

Convergence on a Symmetric Accelerated Stochastic ADMM with Larger Stepsizes

Jianchao Bai¹, Deren Han^{2,*}, Hao Sun³ and Hongchao Zhang⁴

¹ School of Mathematics and Statistics, the MIIT Key Laboratory of Dynamics and Control of Complex Systems, Northwestern Polytechnical University, Xi'an, 710129, P.R. China.

² LMIB of the Ministry of Education, School of Mathematical Sciences, Beihang University, Beijing, 100191, P.R. China.

³ School of Mathematics and Statistics, Northwestern Polytechnical University, Xi'an, 710129, P.R. China.

⁴ Department of Mathematics, Louisiana State University, Baton Rouge, LA 70803-4918, USA.

Received 30 March 2021; Accepted 19 December 2021

Abstract. In this paper, we develop a symmetric accelerated stochastic Alternating Direction Method of Multipliers (SAS-ADMM) for solving separable convex optimization problems with linear constraints. The objective function is the sum of a possibly nonsmooth convex function and an average function of many smooth convex functions. Our proposed algorithm combines both ideas of ADMM and the techniques of accelerated stochastic gradient methods possibly with variance reduction to solve the smooth subproblem. One main feature of SAS-ADMM is that its dual variable is symmetrically updated after each update of the separated primal variable, which would allow a more flexible and larger convergence region of the dual variable compared with that of standard deterministic or stochastic ADMM. This new stochastic optimization algorithm is shown to have ergodic converge in expectation with $\mathcal{O}(1/T)$ convergence rate, where T denotes the number of outer iterations. Our preliminary experiments indicate the proposed algorithm is very effective for solving separable optimization problems from big-data applications. Finally, 3-block extensions of the algorithm and its variant of an accelerated stochastic augmented Lagrangian method are discussed in the appendix.

AMS subject classifications: 65K10, 65Y20, 68W40, 90C25

Key words: Convex optimization, stochastic ADMM, symmetric ADMM, larger stepsize, proximal mapping, complexity.

*Corresponding author. Email addresses: jianchaobai@nwpu.edu.cn (J. Bai), handr@buaa.edu.cn (D. Han), hsun@nwpu.edu.cn (H. Sun), hozhang@math.lsu.edu (H. Zhang)

1 Introduction

We consider the following structured composite convex optimization problem with linear equality constraints:

$$\min\{f(\mathbf{x}) + g(\mathbf{y}) \mid \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}, A\mathbf{x} + B\mathbf{y} = \mathbf{b}\}, \quad (1.1)$$

where $\mathcal{X} \subset \mathbb{R}^{n_1}$, $\mathcal{Y} \subset \mathbb{R}^{n_2}$ are closed convex subsets, $A \in \mathbb{R}^{n \times n_1}$, $B \in \mathbb{R}^{n \times n_2}$, $\mathbf{b} \in \mathbb{R}^n$ are given, $g: \mathcal{Y} \rightarrow \mathbb{R} \cup \{+\infty\}$ is a convex but possibly nonsmooth function, and f is an average of N real-valued convex functions:

$$f(\mathbf{x}) = \frac{1}{N} \sum_{j=1}^N f_j(\mathbf{x}).$$

We assume that each f_j defined on an open set containing \mathcal{X} is Lipschitz continuously differentiable on \mathcal{X} . Problem (1.1) is also referred as the regularized empirical risk minimization in big-data applications [26, 35], including classification and regression models in machine learning, where N denotes the sample size and f_j corresponds to the empirical loss. A major difficulty for solving (1.1) is that the sample size N can be very large such that it is often computationally prohibitive to evaluate either the full function value or the gradient of f at each iteration of an algorithm. Hence, it is essential for an effective algorithm, e.g., a stochastic gradient method, to explore the summation structure of f in the objective function.

The augmented Lagrangian function of (1.1) is

$$\mathcal{L}_\beta(\mathbf{x}, \mathbf{y}, \boldsymbol{\lambda}) = \mathcal{L}(\mathbf{x}, \mathbf{y}, \boldsymbol{\lambda}) + \frac{\beta}{2} \|A\mathbf{x} + B\mathbf{y} - \mathbf{b}\|^2, \quad (1.2)$$

where $\beta > 0$ is a penalty parameter, $\boldsymbol{\lambda}$ is the Lagrange multiplier and the Lagrangian of (1.1) is defined as

$$\mathcal{L}(\mathbf{x}, \mathbf{y}, \boldsymbol{\lambda}) = f(\mathbf{x}) + g(\mathbf{y}) - \boldsymbol{\lambda}^\top (A\mathbf{x} + B\mathbf{y} - \mathbf{b}). \quad (1.3)$$

Although the Augmented Lagrangian Method (ALM) can be applied to solve (1.1), it does not take full advantage of the separable structure of (1.1). As a splitting version of ALM, the standard Alternating Direction Method of Multipliers (ADMM, [11, 12]) exploits the separable structure of the objective function and performs the following iterations:

$$\begin{cases} \mathbf{x}^{k+1} \in \arg\min_{\mathbf{x} \in \mathcal{X}} \mathcal{L}_\beta(\mathbf{x}, \mathbf{y}^k, \boldsymbol{\lambda}^k), \\ \mathbf{y}^{k+1} \in \arg\min_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}_\beta(\mathbf{x}^{k+1}, \mathbf{y}, \boldsymbol{\lambda}^k), \\ \boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^k - s\beta (A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}), \end{cases}$$

where $s \in (0, \frac{1+\sqrt{5}}{2})$ is the stepsize for updating the dual variable $\boldsymbol{\lambda}$.

If the Peaceman-Rachford Splitting Method (PRSM, [27]) is applied to the dual of (1.1), then we obtain a variation of ADMM, whose iteration reads

$$\begin{cases} \mathbf{x}^{k+1} \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \mathcal{L}_\beta(\mathbf{x}, \mathbf{y}^k, \boldsymbol{\lambda}^k), \\ \boldsymbol{\lambda}^{k+\frac{1}{2}} = \boldsymbol{\lambda}^k - \beta (A\mathbf{x}^{k+1} + B\mathbf{y}^k - \mathbf{b}), \\ \mathbf{y}^{k+1} \in \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}_\beta(\mathbf{x}^{k+1}, \mathbf{y}, \boldsymbol{\lambda}^{k+\frac{1}{2}}), \\ \boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^{k+\frac{1}{2}} - \beta (A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}). \end{cases}$$

PRSM is also called Symmetric ADMM (S-ADMM) since the Lagrange multipliers are symmetrically updated twice in each loop. Note that both updates of dual variable in PRSM use the same constant stepsize 1. Recently, Luo-Yang [25] proposed a fast S-ADMM for solving (1.1) with only equality constraints, where the Nesterov's acceleration technique was applied for an additional update of $\boldsymbol{\lambda}^{k+1}$ and the \mathbf{y} variable was updated again by solving

$$\min_{\mathbf{y}} (\mathbf{y} - (\hat{\boldsymbol{\lambda}}^{k+1})^\top B \mathbf{y}) \quad \text{with} \quad \hat{\boldsymbol{\lambda}}^{k+1} = \boldsymbol{\lambda}^{k+1} + \frac{\theta^{k+1}(1-\theta^k)}{\theta^k} (\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k)$$

and $\theta^{k+1} = 2/(k+1)$. Motivated from the ideas of enlarging the dual stepsize in [18], Gu, et al. [14] proposed a symmetric proximal ADMM whose dual variable is updated twice with different stepsizes. Meanwhile, the following extension of S-ADMM was developed by He, et al. [19]:

$$\begin{cases} \mathbf{x}^{k+1} \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \mathcal{L}_\beta(\mathbf{x}, \mathbf{y}^k, \boldsymbol{\lambda}^k), \\ \boldsymbol{\lambda}^{k+\frac{1}{2}} = \boldsymbol{\lambda}^k - \tau\beta (A\mathbf{x}^{k+1} + B\mathbf{y}^k - \mathbf{b}), \\ \mathbf{y}^{k+1} \in \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}_\beta(\mathbf{x}^{k+1}, \mathbf{y}, \boldsymbol{\lambda}^{k+\frac{1}{2}}), \\ \boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^{k+\frac{1}{2}} - s\beta (A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}), \end{cases} \quad (1.4)$$

where, for the sake of convergence, the stepsize pair (τ, s) is required to belong to the following region:

$$\Delta_0 = \left\{ (\tau, s) \mid s \in (0, (1 + \sqrt{5})/2), \tau + s > 0, \tau \in (-1, 1), |\tau| < 1 + s - s^2 \right\}.$$

Bai et al. [1] further designed a Generalized Symmetric ADMM (GS-ADMM) for solving a multi-block separable convex optimization and enlarged the above region Δ_0 to Δ defined in (3.1). Numerical experiments show that symmetrically updating the dual variable in a more flexible way often improves the algorithm performance [1, 14, 19]. The sublinear convergence rate of GS-ADMM in the nonergodic sense and its linear convergence rate were shown in [2]. To our knowledge, Δ is currently the largest convergence region

of the dual stepsizes for symmetric ADMM-type algorithms and has been used in the logarithmic-quadratic proximal based ADMM for solving the 2-block problems [17, 29] and the grouped multi-block problems [4].

Another line of developing ADMM is to apply relaxation techniques such as using

$$\chi^{k+1} = \alpha A x^{k+1} + (1 - \alpha)(b - B y^k)$$

to replace $A x^{k+1}$ during the updates of y^{k+1} and λ^{k+1} , where $\alpha \in (0, 2)$ is a relaxation factor. This leads to the classical Generalized ADMM (G-ADMM, [8]):

$$\begin{cases} x^{k+1} \in \arg \min_{x \in \mathcal{X}} \mathcal{L}_\beta(x, y^k, \lambda^k), \\ y^{k+1} \in \arg \min_{y \in \mathcal{Y}} g(y) - (\lambda^k)^\top B y + \frac{\beta}{2} \|\chi^{k+1} + B y - b\|^2, \\ \lambda^{k+1} = \lambda^k - \beta (\chi^{k+1} + B y^{k+1} - b). \end{cases}$$

Clearly, G-ADMM with $\alpha = 1$ would reduce to the standard ADMM with unit dual step-size. ADMM using some additional proximal terms in their subproblems is also called G-ADMM. For instance, the x -subproblem in [9] was proposed as $\min_{x \in \mathcal{X}} \mathcal{L}_\beta(x, y^k, \lambda^k) + \frac{1}{2} \|x - x^k\|_G^2$, where G is a symmetric positive definite matrix. For more recent G-ADMMs using possibly indefinite proximal terms, one may refer to the references [21, 30].

