

First-order methods for constrained convex programming based on linearized augmented Lagrangian function*

Yangyang Xu[†]

Abstract

First-order methods have been popularly used for solving large-scale problems. However, many existing works only consider unconstrained problems or those with simple constraint. In this paper, we develop two first-order methods for constrained convex programs, for which the constraint set is represented by affine equations and smooth nonlinear inequalities. Both methods are based on the classic augmented Lagrangian function. They update the multipliers in the same way as the augmented Lagrangian method (ALM) but employ different primal variable updates. The first method, at each iteration, performs a single proximal gradient step to the primal variable, and the second method is a block update version of the first one.

For the first method, we establish its global iterate convergence as well as global sublinear and local linear convergence, and for the second method, we show a global sublinear convergence result in expectation. Numerical experiments are carried out on the basis pursuit denoising and a convex quadratically constrained quadratic program to show the empirical performance of the proposed methods. Their numerical behaviors closely match the established theoretical results.

Keywords: augmented Lagrangian method (ALM), nonlinearly constrained programming, first-order method, global convergence, iteration complexity

Mathematics Subject Classification: 90C06, 90C25, 90C30, 68W40.

1 Introduction

Recent years have witnessed the surge of first-order methods partly due to the increasingly big data involved in modern applications. Compared to second or higher-order methods, first-order ones only require gradient information and generally have much lower per-iteration complexity. However, many existing works on first-order methods are about problems without constraint or with easy-to-project constraint and/or with affine constraint.

*This work is partly supported by NSF grant DMS-1719549.

[†]xuy21@rpi.edu. Department of Mathematical Sciences, Rensselaer Polytechnic Institute, Troy, New York.

In this paper, we consider the generally constrained convex programming

$$\min_x f_0(x) \equiv g(x) + h(x), \text{ s.t. } Ax = b, f_j(x) \leq 0, j = 1, \dots, m, \quad (1)$$

where g and f_j for $j = 1, \dots, m$ are convex and Lipschitz differentiable functions, and h is a proper closed convex (possibly nondifferentiable) function. For practical efficiency of our algorithms, we will assume h to be simple in the sense that its proximal mapping is easy to compute. However, our convergence results do not require this assumption.

Applications that can be formulated into (1) appear in many areas including operations research, statistics, machine learning, engineering, just to name a few. Towards finding a solution to (1), we design algorithms that only need zeroth and first-order information of g and $f_j, j = 1, \dots, m$, and the proximal mapping of h .

1.1 Augmented Lagrangian method

Our algorithms are based on augmented Lagrangian function of (1). In the literature, there are several different augmented Lagrangian functions (see [1] for example), and we use the classic one. Let

$$\psi_\beta(u, v) = \begin{cases} uv + \frac{\beta}{2}u^2, & \text{if } \beta u + v \geq 0, \\ -\frac{v^2}{2\beta}, & \text{if } \beta u + v < 0, \end{cases}$$

and

$$\Psi_\beta(x, z) = \sum_{j=1}^m \psi_\beta(f_j(x), z_j).$$

Then the classic augmented Lagrangian function of (1) is

$$\mathcal{L}_\beta(x, y, z) = g(x) + h(x) + y^\top(Ax - b) + \frac{\beta}{2}\|Ax - b\|^2 + \Psi_\beta(x, z), \quad (2)$$

where y and z are Lagrangian multipliers, and $\beta > 0$ is the penalty parameter.

The augmented Lagrangian method (ALM) for (1), at each iteration, renews x -variable by minimizing \mathcal{L}_β with respect to x while y and z are fixed and then perform an augmented dual gradient ascent update to the multipliers y and z , namely,

$$x^{k+1} \in \arg \min_x \mathcal{L}_\beta(x, y^k, z^k), \quad (3a)$$

$$y^{k+1} = y^k + \rho_y \nabla_y \mathcal{L}_\beta(x^{k+1}, y^k, z^k), \quad (3b)$$

$$z^{k+1} = z^k + \rho_z \nabla_z \mathcal{L}_\beta(x^{k+1}, y^k, z^k). \quad (3c)$$

In general, it is difficult to solve the x -subproblem exactly or to a high accuracy. In one recent work [19], we show that if (3a) is solved to a certain error tolerance, a global sublinear convergence of the inexact ALM can be established. In this work, we propose to perform one single proximal gradient update to (3a), and a sublinear convergence can still be shown.

1.2 Related work

ALM has been popularly used to solve constrained optimization problems; see books [2,3]. However, most works on first-order methods in the ALM framework consider affinely constrained problems, and only a few study the methods for generally constrained problems in the form of (1). We review these works below.

For smooth affinely constrained convex programs, [10] analyzes the iteration complexity of an inexact ALM, where each primal subproblem is approximately solved by Nesterov’s optimal first-order method [14]. It shows that to reach an ε -optimal solution (see Definition 1.1 below), $O(\varepsilon^{-\frac{7}{4}})$ gradient evaluations are sufficient. In addition, it shows that $O(\varepsilon^{-1}|\log \varepsilon|)$ gradient evaluations can guarantee an ε -optimal solution by an inexact proximal ALM. Although the number of gradient evaluations is not explicitly given, [11, 12] also consider inexact ALM. They specify the accuracy that each primal subproblem need be solved to and estimate the outer iteration complexity of the inexact ALM. Within the ALM framework, [17] perform a single proximal gradient update to primal variable at each iteration and establish $O(\varepsilon^{-1})$ complexity result to have an ε -optimal solution for affinely constrained composite convex programs. This linearized ALM also appears as a special case of the methods in [5–7, 9, 20], which perform Gauss-Seidel or randomized block coordinate update to the primal variable in the ALM framework.

Towards finding solutions of general saddle-point problems, [13] gives a subgradient method. If both primal and dual constraint sets are compact, the method has $O(1/\sqrt{k})$ convergence rate in terms of primal-dual gap, where k is the number of iterations. It also discusses how to apply the subgradient method to convex optimization problems with nonlinear inequality constraint. On smooth constrained convex problems, [21] proposes a primal-dual type first-order method (see (69) in section 6.2). Assuming compactness of the constraint set, it establishes $O(\varepsilon^{-1})$ iteration complexity result to produce an ε -optimal solution. Recently, [19] studies an inexact ALM for (1) and proposes to use Nesterov’s optimal first-order method to approximately solve each x -subproblem. When the constraint set is bounded, it shows that nearly $O(\varepsilon^{-\frac{3}{2}})$ gradient evaluations suffice to obtain an ε -optimal solution, and for the smooth case, the result can be improved to $O(\varepsilon^{-1}|\log \varepsilon|)$. Compared to these works, our iteration complexity results will be better under weaker assumptions.

1.3 Contributions

This paper mainly makes the following contributions.

- We propose a first-order method, named LALM, for solving composite convex problems with both affine equality and smooth nonlinear inequality constraints. The method is based on proximal linearization of the classic augmented Lagrangian function. Under mild assumptions, we show global iterate sequence convergence of LALM to a primal-dual optimal solution.

- Also, we analyze the iteration complexity of the proposed method. We show that to reach an ε -optimal solution, $O(\varepsilon^{-1})$ gradient evaluations are sufficient. In addition, we establish its local linear convergence by assuming the existence of a non-degenerate primal-dual solution and positive definiteness of Hessian of the augmented Lagrangian function near the non-degenerate primal-dual solution.
- Furthermore, as the problem has the so-called coordinate friendly structure, we propose a block update version of LALM. At each iteration, the method renews a single block coordinate while keeping all the other coordinates unchanged and then immediately performs an update to dual variables. We show that in expectation, an ε -optimal solution can be obtained by $O(\varepsilon^{-1})$ gradient evaluations.
- We implement LALM and its block update version and apply them to the basis pursuit denoising problem and a quadratically constrained quadratic program. On both problems, we notice better performance of the block-LALM in terms of iteration number. In addition, when the iterate is far away from optimality, sublinear convergence is observed, and while the iterate approaches to optimality, both methods converge linearly.

1.4 Notation and organization

We focus on finite-dimensional Euclidean space, but our analysis can be directly extended to a Hilbert space. We use $[m]$ as the set $\{1, 2, \dots, m\}$, and $[a]_+ = \max(0, a)$ denotes the positive part of a real number a . We use I as the identity matrix. Given a symmetric positive definite (SPD) matrix P , we define $\|x\|_P = \sqrt{x^\top P x}$, and if $P = I$, we simply write it as $\|x\|$. Also, given a nonnegative vector $\ell = [\ell_1, \dots, \ell_n] \in \mathbb{R}^n$, we define $\|x\|_\ell^2 = \sum_{i=1}^n \ell_i \|x_i\|^2$ if x is partitioned into n blocks (x_1, \dots, x_n) . For any convex function $f(x)$, we use $\tilde{\nabla} f(x)$ as its subgradient and $\partial f(x)$ the subdifferential of f at x , i.e., the set of all subgradients at x . When f is differentiable, $\tilde{\nabla} f(x)$ coincides with the gradient of f , and we simply write it to $\nabla f(x)$. The indicator function of a set \mathcal{X} is defined as $\iota_{\mathcal{X}}(x) = 0$ if $x \in \mathcal{X}$ and $+\infty$ otherwise. $B_\gamma(x)$ represents a ball with radius γ and center x . \mathbb{E}_{i_k} denotes the expectation about i_k conditioned on all previous history.

For ease of notation, we use w as the triple (x, y, z) and denote the smooth part of \mathcal{L}_β as

$$F_\beta(w) = \mathcal{L}_\beta(w) - h(x).$$

In addition, we define

$$\Phi(\bar{x}; x, y, z) = f_0(\bar{x}) - f_0(x) + y^\top (A\bar{x} - b) + \sum_{j=1}^m z_j f_j(\bar{x}). \quad (4)$$

Definition 1.1 (ε -optimal solution) *Let f_0^* be the optimal value of (1). We call \bar{x} an ε -optimal solution to (1) if*

$$|f_0(\bar{x}) - f_0^*| \leq \varepsilon, \quad \|A\bar{x} - b\| + \sum_{j=1}^m [f_j(\bar{x})]_+ \leq \varepsilon.$$

Organization. The rest of the paper is organized as follows. Section 2 gives several technical results that will be used to prove our main theorems. We propose a linearized ALM for (1) in section 3 and a block linearized ALM in section 4. Convergence results are also given. In section 5, we discuss a few applications and how the proposed methods can be applied. Numerical results are given in section 6, and finally section 7 concludes the paper.

2 Technical assumptions and preliminary results

A point $w = (x, y, z)$ satisfies the Karush-Kuhn-Tucker (KKT) conditions for (1) if

$$0 \in \nabla g(x) + \partial h(x) + A^\top y + \sum_{j=1}^m z_j \nabla f_j(x), \quad (5a)$$

$$Ax = b, \quad (5b)$$

$$z_j \geq 0, f_j(x) \leq 0, z_j f_j(x) = 0, \forall j \in [m]. \quad (5c)$$

If w satisfies the above conditions, we call it a KKT point. For convex programs, the conditions in (5) are sufficient for x to be an optimal solution of (1). If a certain qualification condition (e.g., the Slater condition) holds, they are also necessary.

2.1 Technical assumptions

Throughout the paper, we assume the existence of a KKT point.

Assumption 1 *There exists a point $w^* = (x^*, y^*, z^*)$ satisfying the KKT conditions in (5).*

Under the above assumption, it follows from the convexity of f_0 that

$$\Phi(x; w^*) \geq 0, \forall x, \quad (6)$$

where Φ is defined in (4).

In addition, we make the following assumption, which holds if $\text{dom}(h)$ is bounded.

Assumption 2 *There are constants L_g, L_1, \dots, L_m and B_1, \dots, B_m such that*

$$\|\nabla g(\hat{x}) - \nabla g(\tilde{x})\| \leq L_g \|\hat{x} - \tilde{x}\|, \forall \hat{x}, \tilde{x} \in \text{dom}(h), \quad (7)$$

$$\|\nabla f_j(\hat{x}) - \nabla f_j(\tilde{x})\| \leq L_j \|\hat{x} - \tilde{x}\|, \forall \hat{x}, \tilde{x} \in \text{dom}(h), \forall j \in [m], \quad (8)$$

$$\|\nabla f_j(x)\| \leq B_j, \forall x \in \text{dom}(h), \forall j \in [m]. \quad (9)$$

From the mid-point theorem, the boundedness of ∇f_j implies the Lipschitz continuity of f_j , i.e.,

$$|f_j(\hat{x}) - f_j(\tilde{x})| \leq B_j \|\hat{x} - \tilde{x}\|, \forall \hat{x}, \tilde{x} \in \text{dom}(h), \forall j \in [m]. \quad (10)$$

2.2 Preparatory lemmas

In this subsection, we give several lemmas that will be used multiple times in our convergence analysis. First we show the Lipschitz continuity of $\nabla_x \Psi(w)$ with respect to x .

