot)$  is smooth with Lipschitz continuous gradient  $L$ . Moreover,  $\ell(\cdot)$  is convex with respect to each block variable  $x_k$ , but is not necessarily jointly convex with  $x$ .  $X_k$  is a closed convex set. Problem (1.1) is feasible, that is,  $\{x \mid Ax = q\} \cap \bigcap_{k=1}^K \text{relint} X_k \neq \emptyset$ .
- D4. The penalty parameter  $\rho$  is chosen large enough such that each subproblem is strongly convex with modulus  $\gamma_k(\rho)$ , which is a nondecreasing function of  $\rho$ . Further,  $\rho\gamma_k(\rho) > 2\sigma$  for all  $k$ .

Following a similar argument leading to Theorem 2.4, we can show that as long as Assumption D is satisfied, then the primal feasibility gap  $\|q - \sum_{k=1}^K A_k x_k^{t+1}\|$  goes to zero in the limit, and that every limit point of the sequence  $\{\{x_k^{t+1}\}, x_0^{t+1}, y^{t+1}\}$  is a stationary solution of problem (1.1). A few remarks on Assumption D are in order:

1. Assumption D1 is necessary for showing convergence. Without D1, even if one is able to show that the augmented Lagrangian is decreasing, one cannot claim the convergence to stationary solutions. The reason is that the augmented Lagrangian may go to  $-\infty$ <sup>1</sup>, therefore there is no way to guarantee that

---

<sup>1</sup>In fact, it is very easy to modify the algorithm so that the augmented Lagrangian reduces at each iteration – just change the “+” in the dual update (2.12) to “-”. However, it is obvious that

the successive difference of the iterates goes to 0, or the primal feasibility is satisfied in the limit.

2. The main drawback of Assumption D is that it is made on the iterates rather than on the problem. For different linearly constrained optimization problems, one still needs to verify that these conditions are indeed valid, as we have done for the consensus and the sharing problem considered in this paper.

Here we mention one more family of problems for which Assumption D can be verified. Consider

$$\begin{aligned} \min \quad & f(x_1) + g(x_2) \\ \text{s.t.} \quad & Bx_1 + Ax_2 = c, \ x_1 \in X, \end{aligned} \quad (4.1)$$

where  $f(\cdot)$  is a convex possibly nonsmooth function;  $g(\cdot)$  is a possibly nonconvex function, and has Lipschitzian gradient with modulus  $L_g$ ;  $X \subseteq R^N$ ;  $A$  is an invertible matrix;  $g(\cdot)$  and  $f(\cdot)$  are lower bounded over the set  $X$ . Consider the following ADMM method, where the iterate generated at iteration  $t + 1$  is given by

$$\begin{aligned} x_1^{t+1} &= \arg \min_{x_1 \in X} f(x_1) + \langle Bx_1 + Ax_2^t - c, y^t \rangle + \frac{\rho}{2} \|Bx_1 + Ax_2^t - c\|^2 \\ x_2^{t+1} &= \arg \min g(x_2) + \langle Bx_1^{t+1} + Ax_2 - c, y^t \rangle + \frac{\rho}{2} \|Bx_1^{t+1} + Ax_2 - c\|^2 \\ y^{t+1} &= y^t + \rho (Bx_1^{t+1} + Ax_2^{t+1} - c). \end{aligned}$$

By using steps in Lemma 2.1-Lemma 2.3, one can verify that if  $\rho > L_g/\lambda_{\min}(AA^T)$ , then Assumptions D1 holds true. By having  $\rho$  large enough and by using the invertibility of  $A$ , we can make the  $x_2$  subproblem strongly convex, then Assumption D4 holds true. Other assumptions can be verified along similar lines. Note that in this case the convergence can be obtained with a slightly weaker condition in which the  $x_1$  subproblem is convex but not necessarily strongly convex.

---

by doing this the dual variables will become unbounded, and the primal feasibility will never be satisfied.