For convergence rate of ADMM, it is well-known that most of deterministic ADMM algorithms [1, 6, 7, 9, 15, 16, 18, 19, 21, 28, 31, 32] enjoy a global $\mathcal{O}(1/T)$ ergodic convergence rate for separable convex optimization, where T is the iteration number. Under the assumption that the subdifferential of each component objective function is piecewise linear, Yang-Han [34] established linear convergence rate of ADMM for two-block separable convex optimization. Assuming that an error bound condition holds, the dual stepsize is sufficiently small and the coefficient matrices in the equality constraint have full column ranks, Hong-Luo [20] showed a linear convergence rate of their multi-block ADMM. Zhang et al. [36] developed a majorized ADMM with indefinite proximal terms (iPADMM) for a class of composite convex optimization problems, and the authors analyzed the convergence of this iPADMM with a linear convergence rate under a local error bound condition. Moreover, Chang et al. [6] proposed a linearized symmetric ADMM with indefinite proximal regularization and optimal proximal parameter for solving the multi-block separable convex optimization. More recently, Yuan-Zeng-Zhang [33] showed that the local linear convergence of ADMM can be guaranteed by a partial error bound condition. For more details about linear convergence rate under strongly convexity assumption, we refer interested readers to [3, 5, 13, 23, 24] and the references therein.

2 Preliminaries

2.1 Notations and assumptions

Let \mathbb{R} , \mathbb{R}^n , and $\mathbb{R}^{n \times l}$ be the sets of real numbers, n dimensional real column vectors, and $n \times l$ dimensional real matrices, respectively. The I and $\mathbf{0}$ denote the identity matrix and the zero matrix/vector, respectively. For any symmetric matrices A and B of the same dimension, $A \succ B$ ($A \succeq B$) means $A - B$ is a positive definite (semidefinite) matrix. For any symmetric matrix G , define $\|x\|_G^2 := x^\top G x$ and $\|x\|_G := \sqrt{x^\top G x}$ if $G \succeq \mathbf{0}$. We use $\|\cdot\|$ to denote the standard Euclidean norm equipped with inner product $\langle \cdot, \cdot \rangle$, $\nabla f(x)$ to represent the gradient of f at x , and $\mathbb{E}[\cdot]$ to denote mathematical expectation of a random variable. We also define

$$w = \begin{pmatrix} x \\ y \\ \lambda \end{pmatrix}, \quad \mathcal{J}(w) = \begin{pmatrix} -A^\top \lambda \\ -B^\top \lambda \\ Ax + By - b \end{pmatrix}, \quad (2.1)$$

and

$$w^k = \begin{pmatrix} x^k \\ y^k \\ \lambda^k \end{pmatrix}, \quad \mathcal{J}(w^k) = \begin{pmatrix} -A^\top \lambda^k \\ -B^\top \lambda^k \\ Ax^k + By^k - b \end{pmatrix}. \quad (2.2)$$

For convenience of analysis, we simply denote $F(w) = f(x) + g(y)$.

We make the following two assumptions:

Assumption 2.1. *The primal-dual solution set Ω^* of problem (1.1) is nonempty, and the problem $\min_{y \in \mathcal{Y}} \{g(y) + \frac{1}{2} y^\top B^\top B y + z^\top y\}$ has a minimizer for any $z \in \mathbb{R}^{n_2}$.*

Assumption 2.2. *For any $\mathcal{H} \succ \mathbf{0}$, there exists a constant $\nu > 0$ such that the gradients ∇f_j satisfy the Lipschitz condition*

$$\|\nabla f_j(x_1) - \nabla f_j(x_2)\|_{\mathcal{H}^{-1}} \leq \nu \|x_1 - x_2\|_{\mathcal{H}} \quad (2.3)$$

for every $x_1, x_2 \in \mathcal{X}$ and $j = 1, 2, \dots, N$.

The first assumption is a basic assumption to ensure the solvability of the problem. Under Assumption 2.2, it holds that for every $x, y \in \mathcal{X}$, we have

$$f(x_1) \leq f(x_2) + \langle \nabla f(x_2), x_1 - x_2 \rangle + \frac{\nu}{2} \|x_1 - x_2\|_{\mathcal{H}}^2.$$

2.2 Variational characterization of (1.1)

Denote $\Omega = \mathcal{X} \times \mathcal{Y} \times \mathbb{R}^n$. It's well-known in convex optimization that any saddle-point of the Lagrangian (1.3) corresponds to a primal-dual solution of problem (1.1). A point $w^* := (x^*; y^*; \lambda^*) \in \Omega$ is called a saddle-point of $\mathcal{L}(x, y, \lambda)$ if it satisfies

$$\mathcal{L}(x^*, y^*, \lambda) \leq \mathcal{L}(x^*, y^*, \lambda^*) \leq \mathcal{L}(x, y, \lambda^*), \quad \forall w \in \Omega, \quad (2.4)$$

which is equivalent to

$$\begin{cases} f(\mathbf{x}) - f(\mathbf{x}^*) + (\mathbf{x} - \mathbf{x}^*)^\top (-A^\top \boldsymbol{\lambda}^*) \geq 0, \\ g(\mathbf{y}) - g(\mathbf{y}^*) + (\mathbf{y} - \mathbf{y}^*)^\top (-B^\top \boldsymbol{\lambda}^*) \geq 0, \\ A\mathbf{x}^* + B\mathbf{y}^* - \mathbf{b} = \mathbf{0}. \end{cases}$$

Rewriting these inequalities as a more compact form, it gives

$$F(\mathbf{w}) - F(\mathbf{w}^*) + (\mathbf{w} - \mathbf{w}^*)^\top \mathcal{J}(\mathbf{w}^*) \geq 0, \quad \forall \mathbf{w} \in \Omega. \quad (2.5)$$

Notice that the affine mapping $\mathcal{J}(\cdot)$ is skew-symmetric. So, we have

$$(\mathbf{w} - \bar{\mathbf{w}})^\top [\mathcal{J}(\mathbf{w}) - \mathcal{J}(\bar{\mathbf{w}})] \equiv 0, \quad \forall \mathbf{w}, \bar{\mathbf{w}} \in \Omega. \quad (2.6)$$

Hence, (2.5) is also equivalent to

$$F(\mathbf{w}) - F(\mathbf{w}^*) + (\mathbf{w} - \mathbf{w}^*)^\top \mathcal{J}(\mathbf{w}) \geq 0, \quad \forall \mathbf{w} \in \Omega. \quad (2.7)$$

The above discussion shows that the saddle-point \mathbf{w}^* can be also characterized by the variational inequality (2.7).

3 The proposed algorithm

Motivated from the stochastic AS-ADMM developed in [3] and the deterministic GS-ADMM proposed in [1], we now propose a Symmetric Accelerated Stochastic ADMM (SAS-ADMM, i.e., Algorithm 1), which has the similarly dual stepsize region Δ to that of GS-ADMM defined as

$$\Delta = \{(\tau, s) \mid \tau + s > 0, \tau \leq 1, -\tau^2 - s^2 - \tau s + \tau + s + 1 \geq 0\}. \quad (3.1)$$

The main features of SAS-ADMM are summarized as follows:

- (i) SAS-ADMM has many analogous advantages to AS-ADMM developed in [3]. Specifically, SAS-ADMM has low memory requirement since there is no need to save previous stochastic gradients and iterates. The subroutine **xsub** is a variant of deterministic accelerated gradient method where the full gradient is replaced by a stochastic gradient. In addition, users have the flexibility of choosing a zero mean random vector \mathbf{e}_t to reduce the variance of $\hat{\mathbf{g}}_t$. A simple choice is $\mathbf{e}_t = \mathbf{0}$, while faster convergence is observed in the numerical experiments when a variance reduction technique is employed (see (5.3) in Section 5). Under our blanket assumption that g is a proper convex function, the proximal \mathbf{y} -subproblem is solvable. Also, under the assumption that the projection onto the constraint set \mathcal{X} is simple, the iterations in subroutine **xsub** can be performed efficiently when both \mathcal{M}_k and \mathcal{H} are multiples of identity matrix[†].

[†]As explained in Remark 4.1 and experiments, both \mathcal{M}_k and \mathcal{H} could be chosen as multiples of identity matrix. So the $\tilde{\mathbf{x}}_{t+1}$ -subproblem is equivalent to a projection onto \mathcal{X} .

Algorithm 1 Symmetric accelerated stochastic ADMM (SAS-ADMM)

Parameters: $\beta > 0$, $\mathcal{H} \succ \mathbf{0}$, $L \succeq \mathbf{0}$ and $(\tau, s) \in \Delta$ given by (3.1).

Initialization: $(\mathbf{x}^0, \mathbf{y}^0, \boldsymbol{\lambda}^0) \in \mathcal{X} \times \mathcal{Y} \times \mathbb{R}^n$, $\check{\mathbf{x}}^0 = \mathbf{x}^0$.

For $k=0, 1, \dots$

Choose $m_k > 0$, $\eta_k > 0$ and \mathcal{M}_k such that $\mathcal{M}_k - \beta A^\top A \succeq \mathbf{0}$.

$\mathbf{h}^k := -A^\top \left[\boldsymbol{\lambda}^k - \beta(A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b}) \right]$.

$(\mathbf{x}^{k+1}, \check{\mathbf{x}}^{k+1}) = \mathbf{xsub}(\mathbf{x}^k, \check{\mathbf{x}}^k, \mathbf{h}^k)$.

$\boldsymbol{\lambda}^{k+\frac{1}{2}} = \boldsymbol{\lambda}^k - \tau\beta(A\mathbf{x}^{k+1} + B\mathbf{y}^k - \mathbf{b})$.

$\mathbf{y}^{k+1} \in \arg\min_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}_\beta(\mathbf{x}^{k+1}, \mathbf{y}, \boldsymbol{\lambda}^{k+\frac{1}{2}}) + \frac{1}{2} \|\mathbf{y} - \mathbf{y}^k\|_L^2$.

$\boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^{k+\frac{1}{2}} - s\beta(A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b})$.

end

$(\mathbf{x}^+, \check{\mathbf{x}}^+) = \mathbf{xsub}(\mathbf{x}_1, \check{\mathbf{x}}_1, \mathbf{h})$.

For $t=1, 2, \dots, m_k$

Randomly select $\zeta_t \in \{1, 2, \dots, N\}$ with uniform probability.

$\beta_t = 2/(t+1)$, $\gamma_t = 2/(t\eta_k)$, $\hat{\mathbf{x}}_t = \beta_t \check{\mathbf{x}}_t + (1-\beta_t)\mathbf{x}_t$.

$\mathbf{d}_t = \hat{\mathbf{g}}_t + \mathbf{e}_t$, where $\hat{\mathbf{g}}_t = \nabla f_{\zeta_t}(\hat{\mathbf{x}}_t)$ and \mathbf{e}_t is a random vector satisfying $\mathbb{E}[\mathbf{e}_t] = \mathbf{0}$.