Lemma 2.1 *Under Assumption 2, we have*

$$\|\nabla_x \Psi(\hat{x}, z) - \nabla_x \Psi(x, z)\| \leq L_\Psi(x, z) \|\hat{x} - x\|, \quad \forall \hat{x}, x, z, \quad (11)$$

where

$$L_\Psi(x, z) = \sum_{j=1}^m \left(\beta B_j^2 + L_j [\beta f_j(x) + z_j]_+ \right) \quad (12)$$

Proof. First we notice that $\frac{\partial}{\partial u} \psi_\beta(u, v) = [\beta u + v]_+$, and thus for any v ,

$$\left| \frac{\partial}{\partial u} \psi_\beta(\hat{u}, v) - \frac{\partial}{\partial u} \psi_\beta(\tilde{u}, v) \right| \leq \beta |\hat{u} - \tilde{u}|, \quad \forall \hat{u}, \tilde{u}.$$

Let $h_j(x, z_j) = \psi_\beta(f_j(x), z_j)$, $j = 1, \dots, m$. Then

$$\begin{aligned} & \|\nabla_x h_j(\hat{x}, z_j) - \nabla_x h_j(x, z_j)\| \\ &= \left\| \frac{\partial}{\partial u} \psi_\beta(f_j(\hat{x}), z_j) \nabla f_j(\hat{x}) - \frac{\partial}{\partial u} \psi_\beta(f_j(x), z_j) \nabla f_j(x) \right\| \\ &\leq \left\| \frac{\partial}{\partial u} \psi_\beta(f_j(\hat{x}), z_j) \nabla f_j(\hat{x}) - \frac{\partial}{\partial u} \psi_\beta(f_j(x), z_j) \nabla f_j(\hat{x}) \right\| \\ &\quad + \left\| \frac{\partial}{\partial u} \psi_\beta(f_j(x), z_j) \nabla f_j(\hat{x}) - \frac{\partial}{\partial u} \psi_\beta(f_j(x), z_j) \nabla f_j(x) \right\| \\ &\leq \beta |f_j(\hat{x}) - f_j(x)| \cdot \|\nabla f_j(\hat{x})\| + \left| \frac{\partial}{\partial u} \psi_\beta(f_j(x), z_j) \right| \cdot \|\nabla f_j(\hat{x}) - \nabla f_j(x)\| \\ &\leq \beta B_j^2 \|\hat{x} - x\| + L_j [\beta f_j(x) + z_j]_+ \cdot \|\hat{x} - x\|. \end{aligned} \quad (13)$$

Hence,

$$\|\nabla_x \Psi_\beta(\hat{x}, z) - \nabla_x \Psi_\beta(x, z)\| \leq \sum_{j=1}^m \|\nabla_x h_j(\hat{x}, z_j) - \nabla_x h_j(x, z_j)\| \leq L_\Psi(x, z) \|\hat{x} - x\|,$$

which completes the proof. \square

Remark 2.1 *Note that the Lipschitz constant $L_\Psi(x, z)$ in (11) depends on the point (x, z) and is not a universal constant. We will set its value at the iterate of the algorithm. Together with the next lemma, the inequality in (11) implies that a sufficient progress can be obtained after each x -update.*

Lemma 2.2 *For a continuously differentiable function $\phi(u)$ and a given v , if $\|\nabla \phi(u) - \nabla \phi(v)\| \leq L_\phi(v) \|u - v\|$, $\forall u$, then*

$$\phi(u) \leq \phi(v) + \langle \nabla \phi(v), u - v \rangle + \frac{L_\phi(v)}{2} \|u - v\|^2.$$

The following result is easy to show (c.f., [4, Prop. 2.3]). It will be used for establishing iterate convergence of the proposed algorithm.

Lemma 2.3 *Let $\{P^k\}$ be a sequence of SPD matrices, and there are SPD matrices \underline{P} and \overline{P} such that $\underline{P} \succeq P^k \succeq \overline{P}$. Let \mathcal{W} be a nonempty set. If the sequence $\{w^k\}$ satisfies*

$$\|w^{k+1} - w\|_{P^{k+1}}^2 \leq \|w^k - w\|_{P^k}^2, \forall w \in \mathcal{W},$$

and $\{w^k\}$ has a cluster point \bar{w} in \mathcal{W} , then w^k converges to \bar{w} .

The result below will be used to establish convergence rate of our algorithms. It is similar to a deterministic result in [19] and can be shown in the same way. We omit its proof.

Lemma 2.4 *Assume (x^*, y^*, z^*) is a KKT point of (1). Let \bar{x} be a stochastic point such that for any y and any $z \geq 0$,*

$$\mathbb{E}[\Phi(\bar{x}; x^*, y, z)] \leq \alpha + c_1 \|y\|^2 + c_2 \|z\|^2, \quad (14)$$

where α and c_1, c_2 are nonnegative constants independent of y and z . Then

$$-\left(\alpha + 4c_1 \|y^*\|^2 + 4c_2 \sum_{j=1}^m (z_j^*)^2 \right) \leq \mathbb{E}[f_0(\bar{x}) - f_0(x^*)] \leq \alpha, \quad (15)$$

$$\mathbb{E}\|A\bar{x} - b\| + \sum_{j=1}^m \mathbb{E}[f_j(\bar{x})]_+ \leq \alpha + c_1 (1 + \|y^*\|)^2 + c_2 \sum_{j=1}^m (1 + z_j^*)^2. \quad (16)$$

3 Linearized augmented Lagrangian method

In this section, we propose a linearized augmented Lagrangian method (LALM). Different from the step in (3a), it updates x -variable by a single proximal gradient descent of the augmented Lagrangian function. The method is summarized in Algorithm 1, where $\delta \geq 0$ is a constant and

$$L_F^k = L_g + \beta \|A\|^2 + L_\Psi(x^k, z^k)$$

with L_Ψ defined in (11).

Note that the setting of η^k is for simplicity of our analysis. Practically, one can choose it by starting from η^{k-1} and then backtracking such that

$$F_\beta(x^{k+1}, y^k, z^k) \leq F_\beta(w^k) + \langle \nabla_x F_\beta(w^k), x^{k+1} - x^k \rangle + \frac{\eta^k}{2} \|x^{k+1} - x^k\|^2, \quad (17)$$

and all our convergence results can still be shown. When $\eta^k \geq L_F^k$, the above inequality always holds from Lemma 2.2.

Algorithm 1: Linearized augmented Lagrangian method (LALM) for (1)

1 Initialization: choose x^0, y^0, z^0 and $\beta, \rho_y, \rho_z, \delta \geq 0$; set $\eta^{-1} = 0$

2 for $k = 0, 1, \dots$ **do**

3 Let $\eta^k = \max(\eta^{k-1}, L_F^k + \delta)$

4 Perform the updates

$$x^{k+1} = \arg \min_x h(x) + \langle \nabla_x F_\beta(w^k), x \rangle + \frac{\eta^k}{2} \|x - x^k\|^2, \quad (18a)$$

$$y^{k+1} = y^k + \rho_y (Ax^{k+1} - b), \quad (18b)$$

$$z_j^{k+1} = z_j^k + \rho_z \cdot \max\left(-\frac{z_j^k}{\beta}, f_j(x^{k+1})\right), j = 1, \dots, m. \quad (18c)$$

3.1 Global convergence analysis

To show the convergence results of Algorithm 1, we need the following two lemmas, which can be found in [19].

Lemma 3.1 *Let y and z be updated by (18b) and (18c) respectively. Then for any k , it holds*

$$\frac{1}{2\rho_y} [\|y^{k+1} - y\|^2 - \|y^k - y\|^2 + \|y^{k+1} - y^k\|^2] - \langle y^{k+1} - y, r^{k+1} \rangle = 0, \quad (19)$$

$$\frac{1}{2\rho_z} [\|z^{k+1} - z\|^2 - \|z^k - z\|^2 + \|z^{k+1} - z^k\|^2] - \sum_{j=1}^m (z_j^{k+1} - z_j^k) \cdot \max\left(-\frac{z_j^k}{\beta}, f_j(x^{k+1})\right) = 0, \quad (20)$$

where $r^k = Ax^k - b$.

Lemma 3.2 *For any $z \geq 0$, we have*

$$\frac{\beta - 2\rho_z}{2\rho_z^2} \|z^{k+1} - z^k\|^2 \leq \Psi_\beta(x^{k+1}, z^k) - \sum_{j=1}^m z_j f_j(x^{k+1}) - \sum_{j=1}^m (z_j^{k+1} - z_j^k) \cdot \max\left(-\frac{z_j^k}{\beta}, f_j(x^{k+1})\right). \quad (21)$$

Using the above two lemmas, we establish a fundamental result on Algorithm 1.

Theorem 3.1 (One-iteration progress of LALM) *Let $\{w^k\}$ be the sequence generated from Algorithm 1. Then for any x such that $Ax = b$ and $f_j(x) \leq 0, \forall j \in [m]$, any y , and any $z \geq 0$, it*

holds that

$$\begin{aligned}
& \Phi(x^{k+1}; w) + \frac{\eta^k}{2} \|x^{k+1} - x\|^2 + \frac{1}{2\rho_y} \|y^{k+1} - y\|^2 + \frac{1}{2\rho_z} \|z^{k+1} - z\|^2 \\
& + \frac{\beta - \rho_y}{2} \|r^{k+1}\|^2 + \frac{\beta - \rho_z}{2\rho_z^2} \|z^{k+1} - z^k\|^2 + \frac{1}{2} \|x^{k+1} - x^k\|_{(\eta^k - L_g - L_\Psi^k)I - \beta A^\top A}^2 \\
& \leq \frac{\eta^k}{2} \|x^k - x\|^2 - \frac{\beta}{2} \|r^k\|^2 + \frac{1}{2\rho_y} \|y^k - y\|^2 + \frac{1}{2\rho_z} \|z^k - z\|^2,
\end{aligned} \tag{22}$$

where Φ is defined in (4), and $L_\Psi^k = L_\Psi(x^k, z^k)$ with L_Ψ defined in (12).

Proof. From the update in (18a), it follows that

$$0 \in \partial h(x^{k+1}) + \nabla g(x^k) + A^\top y^k + \beta A^\top r^k + \nabla_x \Psi(x^k, z^k) + \eta^k (x^{k+1} - x^k). \tag{23}$$

By the convexity of h , we have

$$\langle x^{k+1} - x, \tilde{\nabla} h(x^{k+1}) \rangle \geq h(x^{k+1}) - h(x), \tag{24}$$

From the convexity of g and $\Psi(\cdot, z)$, and also Lemmas 2.1 and 2.2, we have

$$\begin{aligned}
& \langle x^{k+1} - x, \nabla g(x^k) + \nabla_x \Psi(x^k, z^k) \rangle \\
& = \langle x^{k+1} - x^k, \nabla g(x^k) + \nabla_x \Psi(x^k, z^k) \rangle + \langle x^k - x, \nabla g(x^k) + \nabla_x \Psi(x^k, z^k) \rangle \\
& \geq g(x^{k+1}) + \Psi(x^{k+1}, z^k) - g(x^k) - \Psi(x^k, z^k) - \frac{L_g + L_\Psi^k}{2} \|x^{k+1} - x^k\|^2 \\
& \quad + g(x^k) + \Psi(x^k, z^k) - g(x) - \Psi(x, z^k) \\
& = g(x^{k+1}) + \Psi(x^{k+1}, z^k) - g(x) - \Psi(x, z^k) - \frac{L_g + L_\Psi^k}{2} \|x^{k+1} - x^k\|^2.
\end{aligned} \tag{25}$$

For x such that $Ax = b$, it holds that

$$\begin{aligned}
& \langle x^{k+1} - x, A^\top y^k + \beta A^\top r^k \rangle \\
& = \langle r^{k+1}, y^{k+1} - \rho_y r^{k+1} + \beta r^k \rangle \\
& = y^\top r^{k+1} + \langle y^{k+1} - y, r^{k+1} \rangle + (\beta - \rho_y) \|r^{k+1}\|^2 + \beta \langle r^{k+1}, r^k - r^{k+1} \rangle \\
& = y^\top r^{k+1} + \langle y^{k+1} - y, r^{k+1} \rangle + (\beta - \rho_y) \|r^{k+1}\|^2 - \frac{\beta}{2} [\|r^{k+1}\|^2 - \|r^k\|^2 + \|r^{k+1} - r^k\|^2].
\end{aligned} \tag{26}$$

Adding (24), (25), (26), and the following equation

$$\langle x^{k+1} - x, \eta^k (x^{k+1} - x^k) \rangle = \frac{\eta^k}{2} [\|x^{k+1} - x\|^2 - \|x^k - x\|^2 + \|x^{k+1} - x^k\|^2],$$

we have from (23) that

$$\begin{aligned}
& f_0(x^{k+1}) - f_0(x) + \Psi(x^{k+1}, z^k) - \Psi(x, z^k) + y^\top r^{k+1} + \langle y^{k+1} - y, r^{k+1} \rangle \\
& - \frac{\beta}{2} [\|r^{k+1}\|^2 + \|r^{k+1} - r^k\|^2] + (\beta - \rho_y) \|r^{k+1}\|^2 + \frac{\eta^k}{2} \|x^{k+1} - x\|^2 + \frac{\eta^k - L_g - L_\Psi^k}{2} \|x^{k+1} - x^k\|^2 \\
& \leq \frac{\eta^k}{2} \|x^k - x\|^2 - \frac{\beta}{2} \|r^k\|^2.
\end{aligned} \tag{27}$$

The desired result is obtained by noting $\Psi(x, z^k) \leq 0$, adding (19), (20), and (21) to the above inequality, and rearranging terms. \square

The next lemma shows the upper boundedness of η^k .