## REFERENCES

- [1] M. Hong, Z.-Q. Luo, and M. Razaviyayn. On the convergence of alternating direction method of multipliers for a family of nonconvex problems. In *ICASSP 2015*, 2015.
- [2] R. Glowinski and A. Marroco. Sur l'approximation, par elements finis d'ordre un, et la resolution, par penalisation-dualite, d'une classe de problemes de dirichlet non lineaires. *Revue Francaise d'Automatique, Informatique et Recherche Operationelle*, 9:41–76, 1975.
- [3] D. Gabay and B. Mercier. A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers & Mathematics with Applications*, 2:17–40, 1976.
- [4] J. Eckstein. Splitting methods for monotone operators with applications to parallel optimization. 1989. Ph.D Thesis, Operations Research Center, MIT.
- [5] J. Eckstein and D. P. Bertsekas. On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators. *Mathematical Programming*, 55(1):293–318, 1992.
- [6] R. Glowinski. *Numerical methods for nonlinear variational problems*. Springer-Verlag, New York, 1984.
- [7] D. P. Bertsekas and J. N. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*, 2nd ed. Athena Scientific, Belmont, MA, 1997.
- [8] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3, 2011.
- [9] W. Yin, S. Osher, D. Goldfarb, and J. Darbon. Bregman iterative algorithms for  $\ell_1$ -minimization with applications to compressed sensing. *SIAM Journal on Imaging Science*, 1(1):143–168, March 2008.
- [10] J. Yang, Y. Zhang, and W. Yin. An efficient TVL1 algorithm for deblurring multichannel images corrupted by impulsive noise. *SIAM Journal on Scientific Computing*, 31(4):2842–2865, 2009.
- [11] X. Zhang, M. Burger, and S. Osher. A unified primal-dual algorithm framework based on Bregman iteration. *Journal of Scientific Computing*, 46(1):20–46, 2011.
- [12] K. Scheinberg, S. Ma, and D. Goldfarb. Sparse inverse covariance selection via alternating linearization methods. In *Advanced in Neural Information Processing Systems (NIPS)*, 2010.
- [13] I. Schizas, A. Ribeiro, and G. Giannakis. Consensus in ad hoc wsns with noisy links - part i: Distributed estimation of deterministic signals. *IEEE Transactions on Signal Processing*, 56(1):350 – 364, 2008.
- [14] C. Feng, H. Xu, and B. Li. An alternating direction method approach to cloud traffic management. *submitted to IEEE/ACM Trans. Networking*, 2014.
- [15] W.-C. Liao, M. Hong, Hamid Farmanbar, Xu Li, Z.-Q. Luo, and Hang Zhang. Min flow rate maximization for software defined radio access networks. *IEEE Journal on Selected Areas in Communication*, 32(6):1282–1294, 2014.
- [16] D. P. Bertsekas. *Nonlinear Programming*, 2nd ed. Athena Scientific, Belmont, MA, 1999.
- [17] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [18] A. Nedic and A. Ozdaglar. Cooperative distributed multi-agent optimization. In *Convex Optimization in Signal Processing and Communications*. Cambridge University Press, 2009.
- [19] D. P. Bertsekas. *Constrained Optimization and Lagrange Multiplier Method*. Academic Press, 1982.
- [20] B. He and X. Yuan. On the  $O(1/n)$  convergence rate of the Douglas-Rachford alternating direction method. *SIAM Journal on Numerical Analysis*, 50(2):700–709, 2012.
- [21] R. Monteiro and B. Svaiter. Iteration-complexity of block-decomposition algorithms and the alternating direction method of multipliers. *SIAM Journal on Optimization*, 23(1):475–507, 2013.
- [22] D. Davis and W. Yin. Convergence rate analysis of several splitting schemes. 2014. UCLA CAM Report 14-51.
- [23] T. Goldstein, B. O'Donoghue, S. Setzer, and R. Baraniuk. Fast alternating direction optimization methods. *SIAM Journal on Imaging Sciences*, 7(3):1588–1623, 2014.
- [24] D. Goldfarb, S. Ma, and K. Scheinberg. Fast alternating linearization methods for minimizing the sum of two convex functions. *Mathematical Programming*, 141(1-2):349–382, 2012.
- [25] W. Deng and W. Yin. On the global linear convergence of alternating direction methods. 2012. preprint.
- [26] C. Chen, B. He, X. Yuan, and Y. Ye. The direct extension of ADMM for multi-block convex minimization problems is not necessarily convergent. 2013. *Mathematical Programming*,

to appear.