$\check{\mathbf{x}}_{t+1} = \arg\min \left\{ \langle \mathbf{d}_t + \mathbf{h}, \mathbf{x} \rangle + \frac{\gamma_t}{2} \|\mathbf{x} - \check{\mathbf{x}}_t\|_{\mathcal{H}}^2 + \frac{1}{2} \|\mathbf{x} - \mathbf{x}^k\|_{\mathcal{M}_k}^2 : \mathbf{x} \in \mathcal{X} \right\}$.

$\mathbf{x}_{t+1} = \beta_t \check{\mathbf{x}}_{t+1} + (1-\beta_t)\mathbf{x}_t$.

end

Return $(\mathbf{x}^+, \check{\mathbf{x}}^+) = (\mathbf{x}_{m_k+1}, \check{\mathbf{x}}_{m_k+1})$.

- (ii) Unlike the classical ADMM and AS-ADMM [3], the dual variable of SAS-ADMM is symmetrically updated twice and allowed to use the large stepsize region (3.1). SAS-ADMM will reduce to the aforementioned PRSM if the \mathbf{x} -subproblem is solved deterministically as in S-ADMM, $L = \mathbf{0}$ and $(\tau, s) = (1, 1) \in \Delta$. When the stepsize $\tau = 0$ and $L = \mathbf{0}$, SAS-ADMM reduces to AS-ADMM with stepsize $s \in (0, \frac{1+\sqrt{5}}{2}]$ (the half-open interval). Compared with the standard stepsize region $(0, \frac{1+\sqrt{5}}{2})$ for the dual variable of ADMM, this symmetric updates of dual variable is more balanced, flexible and often lead better numerical performance.
- (iii) If $m_k=1, N=1$, then SAS-ADMM degrades to a linearized symmetric ADMM. When $m_k > 1, N=1$, SAS-ADMM is a multi-step deterministic inexact symmetric ADMM. Hence, the convergence properties developed in this paper also apply to these deterministic algorithms as special cases. Moreover, by taking $L = \gamma I - \beta B^\top B$ for some

$\gamma > 0$, the \mathbf{y} -subproblem would become the following proximal mapping problem:

$$\text{prox}_{\gamma}^g(\mathbf{y}_c^k) := \arg\min_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{y}) + \frac{\gamma}{2} \|\mathbf{y} - \mathbf{y}_c^k\|^2, \quad (3.2)$$

where $\mathbf{y}_c^k = \mathbf{y}^k - \beta B^\top (A\mathbf{x}^{k+1} + B\mathbf{y}^k - \mathbf{b} - \boldsymbol{\lambda}^{k+\frac{1}{2}}/\beta)/\gamma$. In this case, the Assumption 2.1 is not required since strong convexity of the \mathbf{y} -subproblem implies a unique global solution and a closed-form solution may exist when g has certain structure.

- (iv) With the aid of variational analysis, we show that SAS-ADMM has the worst-case $\mathcal{O}(1/T)$ ergodic convergence rate in terms of the expectation of both the objective value gap and the constraint violation, where T is the number of the outer iterations. Preliminary experiments and results show that SAS-ADMM performs competitively well and often slightly better than AS-ADMM [3] for solving a family of separable convex optimization problems arising from big-data applications.

4 Convergence analysis

To establish the convergence of Algorithm 1, we first need the following lemma about the iterates generated by the `xsub` routine in Algorithm 1. The lemma was given in [3] and thus we omit its proof.

Lemma 4.1. [3, Lemma 3.2] *Let $\delta_t = \nabla f(\hat{\mathbf{x}}_t) - \mathbf{d}_t$. Suppose $\eta_k \in (0, 1/\nu)$ and Assumption 2.2 holds. Then, the iterates generated by Algorithm 1 satisfy*

$$f(\mathbf{x}) - f(\mathbf{x}^{k+1}) + \langle \mathbf{x} - \mathbf{x}^{k+1}, -A^\top \tilde{\boldsymbol{\lambda}}^k \rangle \geq \langle \mathbf{x}^{k+1} - \mathbf{x}, \mathcal{D}_k(\mathbf{x}^{k+1} - \mathbf{x}^k) \rangle + \zeta^k \quad (4.1)$$

for all $\mathbf{x} \in \mathcal{X}$, where

$$\tilde{\boldsymbol{\lambda}}^k = \boldsymbol{\lambda}^k - \beta (A\mathbf{x}^{k+1} + B\mathbf{y}^k - \mathbf{b}), \quad \mathcal{D}_k = \mathcal{M}_k - \beta A^\top A \quad (4.2)$$

and

$$\begin{aligned} \zeta^k = & \frac{2}{m_k(m_k+1)} \left[\frac{1}{\eta_k} \left(\|\mathbf{x} - \check{\mathbf{x}}^{k+1}\|_{\mathcal{H}}^2 - \|\mathbf{x} - \check{\mathbf{x}}^k\|_{\mathcal{H}}^2 \right) \right. \\ & \left. - \sum_{t=1}^{m_k} t \langle \delta_t, \check{\mathbf{x}}_t - \mathbf{x} \rangle - \frac{\eta_k}{4(1-\eta_k\nu)} \sum_{t=1}^{m_k} t^2 \|\delta_t\|_{\mathcal{H}^{-1}}^2 \right]. \end{aligned} \quad (4.3)$$

Based on the above lemma, we can immediately establish the following result.

Lemma 4.2. *Suppose $\eta_k \in (0, 1/\nu)$. Then, the iterates generated by Algorithm 1 satisfy*

$$F(\mathbf{w}) - F(\tilde{\mathbf{w}}^k) + \langle \mathbf{w} - \tilde{\mathbf{w}}^k, \mathcal{J}(\mathbf{w}) \rangle \geq (\mathbf{w} - \tilde{\mathbf{w}}^k)^\top Q_k(\mathbf{w}^k - \tilde{\mathbf{w}}^k) + \zeta^k \quad (4.4)$$

for all $\mathbf{w} \in \Omega$, where ζ^k and $\tilde{\lambda}^k$ are defined in (4.3) and (4.2),

$$\tilde{\mathbf{w}}^k = \begin{pmatrix} \tilde{\mathbf{x}}^k \\ \tilde{\mathbf{y}}^k \\ \tilde{\lambda}^k \end{pmatrix} = \begin{pmatrix} \mathbf{x}^{k+1} \\ \mathbf{y}^{k+1} \\ \tilde{\lambda}^k \end{pmatrix} \quad \text{and} \quad Q_k = \begin{bmatrix} \mathcal{D}_k & & \\ & L + \beta B^\top B & -\tau B^\top \\ & -B & \frac{1}{\beta} \mathbf{I} \end{bmatrix}. \quad (4.5)$$

Proof. By the first-order optimality condition of the \mathbf{y} -subproblem, we have

$$g(\mathbf{y}) - g(\mathbf{y}^{k+1}) + \langle \mathbf{y} - \mathbf{y}^{k+1}, \mathbf{p}_k \rangle \geq 0, \quad \forall \mathbf{y} \in \mathcal{Y}, \quad (4.6)$$

where \mathbf{p}_k is the gradient of the smooth terms in the objective function of the \mathbf{y} -subproblem:

$$\begin{aligned} \mathbf{p}_k &= -B^\top \lambda^{k+\frac{1}{2}} + \beta B^\top (A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}) + L(\mathbf{y}^{k+1} - \mathbf{y}^k) \\ &= -B^\top \lambda^{k+\frac{1}{2}} + \beta B^\top (A\mathbf{x}^{k+1} + B\mathbf{y}^k - \mathbf{b}) + [L + \beta B^\top B] (\mathbf{y}^{k+1} - \mathbf{y}^k) \\ &= -B^\top \lambda^{k+\frac{1}{2}} + B^\top (\lambda^k - \tilde{\lambda}^k) + [L + \beta B^\top B] (\mathbf{y}^{k+1} - \mathbf{y}^k) \\ &= -B^\top \tilde{\lambda}^k + \tau B^\top (\lambda^k - \tilde{\lambda}^k) + [L + \beta B^\top B] (\mathbf{y}^{k+1} - \mathbf{y}^k). \end{aligned}$$

The above last equality uses the following relation

$$\lambda^{k+\frac{1}{2}} = \lambda^k - \tau(\lambda^k - \tilde{\lambda}^k). \quad (4.7)$$

By the definition of $\tilde{\lambda}^k$, we have

$$(A\tilde{\mathbf{x}}^k + B\tilde{\mathbf{y}}^k - \mathbf{b}) - B(\tilde{\mathbf{y}}^k - \mathbf{y}^k) + \frac{1}{\beta}(\tilde{\lambda}^k - \lambda^k) = \mathbf{0}. \quad (4.8)$$

Taking inner product of the above equality with $\lambda - \tilde{\lambda}^k$, we get

$$\langle \lambda - \tilde{\lambda}^k, A\tilde{\mathbf{x}}^k + B\tilde{\mathbf{y}}^k - \mathbf{b} \rangle = \left\langle \lambda - \tilde{\lambda}^k, -B(\mathbf{y}^k - \tilde{\mathbf{y}}^k) + \frac{1}{\beta}(\lambda^k - \tilde{\lambda}^k) \right\rangle. \quad (4.9)$$

Then, the inequality (4.4) is achieved by combining (4.1), (4.6), (4.9) together with the property in (2.6). \square

4.1 More technical results

We show the following corollary and lemmas for establishing the main convergence theorem of our Algorithm 1.