Lemma 3.3 *Let $\{w^k\}$ be the sequence generated from Algorithm 1 with $z_j^0 \geq 0, \forall j \in [m]$. If $\rho_y, \rho_z \in (0, \beta]$, then $\eta^k \leq \bar{\eta}, \forall k \geq 0$, where $\bar{\eta}$ is a constant satisfying*

$$\begin{aligned} \bar{\eta} \geq & \delta + L_g + \beta \|A\|^2 + \beta \sum_{j=1}^m B_j^2 + \beta \sum_{j=1}^m B_j L_j \left(\|x^0 - x^*\| + \frac{\|y^0 - y^*\|}{\sqrt{\rho_y \eta^0}} + \frac{\|z^0 - z^*\|}{\sqrt{\rho_z \eta^0}} \right) \\ & + \sqrt{\sum_{j=1}^m L_j^2} \left(\sqrt{\rho_z \bar{\eta}} \|x^0 - x^*\| + \sqrt{\frac{\rho_z \bar{\eta}}{\rho_y \eta^0}} \|y^0 - y^*\| + \|z^*\| + \max(1, \sqrt{\frac{\bar{\eta}}{\eta^0}}) \|z^0 - z^*\| \right). \end{aligned} \quad (28)$$

Proof. Since $z_j^0 \geq 0, \forall j \in [m]$, we have $z_j^k \geq 0, \forall j \in [m]$ from the update of z and the condition $\rho_z \in (0, \beta]$. Note $f_j(x^k) \leq f_j(x^*) + B_j \|x^k - x^*\| \leq B_j \|x^k - x^*\|$. It follows from the increasing monotonicity of $[a]_+$ that

$$\sum_{j=1}^m L_j [\beta f_j(x^k) + z_j^k]_+ \leq \sum_{j=1}^m (\beta B_j L_j \|x^k - x^*\| + L_j z_j^k),$$

and thus

$$\begin{aligned} L_{\Psi}^k & \leq \beta \sum_{j=1}^m B_j^2 + \sum_{j=1}^m (\beta B_j L_j \|x^k - x^*\| + L_j z_j^k) \\ & \leq \beta \sum_{j=1}^m B_j^2 + \beta \sum_{j=1}^m B_j L_j \|x^k - x^*\| + \sqrt{\sum_{j=1}^m L_j^2} (\|z^*\| + \|z^k - z^*\|). \end{aligned} \quad (29)$$

We next show the desired result by induction. First, the result for $k = 0$ directly follows from (28) and (29). Assume $\eta^k \leq \bar{\eta}, \forall k \leq K - 1$. Then letting $w = w^*$ in (22), we have from (6) and by dropping nonnegative terms on the left hand side that

$$\begin{aligned} & \frac{\eta^k}{2} \|x^{k+1} - x^*\|^2 + \frac{1}{2\rho_y} \|y^{k+1} - y^*\|^2 + \frac{1}{2\rho_z} \|z^{k+1} - z^*\|^2 \\ & \leq \frac{\eta^k}{2} \|x^k - x^*\|^2 + \frac{1}{2\rho_y} \|y^k - y^*\|^2 + \frac{1}{2\rho_z} \|z^k - z^*\|^2. \end{aligned} \quad (30)$$

Since $\eta^{k+1} \geq \eta^k$, dividing by η^k on both sides of the above inequality yields

$$\begin{aligned} & \frac{1}{2} \|x^{k+1} - x^*\|^2 + \frac{1}{2\rho_y \eta^{k+1}} \|y^{k+1} - y^*\|^2 + \frac{1}{2\rho_z \eta^{k+1}} \|z^{k+1} - z^*\|^2 \\ & \leq \frac{1}{2} \|x^k - x^*\|^2 + \frac{1}{2\rho_y \eta^k} \|y^k - y^*\|^2 + \frac{1}{2\rho_z \eta^k} \|z^k - z^*\|^2. \end{aligned} \quad (31)$$

Repeatedly using (31) and also from (30), we have

$$\begin{aligned} & \frac{1}{2}\|x^K - x^*\|^2 + \frac{1}{2\rho_y\eta^{K-1}}\|y^K - y^*\|^2 + \frac{1}{2\rho_z\eta^{K-1}}\|z^K - z^*\|^2 \\ & \leq \frac{1}{2}\|x^0 - x^*\|^2 + \frac{1}{2\rho_y\eta^0}\|y^0 - y^*\|^2 + \frac{1}{2\rho_z\eta^0}\|z^0 - z^*\|^2. \end{aligned}$$

The above inequality together with $\eta^{K-1} \leq \bar{\eta}$ implies

$$\|x^K - x^*\| \leq \|x^0 - x^*\| + \frac{\|y^0 - y^*\|}{\sqrt{\rho_y\eta^0}} + \frac{\|z^0 - z^*\|}{\sqrt{\rho_z\eta^0}}$$

and

$$\|z^K - z^*\| \leq \sqrt{\rho_z\bar{\eta}}\|x^0 - x^*\| + \sqrt{\frac{\rho_z\bar{\eta}}{\rho_y\eta^0}}\|y^0 - y^*\| + \sqrt{\frac{\bar{\eta}}{\eta^0}}\|z^0 - z^*\|.$$

Hence, $\bar{\eta} \geq L_F^K + \delta$ from (28), (29), and the above two inequalities. This completes the proof. \square

We are now ready to show our main convergence and rate results.

Theorem 3.2 (Iterate convergence of LALM) *Under Assumptions 1 and 2, let $\{w^k\}$ be the sequence generated from Algorithm 1 with any x^0, y^0 , and $z_j^0 \geq 0, \forall j \in [m]$. If $\rho_y, \rho_z \in (0, \beta)$ and $\delta > 0$, then w^k converges to a KKT point $\bar{w} = (\bar{x}, \bar{y}, \bar{z})$ of (1).*

Proof. Letting $w = w^*$ in (22) and dividing by η^k , we have from (6) that

$$\begin{aligned} & \frac{1}{2}\|x^{k+1} - x^*\|^2 + \frac{1}{2\rho_y\eta^k}\|y^{k+1} - y^*\|^2 + \frac{1}{2\rho_z\eta^k}\|z^{k+1} - z^*\|^2 \\ & + \frac{\beta - \rho_y}{2\eta^k}\|r^{k+1}\|^2 + \frac{1}{2\rho_z^2\eta^k}(\beta - \rho_z)\|z^{k+1} - z^k\|^2 + \frac{1}{2\eta^k}\|x^{k+1} - x^k\|_{(\eta^k - L_g - L_\Psi)I - \beta A^\top A}^2 \\ & \leq \frac{1}{2}\|x^k - x^*\|^2 + \frac{1}{2\rho_y\eta^k}\|y^k - y^*\|^2 + \frac{1}{2\rho_z\eta^k}\|z^k - z^*\|^2. \end{aligned} \quad (32)$$

Summing up (32) over k and noting $\eta^{k+1} \geq \eta^k$, we have from the condition $\delta > 0, \rho_y, \rho_z \in (0, \beta)$ and Lemma 3.3 that

$$\lim_{k \rightarrow \infty} x^{k+1} - x^k = 0, \quad \lim_{k \rightarrow \infty} y^{k+1} - y^k = \lim_{k \rightarrow \infty} \rho_y r^{k+1} = 0, \quad \lim_{k \rightarrow \infty} z^{k+1} - z^k = 0. \quad (33)$$

In addition, it follows from (32) that $\{w^k\}$ is bounded and must have a cluster point \bar{w} . Hence, there is a subsequence $\{w^k\}_{k \in \mathcal{K}}$ convergent to \bar{w} . Since $\{\eta^k\}$ is increasing and bounded, it must converge to a number η^∞ .

Below we show that \bar{w} is a KKT point. First, we have $A\bar{x} - b = 0$ from (33), i.e., \bar{x} satisfies (5b).

Secondly, from the update of z and $\rho_z < \beta$, it follows $z_j^k \geq 0, \forall j \in [m], \forall k$, and thus $\bar{z}_j \geq 0, \forall j \in [m]$. If $f_j(x^{k+1}) > 0$, then $f_j(x^{k+1}) = \frac{1}{\rho_z}(z_j^{k+1} - z_j^k) \rightarrow 0$ that indicates $[f_j(x^{k+1})]_+ \rightarrow 0$. Hence,

$f_j(\bar{x}) \leq 0, \forall j \in [m]$ follows from the continuity of f_j 's. For any $j \in [m]$, if $\bar{z}_j > 0$, then $z_j^k > \frac{\bar{z}_j}{2}$, as $k \in \mathcal{K}$ is sufficiently large. It follows from $\max(-\frac{z_j^k}{\beta}, f_j(x^{k+1})) \rightarrow 0$ that $f_j(x^{k+1}) \rightarrow 0$ as $\mathcal{K} \ni k \rightarrow \infty$. Hence, $f_j(\bar{x}) = 0$. Therefore, (\bar{x}, \bar{z}) satisfies (5c).

Thirdly, from the optimality of x^{k+1} , it holds that

$$\langle \nabla_x F(w^k), x^{k+1} \rangle + h(x^{k+1}) + \frac{\eta^k}{2} \|x^{k+1} - x^k\|^2 \leq \langle \nabla_x F(w^k), x \rangle + h(x) + \frac{\eta^k}{2} \|x - x^k\|^2, \forall x.$$

Taking limit infimum over $k \in \mathcal{K}$ on both sides of the above equation, we have from the lower semicontinuity of h and continuity of g and f_j 's that

$$\langle \nabla_x F(\bar{w}), \bar{x} \rangle + h(\bar{x}) \leq \langle \nabla_x F(\bar{w}), x \rangle + h(x) + \frac{\eta^\infty}{2} \|x - \bar{x}\|^2, \forall x,$$

namely,

$$\bar{x} = \arg \min_x \langle \nabla_x F(\bar{w}), x \rangle + h(x) + \frac{\eta^\infty}{2} \|x - \bar{x}\|^2.$$

Therefore, \bar{w} satisfies (5a) from the optimality condition of the above minimization problem, and thus \bar{w} is a KKT point of (1).

Hence, (31) holds with w^* replaced by \bar{w} , and thus w^k converges to \bar{w} from Lemma 2.3. \square

Theorem 3.3 (Sublinear convergence rate of LALM) *Under Assumptions 1 and 2, let $\{w^k\}$ be the sequence generated from Algorithm 1 with any $x^0, y^0 = 0$ and $z^0 = 0$. If $\rho_y, \rho_z \in (0, \beta]$, then*

$$|f_0(\bar{x}^{k+1}) - f_0(x^*)| \leq \frac{\eta^\infty}{\eta^0(k+1)} \left(\frac{\eta^0}{2} \|x^0 - x^*\|^2 + \frac{2\|y^*\|^2}{\rho_y} + \sum_{j=1}^m \frac{2(z_j^*)^2}{\rho_z} \right), \quad (34a)$$

$$\|A\bar{x}^{k+1} - b\| + \sum_{j=1}^m [f_j(\bar{x}^{k+1})]_+ \leq \frac{\eta^\infty}{\eta^0(k+1)} \left(\frac{\eta^0}{2} \|x^0 - x^*\|^2 + \frac{(1 + \|y^*\|)^2}{2\rho_y} + \sum_{j=1}^m \frac{(1 + z_j^*)^2}{2\rho_z} \right), \quad (34b)$$

where $\eta^\infty \leq \bar{\eta}$ is the limit of $\{\eta^k\}$ and $\bar{x}^{k+1} = \sum_{t=0}^k \frac{x^{t+1}}{\sum_{t=0}^k \frac{1}{\eta^t}}$.