- [27] M. Hong and Z.-Q. Luo. On the linear convergence of the alternating direction method of multipliers. *arXiv preprint arXiv:1208.3922*, 2012.
- [28] B. He, M. Tao, and X. Yuan. Alternating direction method with Gaussian back substitution for separable convex programming. *SIAM Journal on Optimization*, 22:313–340, 2012.
- [29] M. Hong, T.-H. Chang, X. Wang, M. Razaviyayn, S. Ma, and Z.-Q. Luo. A block successive upper bound minimization method of multipliers for linearly constrained convex optimization. 2013. Preprint, available online arXiv:1401.7079.
- [30] X. Wang, M. Hong, S. Ma, and Z.-Q. Luo. Solving multiple-block separable convex minimization problems using two-block alternating direction method of multipliers. *Pacific Journal on Optimization*, 11(4):645–667, 2015.
- [31] D. Han and X. Yuan. A note on the alternating direction method of multipliers. *Journal of Optimization Theory and Applications*, 155(1):227–238, 2012.
- [32] W. Deng, M. Lai, Z. Peng, and W. Yin. Parallel multi-block ADMM with  $o(1/k)$  convergence. Preprint, available online at arXiv: 1312.3040., 2014.
- [33] B. He, H. Xu, and X. Yuan. On the proximal jacobian decomposition of alm for multipleblock separable convex minimization problems and its relationship to ADMM. 2013. Preprint, available on Optimization-Online.
- [34] T. Lin, S. Ma, and S. Zhang. On the global linear convergence of the admm with multi-block variables. 2014. Preprint.
- [35] M. Hong, T.-H. Chang, X. Wang, M. Razaviyayn, S. Ma, and Z.-Q. Luo. A block coordinate descent method of multipliers: Convergence analysis and applications. In *International Conference on Acoustics, Speech and Signal Processing*, 2014.
- [36] X. Gao and S. Zhang. First-order algorithms for convex optimization with nonseparate objective and coupled constraints. 2015. Preprint.
- [37] Y. Zhang. An alternating direction algorithm for nonnegative matrix factorization. 2010. Preprint.
- [38] D. L. Sun and C. Fevotte. Alternating direction method of multipliers for non-negative matrix factorization with the beta-divergence. In *the Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2014.
- [39] Z. Wen, C. Yang, X. Liu, and S. Marchesini. Alternating direction methods for classical and ptychographic phase retrieval. *Inverse Problems*, 28(11):1–18, 2012.
- [40] R. Zhang and J. T. Kwok. Asynchronous distributed admm for consensus optimization. In *Proceedings of the 31st International Conference on Machine Learning*, 2014.
- [41] P.A. Forero, A. Cano, and G.B. Giannakis. Distributed clustering using wireless sensor networks. *IEEE Journal of Selected Topics in Signal Processing*, 5(4):707–724, Aug 2011.
- [42] B. Ames and M. Hong. Alternating directions method of multipliers for l1-penalized zero variance discriminant analysis and principal component analysis. 2014. Preprint.
- [43] B. Jiang, S. Ma, and S. Zhang. Alternating direction method of multipliers for real and complex polynomial optimization models. 2013. Preprint.
- [44] A. P. Liavas and N. D. Sidiropoulos. Parallel algorithms for constrained tensor factorization via the alternating direction method of multipliers. 2014. Preprint, available at arXiv:1409.2383v1.
- [45] Y. Shen, Z. Wen, and Y. Zhang. Augmented lagrangian alternating direction method for matrix separation based on low-rank factorization. *Optimization Methods Software*, 29(2):239–263, March 2014.
- [46] Y. Xu, W. Yin, Z. Wen, and Y. Zhang. An alternating direction algorithm for matrix completion with nonnegative factors. *Journal of Frontiers of Mathematics in China, Special Issues on Computational Mathematics*, pages 365–384, 2011.
- [47] Z. Wen, X. Peng, X. Liu, X. Bai, and X. Sun. Asset allocation under the basel accord risk measures. 2013. Preprint.
- [48] F. Lin, M. Fardad, and M. R. Jovanovic. Design of optimal sparse feedback gains via the alternating direction method of multipliers. *IEEE Transactions on Automatic Control*, 58(9):2426–2431, Sept 2013.
- [49] Y. Zhang. Convergence of a class of stationary iterative methods for saddle point problems. 2010. Preprint.
- [50] M. Razaviyayn, M. Hong, and Z.-Q. Luo. A unified convergence analysis of block successive minimization methods for nonsmooth optimization. *SIAM Journal on Optimization*, 23(2):1126–1153, 2013.
- [51] S. Ghadimi and G. Lan. Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming*, pages 1–41, 2015.
- [52] S. Ghadimi, G. Lan, and H. Zhang. Mini-batch stochastic approximation methods for nonconvex

- stochastic composite optimization. *Mathematical Programming*, pages 1–39, 2014.
- [53] G. Scutari, F. Facchinei, P. Song, D. P. Palomar, and J.-S. Pang. Decomposition by partial linearization: Parallel optimization of multi-agent systems. *IEEE Transactions on Signal Processing*, 63(3):641–656, 2014.
  - [54] J. Bolte, S. Sabach, and M. Teboulle. Proximal alternating linearized minimization for non-convex and nonsmooth problems. *Mathematical Programming*, 146, 2014.
  - [55] E. Wei and A. Ozdaglar. On the  $O(1/k)$  convergence of asynchronous distributed alternating direction method of multipliers. 2013. Preprint, available at arXiv:1307.8254.
  - [56] T.-H. Chang. A proximal dual consensus admm method for multi-agent constrained optimization. 2014. Preprint, available at arXiv:1409.3307.
  - [57] D. P. Bertsekas and J. N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, Belmont, MA, 1996.