Corollary 4.1. Suppose $\eta_k \in (0, 1/\nu)$. Then, the iterates generated by Algorithm 1 satisfy

$$\begin{aligned} & F(\mathbf{w}) - F(\tilde{\mathbf{w}}^k) + (\mathbf{w} - \tilde{\mathbf{w}}^k)^\top \mathcal{J}(\mathbf{w}) \\ & \geq \frac{1}{2} \left\{ \|\mathbf{w} - \mathbf{w}^{k+1}\|_{\tilde{Q}_k}^2 - \|\mathbf{w} - \mathbf{w}^k\|_{\tilde{Q}_k}^2 + \|\mathbf{w}^k - \tilde{\mathbf{w}}^k\|_{\tilde{G}_k}^2 \right\} + \zeta^k \end{aligned} \quad (4.10)$$

for all $\mathbf{w} \in \Omega$, where ζ^k is defined in (4.3) and

$$\begin{aligned} \tilde{Q}_k &= \begin{bmatrix} \mathcal{D}_k & & \\ & L + (1 - \frac{\tau s}{\tau+s})\beta B^\top B & -\frac{\tau}{\tau+s}B^\top \\ & -\frac{\tau}{\tau+s}B & \frac{1}{\beta(\tau+s)}\mathbf{I} \end{bmatrix}, \\ \tilde{G}_k &= \begin{bmatrix} \mathcal{D}_k & & \\ & L + (1-s)\beta B^\top B & (s-1)B^\top \\ & (s-1)B & \frac{2-\tau-s}{\beta}\mathbf{I} \end{bmatrix}. \end{aligned} \quad (4.11)$$

Proof. By (4.7) and the way of generating λ^{k+1} , we have

$$-s\beta B(\mathbf{y}^k - \tilde{\mathbf{y}}^k) + (\tau+s)(\lambda^k - \tilde{\lambda}^k) = \lambda^k - \lambda^{k+1}, \quad (4.12)$$

which, by the definition of $\tilde{\mathbf{w}}^k$ in (4.5), further shows

$$\mathbf{w}^k - \mathbf{w}^{k+1} = P(\mathbf{w}^k - \tilde{\mathbf{w}}^k) \quad \text{with} \quad P = \begin{bmatrix} \mathbf{I} & & \\ & \mathbf{I} & \\ & -s\beta B & (\tau+s)\mathbf{I} \end{bmatrix}. \quad (4.13)$$

Hence, the relation $Q_k(\mathbf{w}^k - \tilde{\mathbf{w}}^k) = Q_k P^{-1}(\mathbf{w}^k - \mathbf{w}^{k+1})$ holds and

$$Q_k P^{-1} = \begin{bmatrix} \mathcal{D}_k & & \\ & L + (1 - \frac{\tau s}{\tau+s})\beta B^\top B & -\frac{\tau}{\tau+s}B^\top \\ & -\frac{\tau}{\tau+s}B & \frac{1}{\beta(\tau+s)}\mathbf{I} \end{bmatrix} = \tilde{Q}_k.$$

For any $\mathbf{w} \in \Omega$, it follows from (4.4) and the above relation that

$$\begin{aligned} & F(\mathbf{w}) - F(\tilde{\mathbf{w}}^k) + (\mathbf{w} - \tilde{\mathbf{w}}^k)^\top \mathcal{J}(\mathbf{w}) \\ & \geq \zeta^k + (\mathbf{w} - \tilde{\mathbf{w}}^k)^\top \tilde{Q}_k(\mathbf{w}^k - \mathbf{w}^{k+1}) \\ & = \zeta^k + \frac{1}{2} \left\{ \|\mathbf{w} - \mathbf{w}^{k+1}\|_{\tilde{Q}_k}^2 - \|\mathbf{w} - \mathbf{w}^k\|_{\tilde{Q}_k}^2 + \|\mathbf{w}^k - \tilde{\mathbf{w}}^k\|_{\tilde{Q}_k}^2 - \|\mathbf{w}^{k+1} - \tilde{\mathbf{w}}^k\|_{\tilde{Q}_k}^2 \right\}, \end{aligned} \quad (4.14)$$

where the equality uses the identity

$$2(a-b)^\top \tilde{Q}_k(c-d) = \|a-d\|_{\tilde{Q}_k}^2 - \|a-c\|_{\tilde{Q}_k}^2 + \|c-b\|_{\tilde{Q}_k}^2 - \|b-d\|_{\tilde{Q}_k}^2$$

with specifications $a := w$, $b := \tilde{w}^k$, $c := w^k$, $d := w^{k+1}$.

Now, by (4.13) again, we deduce

$$\begin{aligned} \left\| w^k - \tilde{w}^k \right\|_{\tilde{Q}_k}^2 - \left\| w^{k+1} - \tilde{w}^k \right\|_{\tilde{Q}_k}^2 &= \left\| w^k - \tilde{w}^k \right\|_{\tilde{Q}_k}^2 - \left\| w^{k+1} - w^k + w^k - \tilde{w}^k \right\|_{\tilde{Q}_k}^2 \\ &= \left\| w^k - \tilde{w}^k \right\|_{\tilde{Q}_k}^2 - \left\| w^k - \tilde{w}^k - P(w^k - \tilde{w}^k) \right\|_{\tilde{Q}_k}^2 \\ &= \left\| w^k - \tilde{w}^k \right\|_{\tilde{G}_k}^2, \end{aligned}$$

where we use the relation $\tilde{G}_k = P^\top \tilde{Q}_k + \tilde{Q}_k P - P^\top \tilde{Q}_k P$ to obtain the last equality. Then, (4.10) follows from (4.14). \square

In the above Corollary 4.1 and its proof, since \tilde{Q}_k is not necessarily positive semidefinite for any parameter τ , we abuse the notation $\|w^k\|_{\tilde{Q}_k}^2 := (w^k)^\top \tilde{Q}_k w^k$. Next, we provide a sufficient condition to ensure the positive semidefiniteness of \tilde{Q}_k .

Lemma 4.3. *Let $L \succeq (\tau-1)\beta B^\top B$. Then, the matrix \tilde{Q}_k given by (4.11) is symmetric positive semidefinite for any $(\tau, s) \in \Delta$.*

Proof. Clearly, we just need to check the lower-upper 2-by-2 block of \tilde{Q}_k , i.e.,

$$\begin{aligned} \tilde{Q}_k^L &= \begin{bmatrix} L + (1 - \frac{\tau s}{\tau+s})\beta B^\top B & -\frac{\tau}{\tau+s}B^\top \\ -\frac{\tau}{\tau+s}B & \frac{1}{\beta(\tau+s)}I \end{bmatrix} \\ &\succeq \begin{bmatrix} (\tau - \frac{\tau s}{\tau+s})\beta B^\top B & -\frac{\tau}{\tau+s}B^\top \\ -\frac{\tau}{\tau+s}B & \frac{1}{\beta(\tau+s)}I \end{bmatrix} \\ &= \begin{bmatrix} \beta^{\frac{1}{2}}B & \\ & \beta^{-\frac{1}{2}}I \end{bmatrix}^\top \begin{bmatrix} \frac{\tau^2}{\tau+s}I & -\frac{\tau}{\tau+s}I \\ -\frac{\tau}{\tau+s}I & \frac{1}{\tau+s}I \end{bmatrix} \begin{bmatrix} \beta^{\frac{1}{2}}B & \\ & \beta^{-\frac{1}{2}}I \end{bmatrix} \end{aligned}$$

is positive semidefinite. Notice that

$$\begin{bmatrix} I & \\ \tau I & I \end{bmatrix}^\top \begin{bmatrix} \frac{\tau^2}{\tau+s}I & -\frac{\tau}{\tau+s}I \\ -\frac{\tau}{\tau+s}I & \frac{1}{\tau+s}I \end{bmatrix} \begin{bmatrix} I & \\ \tau I & I \end{bmatrix} = \begin{bmatrix} 0 & \\ & \frac{1}{\tau+s}I \end{bmatrix}.$$

So, \tilde{Q}_k^L is positive semidefinite since $\tau+s > 0$ for any $(\tau, s) \in \Delta$. \square

To show the global convergence of Algorithm 1, we need to further establish a useful lower bound on the term $\|w^k - \tilde{w}^k\|_{\tilde{G}_k}^2$, since \tilde{G}_k is not necessarily positive definite for any $(\tau, s) \in \Delta$. In the following lemma, we assume $L \succeq 0$, which implies $L \succeq (\tau-1)\beta B^\top B$ since $\tau \leq 1$, and as a consequence, Lemma 4.3 holds.

Lemma 4.4. Let $L \succeq \mathbf{0}$. Then, for any $(\tau, s) \in \Delta$ defined in (3.1), we have

$$\begin{aligned} \|\mathbf{w}^k - \tilde{\mathbf{w}}^k\|_{\tilde{G}_k}^2 &\geq \|\mathbf{x}^k - \mathbf{x}^{k+1}\|_{\mathcal{D}_k}^2 + \omega_0 \|A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}\|^2 \\ &\quad + \omega_1 \left(\|A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}\|^2 - \|A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b}\|^2 \right) \\ &\quad + \omega_2 \left(\|\mathbf{y}^k - \mathbf{y}^{k+1}\|_L^2 - \|\mathbf{y}^{k-1} - \mathbf{y}^k\|_L^2 \right), \end{aligned} \quad (4.15)$$

where $\omega_i \geq 0$, $i=0,1,2$, are given as

$$\omega_0 = \left(2 - \tau - s - \frac{(1-s)^2}{1+\tau} \right) \beta, \quad \omega_1 = \frac{(1-s)^2}{1+\tau} \beta \quad \text{and} \quad \omega_2 = \frac{1-\tau}{1+\tau}. \quad (4.16)$$

Proof. For any $(\tau, s) \in \Omega$, we have $1+\tau > 0$. By the structure of \tilde{G}_k in (4.11), $L \succeq \mathbf{0}$ and (4.8), we have

$$\begin{aligned} &\|\mathbf{w}^k - \tilde{\mathbf{w}}^k\|_{\tilde{G}_k}^2 \\ &= \|\mathbf{x}^k - \mathbf{x}^{k+1}\|_{\mathcal{D}_k}^2 + \|\mathbf{y}^k - \mathbf{y}^{k+1}\|_{L + (1-s)\beta B^\top B}^2 \\ &\quad + 2(s-1) \left(\boldsymbol{\lambda}^k - \tilde{\boldsymbol{\lambda}}^k \right)^\top B \left(\mathbf{y}^k - \mathbf{y}^{k+1} \right) + \frac{2-\tau-s}{\beta} \|\boldsymbol{\lambda}^k - \tilde{\boldsymbol{\lambda}}^k\|^2 \\ &\geq \|\mathbf{x}^k - \mathbf{x}^{k+1}\|_{\mathcal{D}_k}^2 + (2-\tau-s)\beta \|A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}\|^2 \\ &\quad + (1-\tau)\beta \|B(\mathbf{y}^k - \mathbf{y}^{k+1})\|^2 + 2(1-\tau)\beta \left(A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b} \right)^\top B(\mathbf{y}^k - \mathbf{y}^{k+1}). \end{aligned} \quad (4.17)$$

We now estimate the last crossing term in (4.17). Taking $\mathbf{y} = \mathbf{y}^k$ in the first-order optimality condition (4.6) yields

$$g(\mathbf{y}^k) - g(\mathbf{y}^{k+1}) + \left\langle \mathbf{y}^k - \mathbf{y}^{k+1}, -B^\top \boldsymbol{\lambda}^{k+\frac{1}{2}} + \beta B^\top (A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}) + L(\mathbf{y}^{k+1} - \mathbf{y}^k) \right\rangle \geq 0.$$

Similarly, letting $\mathbf{y} = \mathbf{y}^{k+1}$ in the first-order optimality condition of the \mathbf{y} -subproblem at the $(k-1)$ -th iteration gives

$$g(\mathbf{y}^{k+1}) - g(\mathbf{y}^k) + \left\langle \mathbf{y}^{k+1} - \mathbf{y}^k, -B^\top \boldsymbol{\lambda}^{k-\frac{1}{2}} + \beta B^\top (A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b}) + L(\mathbf{y}^k - \mathbf{y}^{k-1}) \right\rangle \geq 0.$$