Proof. Letting $x = x^*$ in (22), dividing by η^k , and summing it up give

$$\left(\sum_{t=0}^k \frac{1}{\eta^t} \right) \Phi(\bar{x}^{k+1}; x^*, y, z) \leq \sum_{t=0}^k \frac{1}{\eta^t} \Phi(x^{t+1}; x^*, y, z) \leq \frac{1}{\eta^0} \left[\frac{\eta^0}{2} \|x^0 - x^*\|^2 + \frac{1}{2\rho_y} \|y\|^2 + \frac{1}{2\rho_z} \|z\|^2 \right].$$

Since $\sum_{t=0}^k \frac{1}{\eta^t} \geq \frac{k+1}{\eta^\infty}$, the desired results are obtained from Lemma 2.4 with $\alpha = \frac{\eta^\infty}{2(k+1)} \|x^0 - x^*\|^2$, $c_1 = \frac{\eta^\infty}{2\rho_y \eta^0(k+1)}$, and $c_2 = \frac{\eta^\infty}{2\rho_z \eta^0(k+1)}$. \square

Theorem 3.3 implies that to reach an ε -optimal solution of (1), it is sufficient to evaluate the gradients of g and f_j , $j \in [m]$ and proximal mapping of h for K times, where

$$K = \left\lceil \frac{\eta^\infty}{\varepsilon\eta^0} \left(\frac{\eta^0}{2} \|x^0 - x^*\|^2 + \frac{[\max(1 + \|y^*\|, 2\|y^*\|)]^2}{2\rho_y} + \sum_{j=1}^m \frac{[\max(1 + z_j^*, 2z_j^*)]^2}{2\rho_z} \right) \right\rceil.$$

3.2 Local linear convergence of LALM for constrained smooth problems

In this subsection, we assume that $h(x) = \iota_{\mathcal{X}}(x)$ for a closed convex set \mathcal{X} and g, f_1, \dots, f_m are twice continuously differentiable. We show local linear convergence of Algorithm 1 under the following assumption.

Assumption 3 *There is a KKT point w^* and a subset $J \subset [m]$ such that $x^* \in \text{int}(\mathcal{X})$, and*

1. $f_j(x^*) = 0, z_j^* > 0, \forall j \in J$ and $f_j(x^*) < 0, z_j^* = 0, \forall j \notin J$;
2. $x^\top \left(\nabla^2 g(x^*) + \sum_{j \in J} z_j^* \nabla^2 f_j(x^*) \right) x > 0$ for any nonzero vector $x \in \text{Null}(D^\top)$, where

$$D = [A^\top, \nabla f_1(x^*), \dots, \nabla f_m(x^*)]$$

is column full-rank.

When item 1 holds in the above assumption, we have

$$\nabla_x^2 \mathcal{L}_\beta(w^*) = \nabla^2 g(x^*) + \beta A^\top A + \sum_{j \in J} \left(z_j^* \nabla^2 f_j(x^*) + \beta \nabla f_j(x^*) [\nabla f_j(x^*)]^\top \right),$$

and thus if in addition item 2 holds, then $\nabla_x^2 \mathcal{L}_\beta(w^*)$ is positive definite. We denote $\mu > 0$ as its smallest eigenvalue.

From the continuity of $\nabla_x^2 \mathcal{L}_\beta$, we have the following result.

Proposition 3.1 *There is $\gamma > 0$ such that if $\max(\|x - x^*\|, \|y - y^*\|, \|z - z^*\|) \leq \gamma$, then*

$$x \in \text{int}(\mathcal{X}); \nabla_x^2 \mathcal{L}_\beta(w) \succeq \frac{\mu}{2} I; \beta f_j(x) + z_j > 0, \forall j \in J; \beta f_j(x) + z_j < 0, \forall j \notin J. \quad (35)$$

Hence, for any $x \in \mathcal{B}_\gamma(x^)$,*

$$f_0(x) - f_0(x^*) + \langle y^*, Ax - b \rangle + \frac{\beta}{2} \|Ax - b\|^2 + \Psi(x, z^*) - \Psi(x^*, z^*) \geq \frac{\mu}{4} \|x - x^*\|^2. \quad (36)$$

The next lemma can be easily verified from the definition of Ψ . We omit its proof.

Lemma 3.4 *If $x^{k+1} \in \mathcal{B}_\gamma(x^*)$ and $z^k \in \mathcal{B}_\gamma(z^*)$, then*

$$\sum_{j \in J} (z_j^k - z_j^*) f_j(x^{k+1}) = \Psi(x^{k+1}, z^k) - \Psi(x^{k+1}, z^*) - \Psi(x^*, z^k) + \Psi(x^*, z^*). \quad (37)$$

From the update rule of z , we have following result.

Lemma 3.5 *If $x^{k+1} \in \mathcal{B}_\gamma(x^*)$ and $z^k \in \mathcal{B}_\gamma(z^*)$, then*

$$\sum_{j=1}^m (z_j^{k+1} - z_j^*) \cdot \max\left(-\frac{z_j^k}{\beta}, f_j(x^{k+1})\right) \leq \frac{1}{\rho_z} \|z^{k+1} - z^k\|^2 + \sum_{j \in J} (z_j^k - z_j^*) f_j(x^{k+1}). \quad (38)$$

Proof. When $x^{k+1} \in \mathcal{B}_\gamma(x^*)$ and $z^k \in \mathcal{B}_\gamma(z^*)$, it follows from (35) that

$$\beta f_j(x^{k+1}) + z_j^k > 0, \forall j \in J; \beta f_j(x^{k+1}) + z_j^k < 0, \forall j \notin J.$$

Hence,

$$\begin{aligned} & \sum_{j=1}^m (z_j^{k+1} - z_j^*) \cdot \max\left(-\frac{z_j^k}{\beta}, f_j(x^{k+1})\right) \\ &= \sum_{j \in J} (z_j^k + \rho_z f_j(x^{k+1}) - z_j^*) f_j(x^{k+1}) + \sum_{j \notin J} \left((1 - \frac{\rho_z}{\beta}) z_j^k\right) \left(-\frac{z_j^k}{\beta}\right) \\ &\leq \sum_{j \in J} (z_j^k + \rho_z f_j(x^{k+1}) - z_j^*) f_j(x^{k+1}) + \rho_z \sum_{j \notin J} \left(\frac{z_j^k}{\beta}\right)^2 \\ &= \frac{1}{\rho_z} \|z^{k+1} - z^k\|^2 + \sum_{j \in J} (z_j^k - z_j^*) f_j(x^{k+1}), \end{aligned}$$

which completes the proof. \square

Let $\underline{\eta} = L_g + \beta \|A\|^2 + \sum_{j=1}^m \beta B_j^2$. Then we have the next theorem.

Theorem 3.4 *Let $\{w^k\}$ be the sequence generated from Algorithm 1 with w^0 satisfying the following condition:*

$$\|x^0 - x^*\|^2 + \frac{1}{\rho_y \underline{\eta}} \|y^0 - y^*\|^2 + \frac{1}{\rho_z \underline{\eta}} \|z^0 - z^*\|^2 \leq \gamma^2 \cdot \min\left(1, \frac{1}{\rho_y \bar{\eta}}, \frac{1}{\rho_z \bar{\eta}}\right), \quad (39)$$

where γ is given in Proposition 3.1. For any $\theta \in (0, 1)$, if $0 < \rho_y \leq \beta$ and $0 < \rho_z \leq \beta(1 - \theta)$, then for any k , it holds that

$$\begin{aligned} & \frac{\theta \mu}{4} \|x^{k+1} - x^*\|^2 + \frac{\eta^k}{2} \|x^{k+1} - x^*\|^2 - \frac{\beta}{2} \|r^{k+1}\|^2 + \frac{1}{2\rho_y} \|y^{k+1} - y^*\|^2 + \frac{1}{2\rho_z} \|z^{k+1} - z^*\|^2 \\ & + \frac{1}{2} \|x^{k+1} - x^k\|_{(\eta^k - L_g - L_\Psi^k)I - \beta A^\top A}^2 + \left(\beta(1 - \frac{\theta}{2}) - \frac{\rho_y}{2}\right) \|r^{k+1}\|^2 \\ & \leq \frac{\eta^k}{2} \|x^k - x^*\|^2 - \frac{\beta}{2} \|r^k\|^2 + \frac{1}{2\rho_y} \|y^k - y^*\|^2 + \frac{1}{2\rho_z} \|z^k - z^*\|^2. \end{aligned} \quad (40)$$

Proof. We first note $\max(\|x^k - x^*\|, \|y^k - y^*\|, \|z^k - z^*\|) \leq \gamma, \forall k \geq 0$ from (31), (39), and the following inequality

$$\begin{aligned} & \min\left(1, \frac{1}{\rho_y \bar{\eta}}, \frac{1}{\rho_z \bar{\eta}}\right) (\|x^k - x^*\|^2 + \|y^k - y^*\|^2 + \|z^k - z^*\|^2) \\ & \leq \|x^k - x^*\|^2 + \frac{1}{\rho_y \eta^k} \|y^k - y^*\|^2 + \frac{1}{\rho_z \eta^k} \|z^k - z^*\|^2 \\ & \leq \|x^0 - x^*\|^2 + \frac{1}{\rho_y \underline{\eta}} \|y^0 - y^*\|^2 + \frac{1}{\rho_z \underline{\eta}} \|z^0 - z^*\|^2. \end{aligned}$$

Adding (19) with $y = y^*$, (20) with $z = z^*$, θ times of (37) and (38), and $1 - \theta$ times of (21) to (27) with $(x, y) = (x^*, y^*)$, we have by rearranging terms that

$$\begin{aligned} & \theta \left[f_0(x^{k+1}) - f_0(x^*) + \langle y^*, r^{k+1} \rangle + \frac{\beta}{2} \|r^{k+1}\|^2 + \Psi(x^{k+1}, z^*) - \Psi(x^*, z^*) \right] \\ & + (1 - \theta) \left[f_0(x^{k+1}) - f_0(x^*) + \langle y^*, r^{k+1} \rangle + \sum_{j=1}^m z_j^* f_j(x^{k+1}) - \Psi(x^*, z^k) \right] \\ & + \frac{\eta^k}{2} \|x^{k+1} - x^*\|^2 - \frac{\beta}{2} \|r^{k+1}\|^2 + \frac{1}{2\rho_y} \|y^{k+1} - y^*\|^2 + \frac{1}{2\rho_z} \|z^{k+1} - z^*\|^2 \\ & + \frac{1}{2} \|x^{k+1} - x^k\|_{(\eta^k - L_g - L_{\Psi}^k)I - \beta A^\top A}^2 + \left(\beta \left(1 - \frac{\theta}{2}\right) - \frac{\rho_y}{2} \right) \|r^{k+1}\|^2 + \frac{1}{2\rho_z} \left(\frac{\beta(1 - \theta)}{\rho_z} - 1 \right) \|z^{k+1} - z^k\|^2 \\ & \leq \frac{\eta^k}{2} \|x^k - x^*\|^2 - \frac{\beta}{2} \|r^k\|^2 + \frac{1}{2\rho_y} \|y^k - y^*\|^2 + \frac{1}{2\rho_z} \|z^k - z^*\|^2. \end{aligned}$$

From (6), (36), the above inequality, and $\Psi(x^*, z^k) \leq 0$, the desired result follows. \square

In addition, we can bound $\|y^k - y^*\|^2$ and $\|z^k - z^*\|^2$ by x -terms.