Summing up the above two inequalities together with the relation

$$\boldsymbol{\lambda}^{k-\frac{1}{2}} - \boldsymbol{\lambda}^{k+\frac{1}{2}} = \tau\beta (A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}) + s\beta (A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b}) + \tau\beta B(\mathbf{y}^k - \mathbf{y}^{k+1}),$$

and noticing that $1+\tau>0$, we obtain

$$\begin{aligned}
& \left(A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b} \right)^\top B \left(\mathbf{y}^k - \mathbf{y}^{k+1} \right) \\
& \geq \frac{1-s}{1+\tau} \left(A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b} \right)^\top B \left(\mathbf{y}^k - \mathbf{y}^{k+1} \right) - \frac{\tau}{1+\tau} \left\| B(\mathbf{y}^k - \mathbf{y}^{k+1}) \right\|^2 \\
& \quad + \frac{1}{\beta(1+\tau)} \left(\mathbf{y}^k - \mathbf{y}^{k+1} \right)^\top L \left[\left(\mathbf{y}^k - \mathbf{y}^{k+1} \right) - \left(\mathbf{y}^{k-1} - \mathbf{y}^k \right) \right] \\
& \geq \frac{1-s}{1+\tau} \left(A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b} \right)^\top B \left(\mathbf{y}^k - \mathbf{y}^{k+1} \right) - \frac{\tau}{1+\tau} \left\| B(\mathbf{y}^k - \mathbf{y}^{k+1}) \right\|^2 \\
& \quad + \frac{1}{2\beta(1+\tau)} \left(\left\| \mathbf{y}^k - \mathbf{y}^{k+1} \right\|_L^2 - \left\| \mathbf{y}^{k-1} - \mathbf{y}^k \right\|_L^2 \right). \tag{4.18}
\end{aligned}$$

Then, combining (4.17) and (4.18) we have

$$\begin{aligned}
\left\| \mathbf{w}^k - \tilde{\mathbf{w}}^k \right\|_{\tilde{G}_k}^2 & \geq \left\| \mathbf{x}^k - \mathbf{x}^{k+1} \right\|_{\mathcal{D}_k}^2 + (2-\tau-s)\beta \left\| A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b} \right\|^2 \\
& \quad + \frac{2\beta}{1+\tau} (1-\tau)(1-s) \left(A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b} \right)^\top B \left(\mathbf{y}^k - \mathbf{y}^{k+1} \right) \\
& \quad + (1-\tau)\beta \left\| B \left(\mathbf{y}^k - \mathbf{y}^{k+1} \right) \right\|^2 - \frac{2\tau(1-\tau)}{1+\tau} \beta \left\| B \left(\mathbf{y}^k - \mathbf{y}^{k+1} \right) \right\|^2 \\
& \quad + \frac{1-\tau}{1+\tau} \left(\left\| \mathbf{y}^k - \mathbf{y}^{k+1} \right\|_L^2 - \left\| \mathbf{y}^{k-1} - \mathbf{y}^k \right\|_L^2 \right) \\
& \geq \left\| \mathbf{x}^k - \mathbf{x}^{k+1} \right\|_{\mathcal{D}_k}^2 + \left(2-\tau-s - \frac{(1-s)^2}{1+\tau} \right) \beta \left\| A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b} \right\|^2 \\
& \quad + \frac{(1-s)^2}{1+\tau} \beta \left(\left\| A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b} \right\|^2 - \left\| A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b} \right\|^2 \right) \\
& \quad + \frac{1-\tau}{1+\tau} \left(\left\| \mathbf{y}^k - \mathbf{y}^{k+1} \right\|_L^2 - \left\| \mathbf{y}^{k-1} - \mathbf{y}^k \right\|_L^2 \right) \\
& \quad + \left(1-\tau - \frac{(1-\tau)^2}{1+\tau} - \frac{2\tau(1-\tau)}{1+\tau} \right) \beta \left\| B(\mathbf{y}^k - \mathbf{y}^{k+1}) \right\|^2 \\
& = \left\| \mathbf{x}^k - \mathbf{x}^{k+1} \right\|_{\mathcal{D}_k}^2 + \left(2-\tau-s - \frac{(1-s)^2}{1+\tau} \right) \beta \left\| A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b} \right\|^2 \\
& \quad + \frac{(1-s)^2}{1+\tau} \beta \left(\left\| A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b} \right\|^2 - \left\| A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b} \right\|^2 \right) \\
& \quad + \frac{1-\tau}{1+\tau} \left(\left\| \mathbf{y}^k - \mathbf{y}^{k+1} \right\|_L^2 - \left\| \mathbf{y}^{k-1} - \mathbf{y}^k \right\|_L^2 \right),
\end{aligned}$$

where the second inequality follows from the Cauchy-Schwartz inequality

$$\begin{aligned} & 2(1-s)(1-\tau) \left(A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b} \right)^\top B \left(\mathbf{y}^k - \mathbf{y}^{k+1} \right) \\ & \geq -(1-s)^2 \left\| A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b} \right\|^2 - (1-\tau)^2 \left\| B \left(\mathbf{y}^k - \mathbf{y}^{k+1} \right) \right\|^2. \end{aligned}$$

So, (4.15) holds with $\omega_i, i=0,1,2$, defined as in (4.16). Moreover, for any $(\tau, s) \in \Delta$, we can derive $\omega_i \geq 0$ for $i=0,1,2$. This completes the whole proof. \square

4.2 Iteration complexity in expectation

We now analyze the global ergodic convergence and the iteration complexity of Algorithm 1.

Theorem 4.1. Suppose $L \succeq \mathbf{0}$ and $(\tau, s) \in \Delta$ defined in (3.1). If for some integers $\kappa, T > 0$, the following conditions hold for all $k \in [\kappa, \kappa + T]$: (I) $\eta_k \in (0, 1/(2\nu)]$ and the sequence $\{\eta_k m_k(m_k + 1)\}$ is nondecreasing; (II) $\mathcal{D}_k \succeq \mathcal{D}_{k+1} \succeq \mathbf{0}$ and $\mathbb{E}(\|\delta_t\|_{\mathcal{H}^{-1}}^2) \leq \sigma^2$ for some $\sigma > 0$, where δ_t and \mathcal{D}_k are defined in Lemma 4.1. Then, for any $\mathbf{w} \in \Omega$, we have

$$\begin{aligned} & \mathbb{E} \left[F(\mathbf{w}_T) - F(\mathbf{w}) + (\mathbf{w}_T - \mathbf{w})^\top \mathcal{J}(\mathbf{w}) \right] \\ & \leq \frac{1}{2T} \left\{ \sigma^2 \sum_{k=\kappa}^{\kappa+T} \eta_k m_k + \frac{4}{m_\kappa(m_\kappa + 1)\eta_\kappa} \|\mathbf{x} - \check{\mathbf{x}}^\kappa\|_{\mathcal{H}}^2 + \|\mathbf{w} - \mathbf{w}^\kappa\|_{\tilde{\mathcal{Q}}_\kappa}^2 \right. \\ & \quad \left. + \omega_1 \|A\mathbf{x}^\kappa + B\mathbf{y}^\kappa - \mathbf{b}\|^2 + \omega_2 \|\mathbf{y}^{\kappa-1} - \mathbf{y}^\kappa\|_L^2 \right\}, \end{aligned} \quad (4.19)$$

where $\mathbf{w}_T = \frac{1}{T} \sum_{k=\kappa}^{\kappa+T} \tilde{\mathbf{w}}^k$, $\omega_1 \geq 0$ and $\omega_2 \geq 0$ are defined in (4.16).

Proof. By the assumption, $\mathcal{D}_k \succeq \mathcal{D}_{k+1} \succeq \mathbf{0}$ implies $\tilde{\mathcal{Q}}_k \succeq \tilde{\mathcal{Q}}_{k+1} \succeq \mathbf{0}$ for the matrix $\tilde{\mathcal{Q}}_k$ given in (4.11). Substituting (4.15) into (4.10) and utilizing the relation $\tilde{\mathcal{Q}}_k \succeq \tilde{\mathcal{Q}}_{k+1}$, it follows from Lemma 4.4 that

$$\begin{aligned} & F(\tilde{\mathbf{w}}^k) - F(\mathbf{w}) + (\tilde{\mathbf{w}}^k - \mathbf{w})^\top \mathcal{J}(\mathbf{w}) \\ & \leq -\zeta^k + \frac{1}{2} \left\{ \|\mathbf{w} - \mathbf{w}^k\|_{\tilde{\mathcal{Q}}_k}^2 - \|\mathbf{w} - \mathbf{w}^{k+1}\|_{\tilde{\mathcal{Q}}_{k+1}}^2 \right\} \\ & \quad + \frac{\omega_1}{2} \left(\|A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b}\|^2 - \|A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}\|^2 \right) \\ & \quad + \frac{\omega_2}{2} \left(\|\mathbf{y}^{k-1} - \mathbf{y}^k\|_L^2 - \|\mathbf{y}^k - \mathbf{y}^{k+1}\|_L^2 \right), \end{aligned}$$

where ω_1, ω_2 are defined in (4.16). Summing the above inequality over k between κ and

$\kappa + T$, we deduce by Lemma 4.3 that

$$\begin{aligned} & \sum_{k=\kappa}^{\kappa+T} F(\tilde{\mathbf{w}}^k) - T \left\{ F(\mathbf{w}) + (\mathbf{w}_T - \mathbf{w})^\top \mathcal{J}(\mathbf{w}) \right\} \\ & \leq - \sum_{k=\kappa}^{\kappa+T} \zeta^k + \frac{1}{2} \left\{ \|\mathbf{w} - \mathbf{w}^\kappa\|_{\tilde{Q}_\kappa}^2 + \omega_1 \|A\mathbf{x}^\kappa + B\mathbf{y}^\kappa - \mathbf{b}\|^2 + \omega_2 \|\mathbf{y}^{\kappa-1} - \mathbf{y}^\kappa\|_L^2 \right\}. \end{aligned} \quad (4.20)$$

Then, it follows from convexity of F and the definition of \mathbf{w}_T that

$$F(\mathbf{w}_T) \leq \frac{1}{T} \sum_{k=\kappa}^{\kappa+T} F(\tilde{\mathbf{w}}^k). \quad (4.21)$$

Dividing (4.20) by T and using (4.21), we obtain

$$\begin{aligned} & F(\mathbf{w}_T) - F(\mathbf{w}) + (\mathbf{w}_T - \mathbf{w})^\top \mathcal{J}(\mathbf{w}) \\ & \leq \frac{1}{T} \left[- \sum_{k=\kappa}^{\kappa+T} \zeta^k + \frac{1}{2} \left\{ \|\mathbf{w} - \mathbf{w}^\kappa\|_{\tilde{Q}_\kappa}^2 + \omega_1 \|A\mathbf{x}^\kappa + B\mathbf{y}^\kappa - \mathbf{b}\|^2 + \omega_2 \|\mathbf{y}^{\kappa-1} - \mathbf{y}^\kappa\|_L^2 \right\} \right]. \end{aligned} \quad (4.22)$$

Let us focus on the terms involving ζ^k . By assumption, the sequence $\{m_k(m_k+1)\eta_k\}$ is nondecreasing for $k \in [\kappa, \kappa+T]$ and $\mathcal{H} \succ \mathbf{0}$, thus we have

$$\sum_{k=\kappa}^{\kappa+T} \frac{2}{m_k(m_k+1)\eta_k} \left(\|\mathbf{x} - \check{\mathbf{x}}^k\|_{\mathcal{H}}^2 - \|\mathbf{x} - \check{\mathbf{x}}^{k+1}\|_{\mathcal{H}}^2 \right) \leq \frac{2\|\mathbf{x} - \check{\mathbf{x}}^\kappa\|_{\mathcal{H}}^2}{m_\kappa(m_\kappa+1)\eta_\kappa}. \quad (4.23)$$

Note that

$$\delta_t = \nabla f(\hat{\mathbf{x}}_t) - \mathbf{d}_t = \nabla f(\hat{\mathbf{x}}_t) - \nabla f_{\zeta_t}(\hat{\mathbf{x}}_t) - \mathbf{e}_t$$

only depends on the index ζ_t . So we have $\mathbb{E}[\delta_t] = \mathbf{0}$ since the random variable $\zeta_t \in \{1, 2, \dots, N\}$ is chosen with uniform probability and $\mathbb{E}[\mathbf{e}_t] = \mathbf{0}$. Also, since $\check{\mathbf{x}}_t$ depends on $\zeta_{t-1}, \zeta_{t-2}, \dots$, we have $\mathbb{E}[\langle \delta_t, \check{\mathbf{x}}_t - \mathbf{x} \rangle] = \mathbf{0}$. By the assumption that $\mathbb{E}(\|\delta_t\|_{\mathcal{H}^{-1}}^2) \leq \sigma^2$, we have

$$\mathbb{E} \left[\sum_{t=1}^{m_k} t^2 \|\delta_t\|_{\mathcal{H}^{-1}}^2 \right] \leq \frac{\sigma^2 m_k(m_k+1)(2m_k+1)}{6} \leq \frac{\sigma^2}{2} m_k^2(m_k+1)$$

since $m_k \geq 1$. Combining these bounds for the terms in ζ^k with the condition $\eta_k \leq 1/(2\nu)$ is to get

$$- \mathbb{E} \left[\sum_{k=\kappa}^{\kappa+T} \zeta^k \right] \leq \frac{2}{m_\kappa(m_\kappa+1)\eta_\kappa} \|\mathbf{x} - \check{\mathbf{x}}^\kappa\|_{\mathcal{H}}^2 + \frac{\sigma^2}{2} \sum_{k=\kappa}^{\kappa+T} \eta_k m_k.$$

Finally, applying the expectation operator to (4.22) and substituting this bound into the ζ^k term complete the proof. \square

By properly setting the algorithm parameters, the following theorem shows the convergence rate of Algorithm 1 in the expectation of both the objective function value gap and the constraint violation.

Theorem 4.2. *Suppose the conditions in Theorem 4.1 hold. Let*

$$\eta_k = \min \left\{ \frac{c_1}{m_k(m_k+1)}, c_2 \right\} \quad \text{and} \quad m_k = \max \{ \lceil c_3 k^\varrho \rceil, m \}, \quad (4.24)$$

where $c_1, c_2, c_3 > 0$, $\varrho \geq 1$ are constants and $m > 0$ is a given integer. Then, for every $\mathbf{w}^* \in \Omega^*$, we have

$$|\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)]| = E_\varrho(T) = \mathbb{E}[\|A\mathbf{x}_T + B\mathbf{y}_T - \mathbf{b}\|], \quad (4.25)$$

where $E_\varrho(T) = \mathcal{O}(1/T)$ for $\varrho > 1$ and $E_1(T) = \mathcal{O}(T^{-1} \log T)$ for $\varrho = 1$.

Proof. The proof is same as that of [3, Theorem 4.2] and thus is omitted here. \square

Remark 4.1. (I) In practice, the matrix \mathcal{M}_k in Algorithm 1 could be adaptively adjusted as $\mathcal{M}_k = \rho_k \mathbf{I}$, where $\rho_k = \max\{\rho_{\min}, \beta \delta_2^k / \delta_1^k\}$ with $\rho_{\min} > 0$,

$$\delta_1^k = \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 \quad \text{and} \quad \delta_2^k = \|A(\mathbf{x}^k - \mathbf{x}^{k-1})\|^2.$$

Since δ_2^k / δ_1^k is an underestimate of the largest eigenvalue of $A^\top A$, to ensure convergence, the safeguard lower bound ρ_{\min} should be increased during the optimization if necessary. One may see [3, Remark 4.2] for more details.

(II) When the set \mathcal{Y} is bounded, we may even use a positive-indefinite proximal matrix

$$L = \tau\chi \mathbf{I} - \beta B^\top B, \quad \text{where} \quad \beta \|B^\top B\| \leq \chi < +\infty \quad \text{and} \quad \tau \in [-1, 1],$$

in the update of \mathbf{y} -subproblem. In this case, denoting $\mathcal{N}_\mathcal{Y} = \sup\{\|\mathbf{y}_1 - \mathbf{y}_2\| : \mathbf{y}_1, \mathbf{y}_2 \in \mathcal{Y}\} < \infty$, analogous to Theorem 4.1 we can show

$$\begin{aligned} & \mathbb{E} \left[F(\mathbf{w}_T) - F(\mathbf{w}) + (\mathbf{w}_T - \mathbf{w})^\top \mathcal{J}(\mathbf{w}) \right] \\ & \leq \frac{1}{2T} \left\{ \sigma^2 \sum_{k=\kappa}^{\kappa+T} \eta_k m_k + \frac{4}{m_\kappa(m_\kappa+1)\eta_\kappa} \|\mathbf{x} - \check{\mathbf{x}}^\kappa\|_{\mathcal{H}}^2 + \|\mathbf{w} - \mathbf{w}^\kappa\|_{\tilde{\mathcal{Q}}_\kappa}^2 \right. \\ & \quad \left. + \omega_1 \|A\mathbf{x}^\kappa + B\mathbf{y}^\kappa - \mathbf{b}\|^2 + \omega_2 \text{const} \right\}, \end{aligned} \quad (4.26)$$

where

$$\text{const} = \begin{cases} \mathcal{N}_\mathcal{Y}^2(\beta \|B\|^2 - \tau\chi), & \text{if } \tau \in [-1, 0], \\ \beta(\mathcal{N}_\mathcal{Y} \|B\|)^2 + \tau\chi \|\mathbf{y}^{\kappa-1} - \mathbf{y}^\kappa\|^2, & \text{if } \tau \in (0, 1]. \end{cases}$$

This means that the results of Theorem 4.2 could still hold even when the proximal matrix L is positive-indefinite.

(III) Similar ideas of Algorithm 1 can be further generalized to solve separable convex optimization with one or multi-block structures. For content focus of the paper, we leave these discussions in the Appendix.

5 Numerical experiments

In this section, we apply the proposed algorithm to solve the following graph-guided fused lasso problem in machine learning:

$$\min_{\mathbf{x}} \frac{1}{N} \sum_{j=1}^N f_j(\mathbf{x}) + \mu \|\mathbf{A}\mathbf{x}\|_1,$$

where $f_j(\mathbf{x}) = \log(1 + \exp(-b_j a_j^T \mathbf{x}))$ denotes the logistic loss function on the feature-label pair $(a_j, b_j) \in \mathbb{R}^l \times \{-1, 1\}$, $N(> l)$ is the data size, $\mu > 0$ is a given regularization parameter, and $\mathbf{A} = [\mathbf{G}; \mathbf{I}]$ is a matrix encoding the feature sparsity pattern. Here, \mathbf{G} is the sparsity pattern of the graph that is obtained by sparse inverse covariance estimation [10]. Introducing an auxiliary variable \mathbf{y} , the above problem is equivalent to the problem

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{y}} \quad & F(\mathbf{x}, \mathbf{y}) := \frac{1}{N} \sum_{j=1}^N f_j(\mathbf{x}) + \mu \|\mathbf{y}\|_1 \\ \text{s.t.} \quad & \mathbf{A}\mathbf{x} - \mathbf{y} = \mathbf{0}, \end{aligned} \quad (5.1)$$

which has the format of our model (1.1). In addition, it can be easily verified that the Assumptions 2.1-2.2 hold. Since the coefficient matrix of the \mathbf{y} variable in the constraints of (5.1) is $-\mathbf{I}$, the \mathbf{y} -subproblem will have a closed-form solution by simply setting $L = \mathbf{0}$ in Algorithm 1. Otherwise, the linearization techniques discussed in (3.2) on choosing L can be applied to obtain a closed-form solution of the \mathbf{y} -subproblem. With $L = \mathbf{0}$, the subproblems in Algorithm 1 would have the following closed-form solution:

$$\begin{cases} \tilde{\mathbf{x}}_{t+1} = [\gamma_t \mathcal{H} + \mathcal{M}_k]^{-1} [\gamma_t \mathcal{H} \tilde{\mathbf{x}}_t + \mathcal{M}_k \mathbf{x}^k - \mathbf{d}_t - \mathbf{h}^k], \\ \mathbf{y}^{k+1} = \text{Shrink}\left(\frac{\mu}{\beta}, \mathbf{A}\mathbf{x}^{k+1} - \frac{\boldsymbol{\lambda}^{k+\frac{1}{2}}}{\beta}\right). \end{cases} \quad (5.2)$$

Here, $\text{Shrink}(\cdot, \cdot)$ denotes the soft shrinkage operator and can be evaluated using the MATLAB built-in function “wthresh”.