Lemma 3.6 *Let $\nu > 0$ be the smallest eigenvalue of $D^\top D$. Under the assumption of Theorem 3.4, we have*

$$\begin{aligned} & \nu (\|y^k - y^*\|^2 + \|z^k - z^*\|^2) \\ & \leq \left(4L_g^2 + 8|J| \sum_{j \in J} (\beta^2 B_j^4 + |z_j^k|^2 L_j^2) \right) \|x^k - x^*\|^2 + 4\beta^2 \|A\|^2 \|r^k\|^2 + 4\bar{\eta}^2 \|x^{k+1} - x^k\|^2. \end{aligned} \quad (41)$$

Proof. Note that

$$\nabla_x \Psi(x^k, z^k) = \sum_{j=1}^m [\beta f_j(x^k) + z_j^k]_+ \nabla f_j(x^k) = \sum_{j \in J} (\beta f_j(x^k) + z_j^k) \nabla f_j(x^k).$$

Hence, from the update of x and the fact $x^{k+1} \in \text{int}(\mathcal{X})$, it follows

$$\nabla g(x^k) + A^\top y^k + \beta A^\top r^k + \sum_{j \in J} (\beta f_j(x^k) + z_j^k) \nabla f_j(x^k) + \eta^k (x^{k+1} - x^k) = 0.$$

In addition, since $x^* \in \text{int}(\mathcal{X})$, it holds that

$$\nabla g(x^*) + A^\top y^* + \sum_{j \in J} z_j^* \nabla f_j(x^*) = 0.$$

From the above two equations, it follows that

$$\begin{aligned} & \|D[y^k; z^k] - D[y^*; z^*]\|^2 \\ = & \|\nabla g(x^k) - \nabla g(x^*) + \beta A^\top r^k + \sum_{j \in J} (\beta f_j(x^k) + z_j^k) \nabla f_j(x^k) - \sum_{j \in J} z_j^k \nabla f_j(x^*) + \eta^k(x^{k+1} - x^k)\|^2 \\ \leq & 4 \left(\|\nabla g(x^k) - \nabla g(x^*)\|^2 + \|\beta A^\top r^k\|^2 + \left\| \sum_{j \in J} [(\beta f_j(x^k) + z_j^k) \nabla f_j(x^k) - z_j^k \nabla f_j(x^*)] \right\|^2 + \|\eta^k(x^{k+1} - x^k)\|^2 \right) \\ \leq & 4L_g^2 \|x^k - x^*\|^2 + 4\beta^2 \|A\|^2 \|r^k\|^2 + 4\bar{\eta}^2 \|x^{k+1} - x^k\|^2 + 4 \left\| \sum_{j \in J} [(\beta f_j(x^k) + z_j^k) \nabla f_j(x^k) - z_j^k \nabla f_j(x^*)] \right\|^2. \end{aligned} \quad (42)$$

Note $f_j(x^*) = 0, \forall j \in J$. Hence,

$$\begin{aligned} & \left\| \sum_{j \in J} [(\beta f_j(x^k) + z_j^k) \nabla f_j(x^k) - z_j^k \nabla f_j(x^*)] \right\|^2 \\ \leq & |J| \sum_{j \in J} \left\| (\beta f_j(x^k) + z_j^k) \nabla f_j(x^k) - z_j^k \nabla f_j(x^*) \right\|^2 \\ = & |J| \sum_{j \in J} \left\| \beta (f_j(x^k) - f_j(x^*)) \nabla f_j(x^k) + z_j^k \nabla f_j(x^k) - z_j^k \nabla f_j(x^*) \right\|^2 \\ \leq & 2|J| \sum_{j \in J} (\beta^2 B_j^4 + |z_j^k|^2 L_j^2) \|x^k - x^*\|^2. \end{aligned}$$

Plugging in the above inequality into (42) and noting $\nu \| [y^k; z^k] - [y^*; z^*] \|^2 \leq \|D[y^k; z^k] - D[y^*; z^*]\|^2$, we obtain the desired result. \square

If necessary, taking a smaller γ , we can assume

$$\left| \sum_{j \in J} L_j (\beta f_j(\hat{x}) + \hat{z}_j) - \sum_{j \in J} L_j (\beta f_j(\tilde{x}) + \tilde{z}_j) \right| \leq \frac{\mu}{8}, \quad \forall \hat{x}, \tilde{x} \in \mathcal{B}_\gamma(x^*), \forall \hat{z}, \tilde{z} \in \mathcal{B}_\gamma(z^*). \quad (43)$$

Then we have the local linear convergence of Algorithm 1 as follows.

Theorem 3.5 (Local linear convergence) *Under Assumptions 2 and 3, let $\{w^k\}$ be the sequence generated from Algorithm 1 with w^0 satisfying (39), $\rho_y = \rho_z = \frac{\beta}{2}$, and $\delta > 0$. Let*

$$C = L_g^2 + 2|J| \left(\sum_{j \in J} \beta^2 B_j^4 + 2L_{\max}^2 (|z^*|^2 + \gamma^2) \right),$$

where $L_{\max} = \max_j L_j$. For any $\alpha > 0$ such that

$$\alpha < \min \left(\frac{\mu}{8C}, \frac{\delta}{\bar{\eta}^2}, \frac{1}{\beta \|A\|^2} \right), \quad (44)$$

it holds $\phi(x^{k+1}, y^{k+1}, z^{k+1}) \leq \sigma \cdot \phi(x^k, y^k, z^k)$, where

$$\phi(x^k, y^k, z^k) = \left(\frac{\mu}{16} + \frac{\eta^k}{2}\right) \|x^k - x^*\|^2 + \frac{1}{\beta} (\|y^k - y^*\|^2 + \|z^k - z^*\|^2)$$

and

$$\sigma = \max\left(\frac{\alpha C + \bar{\eta}}{\frac{\mu}{8} + \bar{\eta}}, 1 - \frac{\alpha\beta\nu}{8}\right) < 1.$$

Proof. Adding $\frac{\alpha}{8}$ of (41) to (40) with $\theta = \frac{1}{2}$, and noting $\alpha\bar{\eta}^2 I \preceq (\eta^k - L_g - L_{\Psi}^k)I - \beta A^\top A$, $\forall k$ and $\alpha\beta^2\|A\|^2 \leq \beta$ gives

$$\begin{aligned} & \left(\frac{\mu}{8} + \frac{\eta^k}{2}\right) \|x^{k+1} - x^*\|^2 + \frac{1}{\beta} (\|y^{k+1} - y^*\|^2 + \|z^{k+1} - z^*\|^2) \\ & \leq \left(\frac{\alpha C}{2} + \frac{\eta^k}{2}\right) \|x^k - x^*\|^2 + \left(\frac{1}{\beta} - \frac{\alpha\nu}{8}\right) (\|y^k - y^*\|^2 + \|z^k - z^*\|^2). \end{aligned} \quad (45)$$

Let

$$\eta_{\max} = \delta + L_g + \beta\|A\|^2 + \sum_{j=1}^m \beta B_j^2 + \max_{\substack{x \in \mathcal{B}_\gamma(x^*) \\ z \in \mathcal{B}_\gamma(z^*)}} \sum_{j \in J} L_j(\beta f_j(x) + z_j).$$

From the setting of η^k , we have that

$$\eta^k = \max(L_F^0, L_F^1, \dots, L_F^k) + \delta \leq \eta_{\max}, \forall k.$$

In addition, (43) indicates $\eta^{k+1} - \eta^k \leq \eta_{\max} - (L_F^k + \delta) \leq \frac{\mu}{8}$. Hence,

$$\frac{\mu}{16} + \frac{\eta^{k+1}}{2} \leq \frac{\mu}{8} + \frac{\eta^k}{2}. \quad (46)$$

Since $\frac{\alpha C}{2} + \frac{\eta^k}{2} \leq \sigma\left(\frac{\mu}{16} + \frac{\eta^k}{2}\right)$ and $\frac{1}{\beta} - \frac{\alpha\nu}{8} \leq \frac{\sigma}{\beta}$, we have the desired result from (45), (46), and the definition of ϕ . \square

Remark 3.1 In Theorem 3.5, the setting of $\rho_y = \rho_z = \frac{\beta}{2}$ is for simplicity of the analysis. The local linear convergence can be obtained for any $\rho_y, \rho_z \in (0, \beta)$. Therefore, from Theorems 3.2 and 3.5, the algorithm may eventually converge linearly. This phenomenon is observed from our numerical experiments; see Figures 1 and 2.

4 Block linearized augmented Lagrangian method

In this section, we assume that in (1), x can be partitioned into n disjoint blocks and the non-differentiable part $h(x)$ is separable, i.e.,

$$x = (x_1, x_2, \dots, x_n), \quad h(x) = \sum_{i=1}^n h_i(x_i).$$

Correspondingly, A can be written as the block matrix format $[A_1, \dots, A_n]$.

4.1 Algorithm

Towards a solution of the block structured problem, we propose a block linearized augmented Lagrangian method (BLALM). At each iteration, it randomly picks one block primal variable to update and then immediately renews the multipliers. The method is summarized in Algorithm 2.

Algorithm 2: Block linearized augmented Lagrangian method for (1)

1 Initialization: choose x^0, y^0, z^0 and $\beta, \rho_y, \rho_z, \boldsymbol{\eta} = [\eta_1, \dots, \eta_n]$; let $r^0 = Ax^0 - b$

2 for $k = 0, 1, \dots$ **do**

3 Pick $i_k \in [n]$ uniformly at random and perform the updates

$$x_i^{k+1} = \begin{cases} \arg \min_{x_i} h_i(x_i) + \langle \nabla_{x_i} F_\beta(w^k), x_i \rangle + \frac{\eta_i}{2} \|x_i - x_i^k\|^2, & \text{if } i = i_k \\ x_i^k, & \text{if } i \neq i_k \end{cases} \quad (47a)$$

$$\begin{aligned} r^{k+1} &= r^k + A_{i_k}(x_{i_k}^{k+1} - x_{i_k}^k), \\ y^{k+1} &= y^k + \rho_y r^{k+1}, \end{aligned} \quad (47b)$$

$$z_j^{k+1} = z_j^k + \rho_z \cdot \max\left(-\frac{z_j^k}{\beta}, f_j(x^{k+1})\right), j = 1, \dots, m. \quad (47c)$$

To make Algorithm 2 efficient, we require (1) to have the so-called coordinate friendly structure [16]. Roughly speaking, computing all n block partial gradients $\nabla_{x_i} F_\beta$ has nearly the same complexity as a full gradient evaluation. In addition, $f(x^{k+1})$ can be easily calculated from x^k , $f(x^k)$ and the change of x_{i_k} .

We let ℓ_i^k be the Lipschitz constant of $\nabla_{x_i} g(x) + \nabla_{x_i} \Psi(x, z^k)$ with respect to x_i for every $i = 1, \dots, n$ and $\boldsymbol{\ell}^k = [\ell_1^k, \dots, \ell_n^k]$. In general, ℓ_i^k can be significantly smaller than the Lipschitz constant of $\nabla g(x) + \nabla_x \Psi(x, z)$, and thus a larger stepsize can be made if a single block is updated instead of all blocks.

4.2 Convergence analysis

To show the convergence results of Algorithm 2, we first establish a fundamental result that is similar to Theorem 3.1.

Theorem 4.1 (One-iteration result of BLALM) *Let $\{w^k\}$ be the sequence generated from Al-*

gorithm 2. Then for any x such that $Ax = b$ and $f_j(x) \leq 0, \forall j \in [m]$, it holds

$$\begin{aligned} & \mathbb{E}_{i_k} \left[f_0(x^{k+1}) - f_0(x) + \langle y^{k+1}, r^{k+1} \rangle + (\beta - \rho_y) \|r^{k+1}\|^2 + \Psi_\beta(x^{k+1}, z^k) \right] \\ & + \frac{1}{2} \mathbb{E}_{i_k} \left[\|x^{k+1} - x\|_\eta^2 - \|x^k - x\|_\eta^2 + \|x^{k+1} - x^k\|_{\eta-\ell^k}^2 \right] - \frac{\beta}{2} \mathbb{E}_{i_k} \left[\|r^{k+1}\|^2 - \|r^k\|^2 + \|x^{k+1} - x^k\|_{A^\top A}^2 \right] \\ & \leq \left(1 - \frac{1}{n}\right) [f_0(x^k) - f_0(x) + \langle y^k, r^k \rangle + \beta \|r^k\|^2 + \Psi_\beta(x^k, z^k)]. \end{aligned} \quad (48)$$

Proof. From the update of x_{i_k} , we have

$$0 \in \partial h_{i_k}(x_{i_k}^{k+1}) + \nabla_{x_{i_k}} g(x^k) + A_{i_k}^\top (y^k + \beta r^k) + \nabla_{x_{i_k}} \Psi_\beta(x^k, z^k) + \eta_{i_k} (x_{i_k}^{k+1} - x_{i_k}^k). \quad (49)$$