In the numerical experiments, the penalty parameter in SAS-ADMM is taken as $\beta = 0.001$, the matrices \mathcal{M}_k are updated adaptively by the strategy explained in Remark 4.1 (I) with initial values $\rho_0 = 1$, $\rho_{\min} = 10^{-5}$ and $\mathcal{H} = 2 \times 10^{-5} \mathbf{I}$. The other parameters as well as the vector \mathbf{e}_t in SAS-ADMM (i.e. Algorithm 1) are chosen the same way as that used in [3, Section 7.1], that is

$$\mathbf{e}_t = \begin{cases} \nabla f(\mathbf{x}_{k-1}) - \nabla f_{\tilde{\zeta}_t}(\mathbf{x}_{k-1}), & \text{if } m_k > l, \\ \mathbf{0}, & \text{otherwise,} \end{cases} \quad (5.3)$$

where \mathbf{x}_k is the ergodic mean of the \mathbf{x} -iterates. Motivated from Theorem 4.2, we use

$$\text{Obj_err} = \frac{|F(\mathbf{x}, \mathbf{y}) - F^*|}{\max\{F^*, 1\}} \quad \text{and} \quad \text{Equ_err} = \|\mathbf{A}\mathbf{x} - \mathbf{y}\|$$

to denote the relative objective value error and the constraint violation error. Here, F^* is the approximate optimal objective function value obtained by running Algorithm 1 for more than 10 minutes. To measure the performance of a algorithm, we plot the maximum of the relative objective error and the constraint error, that is

$$\text{Opt_err} = \max(\text{Obj_err}, \text{Equ_err}),$$

against the CPU time used. All experiments are implemented in MATLAB R2018a (64-bit) with the same starting point $(\mathbf{x}^0, \mathbf{y}^0, \lambda^0) = (\mathbf{0}, \mathbf{0}, \mathbf{0})$ and performed on a PC with Windows 10 operating system, with an Intel i7-8700K CPU and 16GB RAM.

We compare the numerical performance of the proposed algorithm SAS-ADMM[‡] using stepsizes $(\tau, s) = (0.9, 1.09)$, which is suggested in [1] for GS-ADMM, and AS-ADMM [3] for solving problem (5.1) on the dataset *mnist* (including 11,791 samples and 784 features, that is, $(N, I) = (11791, 784)$) downloaded from LIBSVM website. The regularization parameter μ in (5.1) is set as 10^{-5} . For both SAS-ADMM and AS-ADMM, we plot the error associated with the iterates over the first 1/3 of the total CPU time budget, followed by the error associated with the ergodic iterates over the last 2/3 of the budget. We make 10 and 20 successive runs of each algorithm under the CPU time budgets 120s and 200s, respectively. The average comparison results on Opt_err are shown in Fig. 1, and the comparison of the finally obtained iterative solution \mathbf{x}^{k+1} and $\text{hist}(\mathbf{x}^{k+1})$ are shown in Figs. 2-3. Here, we only compare SAS-ADMM with AS-ADMM since in [3] AS-ADMM was shown competitive or better than other state-of-the-art deterministic and stochastic methods. Note that Opt_err has a big drop at around 1/3 of the CPU time budget, the point where the ergodic iterates are started to use for reporting the objective value. From Fig. 1, we can see that SAS-ADMM initially performs worse than AS-ADMM at the beginning iterations. But after the first 1/3 of the total CPU time budget, the SAS-ADMM eventually seems to perform better than AS-ADMM. Finally, Figs. 2-3 show that both the comparison algorithms indeed get sparse solutions.

6 Conclusion

We proposed a symmetric accelerated stochastic alternating direction method of multipliers, called SAS-ADMM, whose dual variables are symmetrically updated. We gave the specific dual stepsizes region ensuring the global convergence, which is larger than those in the literature. Under proper choice of the algorithm parameters, we proved the convergence of SAS-ADMM in expectation with the worst-case $\mathcal{O}(1/T)$ convergence rate, where T represents the number of iterations. Our preliminary experiments showed that by symmetrically updating the dual variables using a more flexible region, SAS-ADMM could outperform AS-ADMM, which only updates the dual variable once, for solving some structured optimization problems arising in machine learning.

[‡]All codes are available at <https://github.com/bjc1987/bjc1987.github.io>

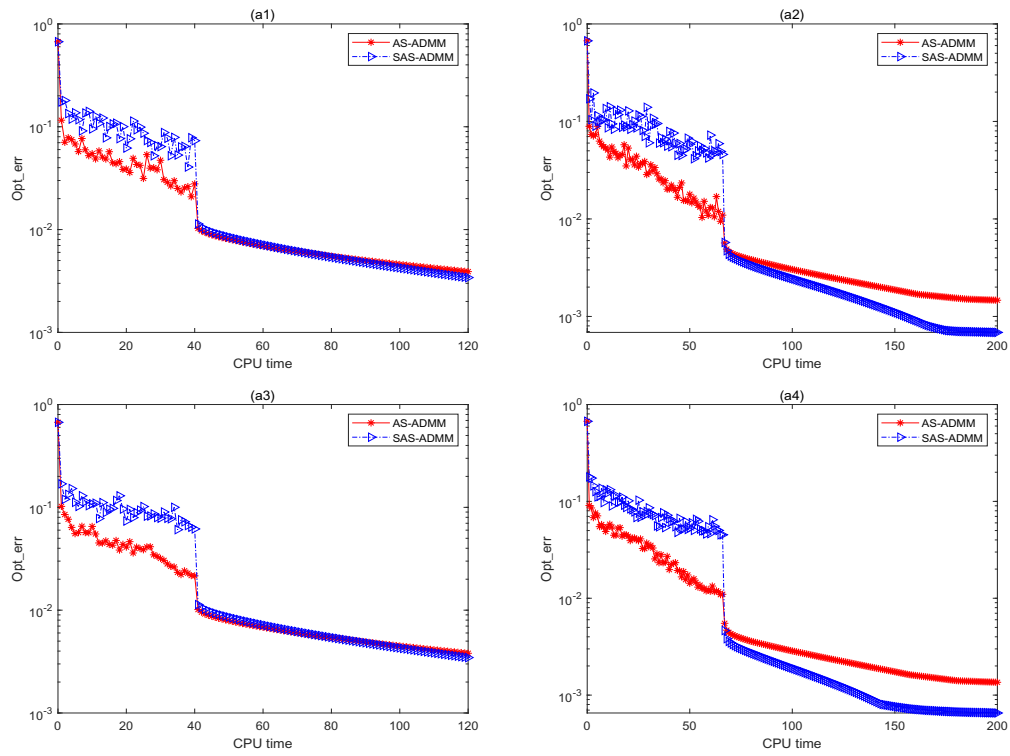


Figure 1: Comparison of Opt_err vs CPU time for Problem (5.1) on the *mnist* dataset: (a1)-(a2) are the results after 10 successive runs; (a3)-(a4) are the results after 20 successive runs.

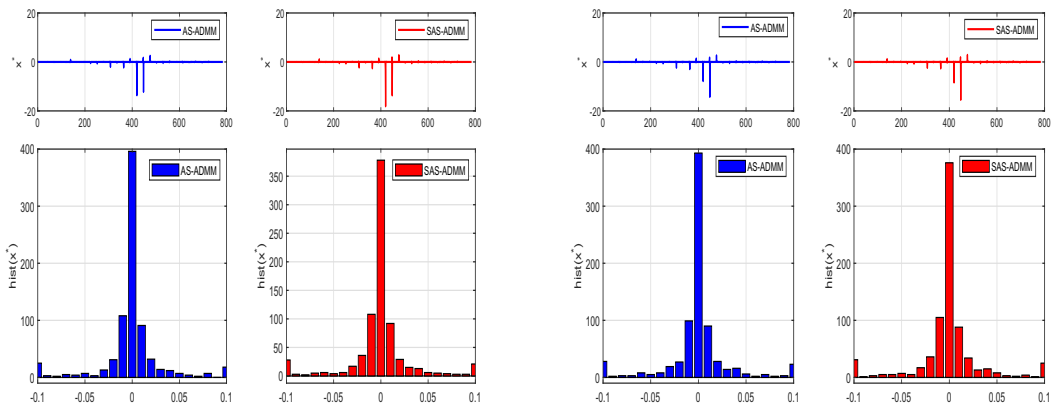


Figure 2: Comparison of the finally obtained iterate x^{k+1} and $\text{hist}(x^{k+1})$ after 10 successive runs: the left two subfigures correspond to (a1); the right two subfigures correspond to (a2).

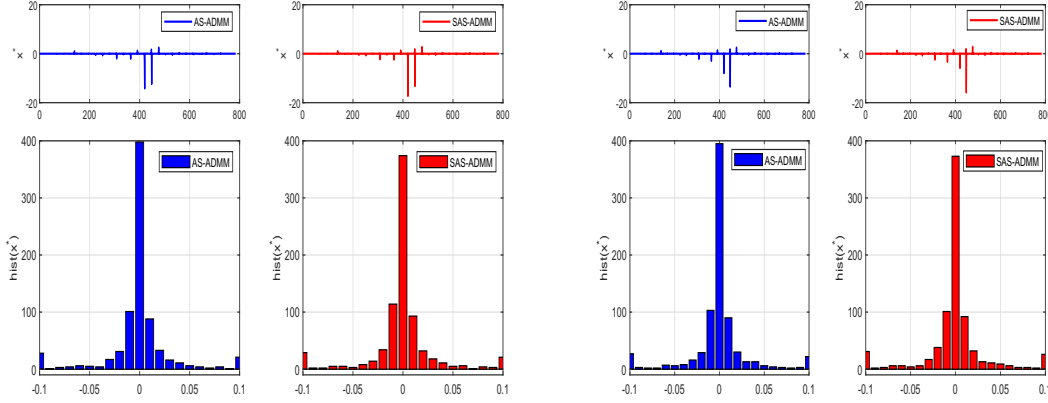


Figure 3: Comparison of the finally obtained iterate x^{k+1} and $\text{hist}(x^{k+1})$ after 20 successive runs: the left two subfigures correspond to (a3); the right two subfigures correspond to (a4).

Acknowledgments

This work was supported by the National Natural Science Foundation of China (12001430, 12131004, 12126603, 72071158), the China Postdoctoral Science Foundation (2020M683545), the Foundation of National Key Laboratory of Science and Technology on Aerodynamic Design and Research (614220119040101) and the USA National Science Foundation (1819161, 2110722).

Appendix: Further discussions

In this section, we discuss 3-block extensions of Algorithm 1 and its variance of a stochastic augmented Lagrangian method.