Note that for any x ,

$$\begin{aligned} \mathbb{E}_{i_k} \langle x_{i_k}^{k+1} - x_{i_k}, \tilde{\nabla} h_{i_k}(x_{i_k}^{k+1}) \rangle & \geq \mathbb{E}_{i_k} [h_{i_k}(x_{i_k}^{k+1}) - h_{i_k}(x_{i_k})] \\ & = \mathbb{E}_{i_k} [h_{i_k}(x_{i_k}^{k+1}) - h_{i_k}(x_{i_k}^k) + h_{i_k}(x_{i_k}^k) - h_{i_k}(x_{i_k})] \\ & = \mathbb{E}_{i_k} [h(x^{k+1}) - h(x^k)] + \frac{1}{n} [h(x^k) - h(x)], \end{aligned} \quad (50)$$

and

$$\begin{aligned} & \mathbb{E}_{i_k} \langle x_{i_k}^{k+1} - x_{i_k}, \nabla_{x_{i_k}} g(x^k) + \nabla_{x_{i_k}} \Psi_\beta(x^k, z^k) \rangle \\ & = \mathbb{E}_{i_k} \langle x_{i_k}^{k+1} - x_{i_k}^k + x_{i_k}^k - x_{i_k}, \nabla_{x_{i_k}} g(x^k) + \nabla_{x_{i_k}} \Psi_\beta(x^k, z^k) \rangle \\ & \geq \mathbb{E}_{i_k} [g(x^{k+1}) + \Psi_\beta(x^{k+1}, z^k) - g(x) - \Psi_\beta(x^k, z^k) - \frac{1}{2} \|x^{k+1} - x^k\|_{\ell^k}^2] \\ & \quad + \frac{1}{n} [g(x^k) + \Psi_\beta(x^k, z^k) - g(x) - \Psi_\beta(x, z^k)]. \end{aligned} \quad (51)$$

In addition, for any x such that $Ax = b$, we have from [18, Lemma 3.2] that

$$\begin{aligned} \mathbb{E}_{i_k} \langle x_{i_k}^{k+1} - x_{i_k}, A_{i_k}^\top (y^k + \beta r^k) \rangle & = - \left(1 - \frac{1}{n}\right) (\langle y^k, r^k \rangle + \beta \|r^k\|^2) + \mathbb{E}_{i_k} [\langle y^k, r^{k+1} \rangle + \beta \|r^{k+1}\|^2] \\ & \quad - \frac{\beta}{2} \mathbb{E}_{i_k} [\|r^{k+1}\|^2 - \|r^k\|^2 + \|x^{k+1} - x^k\|_{A^\top A}^2]. \end{aligned} \quad (52)$$

Furthermore,

$$\mathbb{E}_{i_k} \langle x_{i_k}^{k+1} - x_{i_k}, \eta_{i_k} (x_{i_k}^{k+1} - x_{i_k}^k) \rangle = \frac{1}{2} \mathbb{E}_{i_k} [\|x^{k+1} - x\|_\eta^2 - \|x^k - x\|_\eta^2 + \|x^{k+1} - x^k\|_\eta^2]. \quad (53)$$

Adding (50) through (53), we have from (49) that

$$\begin{aligned} & \mathbb{E}_{i_k} \left[f_0(x^{k+1}) - f_0(x) + \langle y^k, r^{k+1} \rangle + \beta \|r^{k+1}\|^2 + \Psi_\beta(x^{k+1}, z^k) \right] \\ & + \frac{1}{2} \mathbb{E}_{i_k} \left[\|x^{k+1} - x\|_\eta^2 - \|x^k - x\|_\eta^2 + \|x^{k+1} - x^k\|_{\eta-\ell^k}^2 \right] - \frac{\beta}{2} \mathbb{E}_{i_k} \left[\|r^{k+1}\|^2 - \|r^k\|^2 + \|x^{k+1} - x^k\|_{A^\top A}^2 \right] \\ & \leq \left(1 - \frac{1}{n}\right) [f_0(x^k) - f_0(x) + \langle y^k, r^k \rangle + \beta \|r^k\|^2 + \Psi_\beta(x^k, z^k)] + \frac{1}{n} \Psi_\beta(x, z^k). \end{aligned}$$

Since $y^{k+1} = y^k + \rho_y r^{k+1}$ and $\Psi_\beta(x, z^k) \leq 0$, (48) is obtained from the above inequality. \square

We also need the next lemma.

Lemma 4.1 For any $\rho_z \leq \beta$,

$$-\frac{1}{\rho_z} \|z^{k+1} - z^k\|^2 \leq \Psi_\beta(x^{k+1}, z^k) - \Psi_\beta(x^{k+1}, z^{k+1}) \quad (54)$$

Proof. Since $\rho_z \leq \beta$, we have $z_j^k \geq 0, \forall k$. Let

$$J_1^k = \{j \in [m] : \beta f_j(x^{k+1}) + z_j^k \geq 0, \beta f_j(x^{k+1}) + z_j^{k+1} \geq 0\}, \quad (55a)$$

$$J_2^k = \{j \in [m] : \beta f_j(x^{k+1}) + z_j^k \geq 0, \beta f_j(x^{k+1}) + z_j^{k+1} < 0\}, \quad (55b)$$

$$J_3^k = \{j \in [m] : \beta f_j(x^{k+1}) + z_j^k < 0\}. \quad (55c)$$

For any $j \in J_1^k \cup J_2^k$, $z_j^{k+1} = z_j^k + \rho_z f_j(x^{k+1})$, and for any $j \in J_3^k$, $z_j^{k+1} = (1 - \frac{\rho_z}{\beta})z_j^k$ and $\beta f_j(x^{k+1}) + z_j^{k+1} < 0$. Hence,

$$\begin{aligned} & \Psi_\beta(x^{k+1}, z^k) - \Psi_\beta(x^{k+1}, z^{k+1}) \\ &= \sum_{j \in J_1^k \cup J_2^k} \left(z_j^k f_j(x^{k+1}) + \frac{\beta}{2} [f_j(x^{k+1})]^2 \right) - \sum_{j \in J_3^k} \frac{(z_j^k)^2}{2\beta} \\ & \quad - \sum_{j \in J_1^k} \left(z_j^{k+1} f_j(x^{k+1}) + \frac{\beta}{2} [f_j(x^{k+1})]^2 \right) + \sum_{j \in J_2^k \cup J_3^k} \frac{(z_j^{k+1})^2}{2\beta} \\ &= - \sum_{j \in J_1^k} \rho_z [f_j(x^{k+1})]^2 + \sum_{j \in J_2^k} \left(z_j^k f_j(x^{k+1}) + \frac{\beta}{2} [f_j(x^{k+1})]^2 + \frac{(z_j^{k+1})^2}{2\beta} \right) \\ & \quad - \sum_{j \in J_3^k} \frac{(z_j^k)^2 - (z_j^{k+1})^2}{2\beta} \end{aligned} \quad (56)$$

For $j \in J_2^k$, we have

$$\begin{aligned} & z_j^k f_j(x^{k+1}) + \frac{\beta}{2} [f_j(x^{k+1})]^2 + \frac{(z_j^{k+1})^2}{2\beta} \\ &= z_j^k f_j(x^{k+1}) + \frac{\beta}{2} [f_j(x^{k+1})]^2 + \frac{(z_j^k + \rho_z f_j(x^{k+1}))^2}{2\beta} \\ &= \frac{(z_j^k)^2}{2\beta} + \left(1 + \frac{\rho_z}{\beta}\right) z_j^k f_j(x^{k+1}) + \left(\frac{\beta}{2} + \frac{\rho_z^2}{2\beta}\right) [f_j(x^{k+1})]^2 \\ &\geq -\rho_z [f_j(x^{k+1})]^2, \end{aligned} \quad (57)$$

where the inequality follows from the Young's inequality. For $j \in J_3^k$, we have

$$-\frac{(z_j^k)^2 - (z_j^{k+1})^2}{2\beta} = -\left(1 - \left(1 - \frac{\rho_z}{\beta}\right)^2\right) \frac{(z_j^k)^2}{2\beta} \geq -\frac{\rho_z}{\beta^2} (z_j^k)^2. \quad (58)$$

Plugging (57) and (58) into (56) gives

$$\Psi_\beta(x^{k+1}, z^k) - \Psi_\beta(x^{k+1}, z^{k+1}) \geq - \sum_{j \in J_1^k \cup J_2^k} \rho_z [f_j(x^{k+1})]^2 - \sum_{j \in J_3^k} \frac{\rho_z}{\beta^2} (z_j^k)^2 = -\frac{1}{\rho_z} \|z^{k+1} - z^k\|^2,$$

which completes the proof. \square

The following results are easy to show from the Young's inequality and the update rule of z .

Lemma 4.2 *For any y and $z \geq 0$,*

$$\langle y, r^{k+1} \rangle \leq \langle y^k, r^{k+1} \rangle + \frac{\beta}{2} \|r^{k+1}\|^2 + \frac{1}{2\beta} \|y^k - y\|^2, \quad (59)$$

and

$$0 \leq \Psi_\beta(x^{k+1}, z^k) - \sum_{j=1}^m z_j f_j(x^{k+1}) + \frac{1}{2\beta} \|z^k - z\|^2 \quad (60)$$

Proof. The inequality in (59) directly follows from the Young's inequality.

Let $J_+^k = J_1^k \cup J_2^k$ and $J_-^k = J_3^k$, where J_1^k, J_2^k and J_3^k are defined in (55). Then

$$\begin{aligned} & \Psi_\beta(x^{k+1}, z^k) - \sum_{j=1}^m z_j f_j(x^{k+1}) + \frac{1}{2\beta} \|z^k - z\|^2 \\ &= \sum_{j \in J_+^k} \left[(z_j^k - z_j) f_j(x^{k+1}) + \frac{\beta}{2} [f_j(x^{k+1})]^2 + \frac{1}{2\beta} (z_j^k - z_j)^2 \right] + \sum_{j \in J_-^k} \left[-\frac{(z_j^k)^2}{2\beta} - z_j f_j(x^{k+1}) + \frac{1}{2\beta} (z_j^k - z_j)^2 \right] \\ &\geq \sum_{j \in J_-^k} \left[-\frac{(z_j^k)^2}{2\beta} - z_j f_j(x^{k+1}) + \frac{1}{2\beta} (z_j^k - z_j)^2 \right] \geq \sum_{j \in J_-^k} \frac{1}{2\beta} (z_j^k)^2, \end{aligned}$$

where the first inequality follows from the Young's inequality, and the second one holds because $f_j(x^{k+1}) \leq -\frac{z_j^k}{\beta}, \forall j \in J_-^k$ and $z_j \geq 0, \forall j$. This completes the proof. \square

Using the previous establish results, we are now able to show the convergence rate of Algorithm 2.

Theorem 4.2 (Sublinear convergence of BLALM) *Under Assumptions 1 and 2, let $\{w^k\}$ be the sequence from Algorithm 2 with $y^0 = 0$ and $z^0 = 0$. Assume ℓ_i^k is upper bounded by $\bar{\ell}_i$ for any $i \in [n]$ and any k . If $\rho_y \in (0, \frac{\beta}{n}]$, $\rho_z \in (0, \frac{\beta}{2n}]$, and $\eta_i \geq \bar{\ell}_i + \beta \|A_i\|^2, \forall i \in [n]$, then*

$$|\mathbb{E}[f_0(\bar{x}^{k+1}) - f_0(x^*)]| \leq \frac{1}{1 + \frac{k}{n}} \left(C_{x^0} + \frac{2\|y^*\|^2}{n\rho_y} + \frac{2\|z^*\|^2}{n\rho_z} \right), \quad (61a)$$

$$\mathbb{E} \left[\|A\bar{x}^{k+1} - b\| + \sum_{j=1}^m [f_j(\bar{x}^{k+1})]_+ \right] \leq \frac{1}{1 + \frac{k}{n}} \left(C_{x^0} + \frac{(1 + \|y^*\|)^2}{2n\rho_y} + \frac{1}{2n\rho_z} \sum_{j=1}^m (1 + z_j^*)^2 \right), \quad (61b)$$

where $\bar{x}^{k+1} = \frac{1}{1+\frac{k}{n}} \sum_{t=0}^k x^{t+1}$, and

$$C_{x^0} = \left(1 - \frac{1}{n}\right) \left[f_0(x^0) - f_0(x^*) + \frac{\beta}{2} \left(\|r^0\|^2 + \sum_{j=1}^m [f_j(x^0)]_+^2 \right) \right] + \frac{1}{2} \|x^0 - x^*\|_{\boldsymbol{\eta}}^2.$$

Proof. Since $\eta_i \geq \bar{\ell}_i + \beta \|A_i\|^2$, $\forall i \in [n]$ and $x_i^{k+1} = x_i^k$, $\forall i \neq i_k$, it holds

$$\|x^{k+1} - x^k\|_{\boldsymbol{\eta} - \boldsymbol{\ell}^k}^2 \geq \beta \|x^{k+1} - x^k\|_{A^\top A}^2.$$

Hence, taking expectation on both sides of (48) with $x = x^*$ and summing it up give

$$\begin{aligned} & \mathbb{E} \left[f_0(x^{k+1}) - f_0(x^*) + \langle y^k, r^{k+1} \rangle + \Psi_\beta(x^{k+1}, z^k) \right] + \frac{\beta}{2} \mathbb{E} \|r^{k+1}\|^2 \\ & + \frac{1}{n} \sum_{t=0}^{k-1} \mathbb{E} \left[f_0(x^{t+1}) - f_0(x^*) + \langle y^{t+1}, r^{t+1} \rangle + \Psi_\beta(x^{t+1}, z^t) \right] + \left(\frac{\beta}{n} - \rho_y \right) \sum_{t=0}^{k-1} \mathbb{E} \|r^{t+1}\|^2 \\ & + \left(1 - \frac{1}{n}\right) \sum_{t=0}^{k-1} \mathbb{E} \left[\Psi_\beta(x^{t+1}, z^t) - \Psi_\beta(x^{t+1}, z^{t+1}) \right] + \frac{1}{2} \mathbb{E} \|x^{k+1} - x^*\|_{\boldsymbol{\eta}}^2 \\ & \leq \left(1 - \frac{1}{n}\right) \left[f_0(x^0) - f_0(x^*) + \langle y^0, r^0 \rangle + \frac{\beta}{2} \|r^0\|^2 + \Psi_\beta(x^0, z^0) \right] + \frac{1}{2} \|x^0 - x^*\|_{\boldsymbol{\eta}}^2 - \frac{\beta}{2} \|r^0\|^2, \quad (62) \end{aligned}$$

where in the first line, we have used $y^{k+1} = y^k - \rho_y r^{k+1}$. Summing (19), (20), and (21) gives

$$\begin{aligned} & \frac{1}{2\rho_y} \left[\|y^k - y\|^2 - \|y^0 - y\|^2 + \sum_{t=0}^{k-1} \|y^{t+1} - y^t\|^2 \right] - \sum_{t=0}^{k-1} \langle y^{t+1} - y, r^{t+1} \rangle \\ & + \frac{1}{2\rho_z} \left[\|z^k - z\|^2 - \|z^0 - z\|^2 + \sum_{t=0}^{k-1} \|z^{t+1} - z^t\|^2 \right] + \frac{\beta - 2\rho_z}{2\rho_z^2} \sum_{t=0}^{k-1} \|z^{t+1} - z^t\|^2 \\ & \leq \sum_{t=0}^{k-1} \left[\Psi_\beta(x^{t+1}, z^t) - \sum_{j=1}^m z_j f_j(x^{t+1}) \right]. \end{aligned}$$

Since $y^0 = 0$ and $z^0 = 0$, adding $\frac{1}{n}$ of the above inequality to (62), using Lemma 4.1, and noting $\frac{\beta}{n} \geq \rho_y$, $\frac{\beta - \rho_z}{2n\rho_z^2} \geq \frac{1 - \frac{1}{n}}{\rho_z}$ from the choice of ρ_y, ρ_z , we have

$$\begin{aligned} & \mathbb{E} \left[f_0(x^{k+1}) - f_0(x^*) + \langle y^k, r^{k+1} \rangle + \Psi_\beta(x^{k+1}, z^k) \right] + \frac{\beta}{2} \mathbb{E} \|r^{k+1}\|^2 + \frac{1}{2n\rho_y} \mathbb{E} \|y^k - y\|^2 \\ & + \frac{1}{2} \mathbb{E} \|x^{k+1} - x^*\|_{\boldsymbol{\eta}}^2 + \frac{1}{2n\rho_z} \mathbb{E} \|z^k - z\|^2 + \frac{1}{n} \sum_{t=0}^{k-1} \mathbb{E} [\Phi(x^{t+1}; x^*, y, z)] \\ & \leq \left(1 - \frac{1}{n}\right) \left[f_0(x^0) - f_0(x^*) + \beta \|r^0\|^2 + \Psi_\beta(x^0, 0) \right] + \frac{1}{2} \|x^0 - x^*\|_{\boldsymbol{\eta}}^2 - \frac{\beta}{2} \|r^0\|^2 \\ & + \frac{1}{2n\rho_y} \|y\|^2 + \frac{1}{2n\rho_z} \|z\|^2. \quad (63) \end{aligned}$$

Note $\Psi_\beta(x^0, 0) = \sum_{j=1}^m [f_j(x^0)]_+^2$. Since $\rho_y, \rho_z \leq \frac{\beta}{n}$, plugging (59) and (60) into (63) and using the convexity of f_i 's yield

$$\mathbb{E}[\Phi(\bar{x}^{k+1}; x^*, y, z)] \leq \frac{1}{1 + \frac{k}{n}} \left(C_{x^0} + \frac{1}{2n\rho_y} \|y\|^2 + \frac{1}{2n\rho_z} \|z\|^2 \right).$$

Therefore, we complete the proof by Lemma 2.4. \square

Remark 4.1 *If n block updates of x costs roughly the same as one full update to x , then the results in (61) are comparable to those in (34) by noting their differences in choosing ρ_y, ρ_z . One drawback of Theorem 4.2 is the assumption on the upper bound of ℓ^k . From (13), we see that the upper bound can be pre-calculated if $f_j(x), \forall j \in [m]$ are affine. However, in general, it is unknown and dependent on the iterates. Numerically, we can gradually increase η_i by a fixed amount or ratio if $\eta_i < \ell_i^k + \beta \|A_i\|^2$ is detected or by backtracking until the following inequality holds:*

$$F_\beta(x^{k+1}, y^k, z^k) \leq F_\beta(w^k) + \langle \nabla_{x_{i_k}} F_\beta(w^k), x_{i_k}^{k+1} - x_{i_k}^k \rangle + \frac{\eta_{i_k}}{2} \|x_{i_k}^{k+1} - x_{i_k}^k\|^2. \quad (64)$$

After finitely many increases, $\ell_i^k + \beta \|A_i\|^2 \leq \eta_i, \forall i$, will hold in high probability for every k . This can be explained by the following arguments.

Let $\eta_i = \zeta \geq 1, \forall i \in [n]$. Since $n\rho_z \leq \frac{\beta}{2}$, then from (6), (63) with $(y, z) = (y^*, z^*)$, (59), and (60), it follows that

$$\begin{aligned} \mathbb{E}\|x^k - x^*\| &\leq \sqrt{\frac{C_{x^0}}{\zeta}} + \|x^0 - x^*\| + \frac{\|y^*\|}{\sqrt{n\rho_y\zeta}} + \frac{\|z^*\|}{\sqrt{n\rho_z\zeta}}, \\ \mathbb{E}\|z^k - z^*\| &\leq \sqrt{\beta C_{x^0}} + \sqrt{\beta\zeta} \|x^0 - x^*\| + \sqrt{\frac{\beta}{n\rho_y}} \|y^*\| + \sqrt{\frac{\beta}{n\rho_z}} \|z^*\|. \end{aligned}$$

Hence, we have from (29) that $\mathbb{E}[\ell_i^k] = O(\sqrt{\zeta})$. By the Markov inequality, for every $k, \ell_i^k \leq \eta_i, \forall i \in [n]$ holds in high probability if $\eta_i \gg \zeta, \forall i \in [n]$, where $\zeta \geq 1$ satisfies

$$\begin{aligned} \zeta &\geq L_g + \beta \|A\|^2 + \beta \sum_{j=1}^m B_j^2 + \beta \sum_{j=1}^m B_j L_j \left(\sqrt{C_{x^0}} + \|x^0 - x^*\| + \frac{\|y^*\|}{\sqrt{n\rho_y}} + \frac{\|z^*\|}{\sqrt{n\rho_z}} \right) \\ &\quad + \sqrt{\sum_{j=1}^m L_j^2} \left(\|z^*\| + \sqrt{2\beta C_{x^0}} + \sqrt{2\beta\zeta} \|x^0 - x^*\| + \sqrt{\frac{2\beta}{n\rho_y}} \|y^*\| + \sqrt{\frac{2\beta}{n\rho_z}} \|z^*\| \right). \end{aligned}$$

Remark 4.2 *If Assumption 3 is satisfied, we can also show a local linear convergence result of Algorithm 2 following the analysis in section 3.2 and that in [20]. We do not expand details here but leave it to interested readers.*

5 Applications

In this section, we give a few applications that can be formulated in the form of (1) and discuss how Algorithm 1 and/or Algorithm 2 can be applied.

5.1 Basis pursuit denoising

Suppose we observe a noisy measurement $b = A\theta^o + \xi$ of a signal θ^o , where A is a measuring matrix, and ξ is a noise vector. Assume θ^o can be sparsely represented by a dictionary D . Then we can recover the signal through solving the so-called basis pursuit denoising (BPDN) problem:

$$\min_x \|x\|_1, \text{ s.t. } \|ADx - b\|^2 \leq \delta, \quad (65)$$

where δ measures the noise level. Upon obtaining a solution x^* to (65), we let $\theta^r = Dx^*$ be the recovered signal. Depending on the application, one can impose certain bounds on x to make the recovered signal physically meaningful. In this case, all conditions in Assumption 2 holds. In addition, assuming $b \in \text{Range}(AD)$, then Slater condition holds, and thus Assumption 1 is satisfied. Hence, Algorithm 1 is applicable, and the x -subproblem (18a) has closed-form solution by shrinkage or soft-thresholding. If A and D are stored as matrices, (65) is coordinate friendly, and we can also use Algorithm 2. However, for certain signal processing problems, evaluating $A\theta$ and/or Dx may not require explicit form of A or D but can be efficiently realized, such as a partial circulant A and/or a discrete cosine dictionary D . For this case, Algorithm 2 will not be as efficient as Algorithm 1 since evaluating coordinate gradient of $\|ADx - b\|^2$ may require full gradient.

5.2 Quadratically constrained quadratic programming

The quadratically constrained quadratic programming (QCQP) can be formulated as

$$\begin{aligned} \min_{x \in \mathbb{R}^p} \quad & \frac{1}{2}x^\top Q_0x + c_0^\top x + d_0 \\ \text{s.t.} \quad & \frac{1}{2}x^\top Q_jx + c_j^\top x + d_j \leq 0, \forall j \in [m], \\ & Ax = b, \\ & l_i \leq x_i \leq u_i, \forall i \in [p]. \end{aligned} \quad (66)$$

Let $\mathcal{X} = [l_1, u_1] \times \cdots \times [l_p, u_p]$ and $h(x) = \iota_{\mathcal{X}}(x)$. Then (66) can be written as (1) by adding $h(x)$ into the objective. When every Q_j is positive semidefinite, the problem is convex, and if all l_i 's and u_i 's are finite, then \mathcal{X} is bounded and all conditions in Assumption 2 hold. Hence, we can apply Algorithm 1 to find a solution of (66), and the solution of x -subproblem (18a) can be explicitly given by performing projection to a box constraint. In addition, the problem is coordinate friendly since evaluating the partial derivative of $\frac{1}{2}x^\top Q_jx + c_j^\top x + d_j$ about each x_i costs roughly $\frac{1}{p}$ of computing the full gradient. Furthermore, if we maintain Q_jx , then calculating the function value is negligible compared to the gradient computation. Therefore, we can also apply Algorithm 2 to the QCQP.

5.3 Finite minimax problems

Many applications can be formulated as a finite minimax problem (e.g., see [15] and the references therein):

$$\min_{x \in \mathcal{X}} \max_{1 \leq j \leq m} f_j(x), \quad (67)$$

where each f_j is a smooth convex function. Although all f_j 's are differentiable, the objective of (67) is generally not differentiable due to the max operation. Introducing variable t and requiring $\max_{1 \leq j \leq m} f_j(x) \leq t$, one can express the minimax problem equivalently to

$$\min_{x \in \mathcal{X}, t} t, \text{ s.t. } f_j(x) - t \leq 0, \forall j \in [m]. \quad (68)$$

For any $x \in \text{int}(\mathcal{X})$, each inequality constraint holds strictly at $(x, \max_j f_j(x) + 1)$, and thus the Slater condition holds. Hence, Assumption 1 is satisfied. In addition, if \mathcal{X} is bounded, then all conditions in Assumption 2 also hold. Therefore, we can use Algorithm 1 to find a solution of (68) and equivalently (67), and every iteration requires performing a projection to \mathcal{X} . Depending on applications, one may also apply Algorithm 2 if the problem is coordinate friendly, for example, every f_j is a quadratic function.

6 Numerical experiments

In this section, we test Algorithms 1 and 2 on BPDN (65) and QCQP (66) to show their numerical performance. The two algorithms are named as LALM and BLALM respectively. For both algorithms, we choose the parameter η by backtracking. More precisely, at each iteration k , for LALM, we start from $\eta^k = \eta^{k-1}$ and multiply it by 1.5 if (17) fails, and for BLALM, we initialize $\eta_{i_k}^k = \eta_{i_k}^{k-1}$ and multiply it by 1.5 if (64) does not hold. For both tests, we run the compared methods to 10^5 epochs, where one epoch is equivalent to n block updates. Optimal solutions to both tested problems are computed by CVX [8] with high precision.

6.1 Basis pursuit denoising

In this test, we show the convergence speed of LALM and BLALM on solving BPDN (65). For simplicity, we set $D = I$. The matrix $A \in \mathbb{R}^{50 \times 100}$ is randomly generated according to the standard Gaussian distribution, and the underlying sparse signal x^o has 5 nonzero components following the standard Gaussian distribution. Then we let $b = Ax^o + 0.1\xi$, where ξ is a unit Gaussian noise vector. For BLALM, we evenly partition the variable x into 10 blocks. The parameter β is simply set to 1 for both methods, and $\rho_z = \beta$ is set for LALM and $\rho_z = \frac{\beta}{10}$ for BLALM. Note that for the latter, the value of ρ_z is larger than that given by the theorem, and the algorithm still works well. This may indicate that our analysis is not tight.

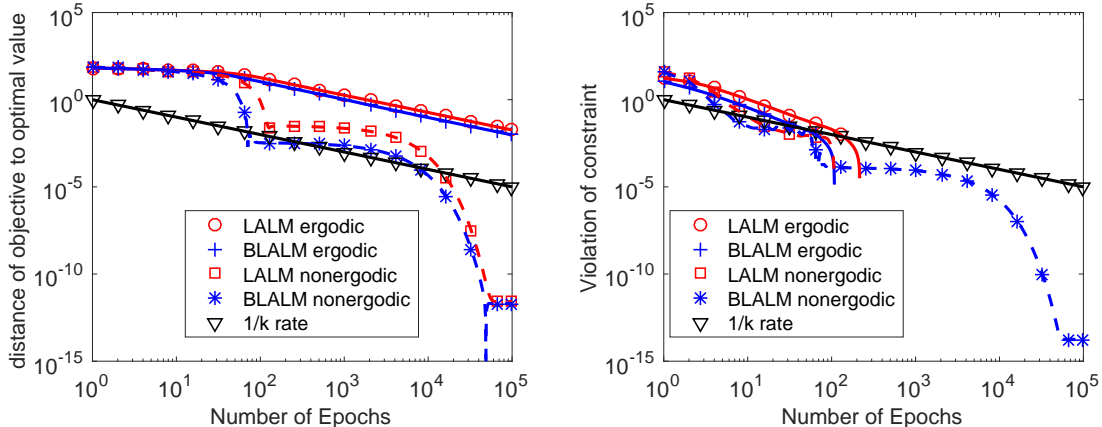


Figure 1: Convergence behaviors of Algorithm 1 (named LALM) and Algorithm 2 (named BLALM) on the BPDN problem (65). Left: distance of objective value to the optimal value $|f_0(x) - f_0(x^*)|$; Right: constraint residual $\sum_{j=1}^m [f_j(x)]_+$. “ergodic” curves are measured by averaged iterate \bar{x}^k and “nonergodic” ones by actual iterate x^k . The missing part on each constraint violation curve corresponds to zero residual.

Figure 1 plots the objective values and constraint residuals produced by both algorithms, where the curve corresponding to “ergodic” is obtained by using the averaged iterates \bar{x}^k and “nonergodic” by the actual iterate x^k . The missing part on each constraint violation curve corresponds to zero residual. Since LALM and BLALM have similar per-epoch complexity, their comparison in terms of running time is similar to that in Figure 1. From the figure, we see that BLALM is better than LALM in terms of both ergodic and nonergodic iterates. The ergodic convergence speed of both methods is precisely the order of $\frac{1}{k}$ and matches our theorems. However, the nonergodic convergence is significantly faster, especially as the iterate approaches to optimality. This is possibly because the iterate enters a region where the algorithm has linear convergence as indicated by the analysis in section 3.2. For this reason, we use the actual iterate in the next test.

6.2 Quadratically constrained quadratic programming

In this subsection, we test LALM and BLALM on the QCQP problem (66) and compare them to the recently proposed first-order primal-dual type method by Yu&Neely [21]. They consider a smooth constrained convex program in the form of (1) without an explicit linear equality constraint. Their method that we name as PD-YN iteratively performs the updates:

$$x^{k+1} = \mathcal{P}_{\mathcal{X}} \left(x^k - \frac{1}{\eta} \left[\nabla f_0(x^k) + \sum_{j=1}^m (\lambda_j^k + f_j(x^k)) \nabla f_j(x^k) \right] \right), \quad (69a)$$

$$\lambda_j^{k+1} = \max(-f_j(x^k), \lambda_j^k + f_j(x^k)), \forall j \in [m], \quad (69b)$$

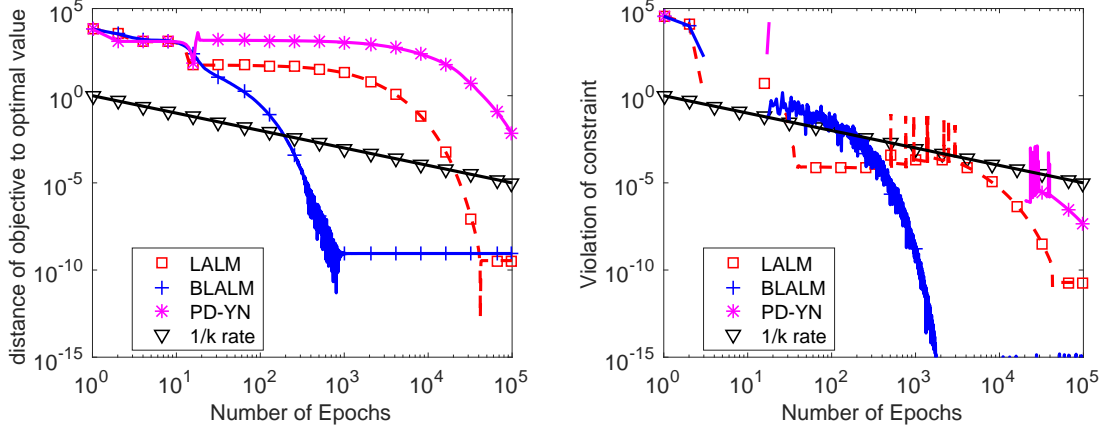


Figure 2: Convergence behaviors of Algorithm 1 (named LALM), Algorithm 2 (named BLALM), and the primal-dual type method (named PD-YN) in [21] on the QCQP problem (66). Left: distance of objective value to the optimal value $|f_0(x) - f_0(x^*)|$; Right: constraint residual $\sum_{j=1}^m [f_j(x)]_+$. The missing part on each constraint violation curve corresponds to zero residual.

where $\lambda_j^0 = \max(0, -f_j(x^0))$, $\forall j \in [m]$, and η is the step size. In the test, we also choose η adaptively by backtracking such that

$$\phi(x^{k+1}, z^k) \leq \phi(x^k, z^k) + \langle \nabla_x \phi(x^k, z^k), x^{k+1} - x^k \rangle + \frac{\eta^k}{2} \|x^{k+1} - x^k\|^2,$$

where $\phi(x, z) = f_0(x) + \sum_{j=1}^m z_j f_j(x)$ and $z_j^k = \lambda_j^k + f_j(x^k)$. Although [21] does not show the convergence of PD-YN with the above adaptive η^k , we observe its better performance than that with a fixed η .

The problem size is set to $m = 10$ and $p = 2000$. We randomly generate SPD matrices $Q_j, j = 0, 1, \dots, m$. A is set to zero, i.e., there is no linear equality constraint. The vector c_j 's are generated according to the standard Gaussian distribution, and $d_0 = 0$ and each d_j is a negative number for $j \in [m]$. Also we set $l_i = -10$ and $u_i = 10$ for each $i \in [p]$. Hence, the zero vector is an interior point of \mathcal{X} and makes every inequality hold strictly, namely, the Slater condition holds. We set $\beta = 0.1$ for both LALM and BLALM, and for the latter, we evenly partition the variable into 200 blocks. The parameter ρ_z is set to β and $\frac{\beta}{200}$ for the two algorithms respectively.

Figure 2 plots the results by the three compared algorithms. Both the proposed methods perform significantly better than PD-YN, and BLALM is the best among the three. In addition, we notice that LALM and BLALM converge linearly when the iterate approaches to optimality, as indicated by Theorem 3.5.

7 Conclusions

We have presented a first-order method for solving composite convex programming with both equality and smooth nonlinear inequality constraints. The method is derived from proximal linearization of the classic augmented Lagrangian function. Its global iterate convergence and global sublinear and local linear convergence results have been established. For the problem that has coordinate friendly structure, we have also proposed a first-order randomized block update method and shown its global sublinear convergence in expectation. In addition, we have implemented the two methods on solving the basis pursuit denoising problem and the convex quadratically constrained quadratic programming. Global sublinear and local linear convergence are both observed in the numerical experiments.

References

- [1] A. Ben-Tal and M. Zibulevsky. Penalty/barrier multiplier methods for convex programming problems. *SIAM Journal on Optimization*, 7(2):347–366, 1997. [2](#)
- [2] D. P. Bertsekas. *Nonlinear programming*. Athena scientific Belmont, 1999. [3](#)
- [3] D. P. Bertsekas. *Constrained optimization and Lagrange multiplier methods*. Academic press, 2014. [3](#)
- [4] P. L. Combettes and J.-C. Pesquet. Stochastic quasi-fejér block-coordinate fixed point iterations with random sweeping. *SIAM Journal on Optimization*, 25(2):1221–1248, 2015. [7](#)
- [5] Y. Cui, X. Li, D. Sun, and K.-C. Toh. On the convergence properties of a majorized alternating direction method of multipliers for linearly constrained convex optimization problems with coupled objective functions. *Journal of Optimization Theory and Applications*, 169(3):1013–1041, 2016. [3](#)
- [6] X. Gao, Y. Xu, and S. Zhang. Randomized primal-dual proximal block coordinate updates. *arXiv preprint arXiv:1605.05969*, 2016. [3](#)
- [7] X. Gao and S.-Z. Zhang. First-order algorithms for convex optimization with nonseparable objective and coupled constraints. *Journal of the Operations Research Society of China*, 5(2):131–159, 2017. [3](#)
- [8] M. Grant, S. Boyd, and Y. Ye. CVX: Matlab software for disciplined convex programming, 2008. [25](#)
- [9] 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. *arXiv preprint arXiv:1401.7079*, 2014. [3](#)

- [10] G. Lan, D. Renato, and C. Monteiro. Iteration-complexity of first-order augmented lagrangian methods for convex programming. *Mathematical Programming*, 155(1-2):511–547, 2016. [3](#)
- [11] Y.-F. Liu, X. Liu, and S. Ma. On the non-ergodic convergence rate of an inexact augmented lagrangian framework for composite convex programming. *arXiv preprint arXiv:1603.05738*, 2016. [3](#)
- [12] V. Nedelcu, I. Necoara, and Q. Tran-Dinh. Computational complexity of inexact gradient augmented lagrangian methods: application to constrained mpc. *SIAM Journal on Control and Optimization*, 52(5):3109–3134, 2014. [3](#)
- [13] A. Nedić and A. Ozdaglar. Subgradient methods for saddle-point problems. *Journal of optimization theory and applications*, 142(1):205–228, 2009. [3](#)
- [14] Y. Nesterov. Gradient methods for minimizing composite functions. *Mathematical Programming*, 140(1):125–161, 2013. [3](#)
- [15] E. Pee and J. O. Royset. On solving large-scale finite minimax problems using exponential smoothing. *Journal of optimization theory and applications*, 148(2):390–421, 2011. [25](#)
- [16] Z. Peng, T. Wu, Y. Xu, M. Yan, and W. Yin. Coordinate-friendly structures, algorithms and applications. *Annals of Mathematical Sciences and Applications*, 1(1):57–119, 2016. [18](#)
- [17] Y. Xu. Accelerated first-order primal-dual proximal methods for linearly constrained composite convex programming. *SIAM Journal on Optimization*, 27(3):1459–1484, 2017. [3](#)
- [18] Y. Xu. Asynchronous parallel primal-dual block update methods. *arXiv preprint arXiv:1705.06391*, 2017. [19](#)
- [19] Y. Xu. Global convergence rates of augmented lagrangian methods for constrained convex programming. *preprint*, 2017. [2](#), [3](#), [7](#), [8](#)
- [20] Y. Xu and S. Zhang. Accelerated primal-dual proximal block coordinate updating methods for constrained convex optimization. *arXiv preprint arXiv:1702.05423*, 2017. [3](#), [23](#)
- [21] H. Yu and M. J. Neely. A primal-dual type algorithm with the $O(1/t)$ convergence rate for large scale constrained convex programs. In *Decision and Control (CDC), 2016 IEEE 55th Conference on*, pages 1900–1905. IEEE, 2016. [3](#), [26](#), [27](#)