A.1 A stochastic ALM

We first consider a stochastic augmented Lagrangian method, a variant of SAS-ADMM, to solve

$$\min\{f(x) \mid Ax = b, x \in \mathcal{X}\}, \quad (\text{A.1})$$

where $\mathcal{X} \subset \mathbb{R}^{n_1}$ is a closed convex subset, and f is an average of N smooth convex functions as defined in (1.1). Now, the augmented Lagrangian of (A.1) is

$$\mathcal{L}_\beta(x, \lambda) := \mathcal{L}(x, \lambda) + \frac{\beta}{2} \|Ax - b\|^2,$$

where $\mathcal{L}(x, \lambda) = f(x) - \lambda^\top (Ax - b)$. Then, based on Algorithm 1, we can propose the following Accelerated Stochastic ALM (AS-ALM), Algorithm 2). Similar to SAS-ADMM, we

Algorithm 2 Accelerated Stochastic ALM (AS-ALM)**Parameters:** $\beta > 0, s \in (0, 2]$ and $\mathcal{H} \succ \mathbf{0}$.**Initialization:** $(\mathbf{x}^0, \boldsymbol{\lambda}^0) \in \mathcal{X} \times \mathbb{R}^n := \Omega$ and $\tilde{\mathbf{x}}^0 = \mathbf{x}^0$.For $k=0, 1, \dots$ Choose $m_k > 0$, $\eta_k > 0$ and \mathcal{M}_k such that $\mathcal{M}_k - \beta A^\top A \succeq \mathbf{0}$. $\mathbf{h}^k := -A^\top [\boldsymbol{\lambda}^k - \beta(A\mathbf{x}^k - \mathbf{b})]$. $(\mathbf{x}^{k+1}, \tilde{\mathbf{x}}^{k+1}) = \mathbf{xsub}(\mathbf{x}^k, \tilde{\mathbf{x}}^k, \mathbf{h}^k)$ with \mathbf{xsub} given in Algorithm 1. $\boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^k - s\beta(A\mathbf{x}^{k+1} - \mathbf{b})$.

end

can easily establish the following lemmas on AS-ALM. However, in this case, the convergence region for the dual stepsize can be enlarged from $(0, (\sqrt{5}+1)/2]$ of AS-ADMM [3] to $(0, 2]$.

Lemma A.1. *Let $\{\mathbf{x}^k\}$ be generated by Algorithm 2 and $\eta_k \in (0, 1/\nu)$. Then, the inequality (4.1) holds with*

$$\tilde{\boldsymbol{\lambda}}^k = \boldsymbol{\lambda}^k - \beta(A\mathbf{x}^{k+1} - \mathbf{b}). \quad (\text{A.2})$$

For the iterates generated by Algorithm 2, in this subsection let $\mathbf{w}^k = (\mathbf{x}^k; \boldsymbol{\lambda}^k)$ and $\tilde{\mathbf{w}}^k = (\mathbf{x}^{k+1}; \tilde{\boldsymbol{\lambda}}^k)$, where $\tilde{\boldsymbol{\lambda}}^k$ is defined in (A.2). Then, we have the following lemma.

Lemma A.2. *Let $\{\mathbf{w}^k\}$ be generated by Algorithm 2 and $\eta_k \in (0, 1/\nu)$. Then, we have $\tilde{\mathbf{w}}^k \in \Omega$ and*

$$f(\mathbf{x}) - f(\mathbf{x}^{k+1}) + \langle \mathbf{w} - \tilde{\mathbf{w}}^k, \mathcal{J}(\mathbf{w}) \rangle \geq (\mathbf{w} - \tilde{\mathbf{w}}^k)^\top Q_k (\mathbf{w}^k - \tilde{\mathbf{w}}^k) + \zeta^k$$

for all $\mathbf{w} \in \Omega$, where ζ^k is given by (4.3),

$$\mathcal{J}(\mathbf{w}) = \begin{pmatrix} -A^\top \boldsymbol{\lambda} \\ A\mathbf{x} - \mathbf{b} \end{pmatrix} \quad \text{and} \quad Q_k = \begin{bmatrix} \mathcal{D}_k & \\ & \frac{1}{\beta} \mathbf{I} \end{bmatrix}.$$

Proof. Combining the inequality (4.1) and the relation $A\mathbf{x}^{k+1} - \mathbf{b} = \frac{1}{\beta}(\boldsymbol{\lambda}^k - \tilde{\boldsymbol{\lambda}}^k)$ gives the results. \square

Lemma A.3. *Let $\{\mathbf{w}^k\}$ be generated by Algorithm 2 and $\eta_k \in (0, 1/\nu)$. Then, for any $\mathbf{w} \in \Omega$, we have*

$$\begin{aligned} & f(\mathbf{x}) - f(\mathbf{x}^{k+1}) + (\mathbf{w} - \tilde{\mathbf{w}}^k)^\top \mathcal{J}(\mathbf{w}) \\ & \geq \frac{1}{2} \left\{ \left\| \mathbf{w} - \mathbf{w}^{k+1} \right\|_{\tilde{Q}_k}^2 - \left\| \mathbf{w} - \mathbf{w}^k \right\|_{\tilde{Q}_k}^2 + \left\| \mathbf{w}^k - \tilde{\mathbf{w}}^k \right\|_{\tilde{G}_k}^2 \right\} + \zeta^k, \end{aligned}$$

where ζ^k is given by (4.3),

$$\tilde{Q}_k = \begin{bmatrix} \mathcal{D}_k & \\ & \frac{1}{s\beta} \mathbf{I} \end{bmatrix} \quad \text{and} \quad \tilde{G}_k = \begin{bmatrix} \mathcal{D}_k & \\ & \frac{2-s}{\beta} \mathbf{I} \end{bmatrix}.$$

Proof. The proof is similar to that of Corollary 4.1 and is omitted. \square

Finally, under the conditions of Theorem 4.2, by Lemma A.3 and a similar proof of Theorem 4.2, we can deduce that for any $\mathbf{w}_T := \frac{1}{1+T} \sum_{k=\kappa}^{\kappa+T} \tilde{\mathbf{w}}^k$ and $\kappa \geq 0$, it has

$$|\mathbb{E}[f(\mathbf{x}_T) - f(\mathbf{x}^*)]| = E_\varrho(T) = \mathbb{E}[\|\mathbf{A}\mathbf{x}_T - \mathbf{b}\|],$$

where $E_\varrho(T) = \mathcal{O}(1/T)$ for $\varrho > 1$ and $E_\varrho(T) = \mathcal{O}(T^{-1} \log T)$ for $\varrho = 1$.

A.2 Three-block extensions

Consider a 3-block extension of problem (1.1)

$$\begin{aligned} \min \quad & F(\mathbf{w}) := f(\mathbf{x}) + g(\mathbf{y}) + l(\mathbf{z}) \\ \text{s.t.} \quad & \mathcal{K}\mathbf{w} := \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{C}\mathbf{z} = \mathbf{b}, \\ & \mathbf{x} \in \mathcal{X}, \quad \mathbf{y} \in \mathcal{Y}, \quad \mathbf{z} \in \mathcal{Z}, \end{aligned} \quad (\text{A.3})$$

where l is a closed convex function, $\mathbf{C} \in \mathbb{R}^{n \times n_3}$ is a given matrix, $\mathcal{Z} \subset \mathbb{R}^{n_3}$ is a simple closed convex subset, and the other functions and variables remain the same definitions as those in problem (1.1). Here, the additional function l can be possibly used to promote some data structure different from the structure promoted by g . For convenience, in this subsection, let us define $\mathcal{K}\mathbf{w} := \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{C}\mathbf{z}$, denote $\mathbf{w} = (\mathbf{z}; \mathbf{x}; \mathbf{y}; \boldsymbol{\lambda})$,

$$\mathbf{w}^k = \begin{pmatrix} \mathbf{x}^k \\ \mathbf{y}^k \\ \mathbf{z}^k \\ \boldsymbol{\lambda}^k \end{pmatrix}, \quad \tilde{\mathbf{w}}^k := \begin{pmatrix} \tilde{\mathbf{x}}^k \\ \tilde{\mathbf{y}}^k \\ \tilde{\mathbf{z}}^k \\ \tilde{\boldsymbol{\lambda}}^k \end{pmatrix} = \begin{pmatrix} \mathbf{x}^{k+1} \\ \mathbf{y}^{k+1} \\ \mathbf{z}^{k+1} \\ \tilde{\boldsymbol{\lambda}}^k \end{pmatrix} \quad \text{and} \quad \mathcal{J}(\mathbf{w}) = \begin{pmatrix} -\mathbf{A}^\top \boldsymbol{\lambda} \\ -\mathbf{B}^\top \boldsymbol{\lambda} \\ -\mathbf{C}^\top \boldsymbol{\lambda} \\ \mathcal{K}\mathbf{w} - \mathbf{b} \end{pmatrix}, \quad (\text{A.4})$$

where $\tilde{\boldsymbol{\lambda}}^k$ will be specified differently in the following two discussion cases.

A.2.1 Extension in Gauss-Seidel update

For this case, we need an assumption that $\mathbf{C}^\top \mathbf{A} = \mathbf{0}$. Then SAS-ADMM can be directly extended to Algorithm 3 for solving the 3-block problem (A.3), where the variable updating order is $\mathbf{z}^{k+1} \rightarrow \mathbf{x}^{k+1} \rightarrow \mathbf{y}^{k+1} \rightarrow \boldsymbol{\lambda}^{k+1}$ in a Gauss-Seidel scheme. Now, let

$$\tilde{\boldsymbol{\lambda}}^k = \boldsymbol{\lambda}^k - \beta (\mathbf{C}\mathbf{z}^{k+1} + \mathbf{A}\mathbf{x}^{k+1} + \mathbf{B}\mathbf{y}^k - \mathbf{b}), \quad (\text{A.5})$$

$$\mathbf{v}^{k+1} = (\mathbf{x}^{k+1}; \mathbf{y}^{k+1}; \boldsymbol{\lambda}^{k+1}) \quad \text{and} \quad \tilde{\mathbf{v}}^k = (\tilde{\mathbf{x}}^k; \tilde{\mathbf{y}}^k; \tilde{\boldsymbol{\lambda}}^k). \quad (\text{A.6})$$

Then, we have the following main lemma for the convergence of Algorithm 3.

Algorithm 3 Extension of SAS-ADMM in Gauss-Seidel update**Parameters:** $\beta > 0$, $\mathcal{H} \succ \mathbf{0}$, $L \succeq \mathbf{0}$ and $(\tau, s) \in \Delta$.**Initialization:** $(\mathbf{x}^0, \mathbf{y}^0, \mathbf{z}^0, \boldsymbol{\lambda}^0) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z} \times \mathbb{R}^n := \Omega$, $\tilde{\mathbf{x}}^0 = \mathbf{x}^0$.For $k=0, 1, \dots$ Choose $m_k > 0$, $\eta_k > 0$ and \mathcal{M}_k such that $\mathcal{M}_k - \beta A^\top A \succeq \mathbf{0}$.

$$\mathbf{z}^{k+1} \in \operatorname{argmin}_{\mathbf{z} \in \mathcal{Z}} l(\mathbf{z}) + \frac{\beta}{2} \left\| C\mathbf{z} + A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b} - \frac{\boldsymbol{\lambda}^k}{\beta} \right\|^2.$$

$$\mathbf{h}^k := -A^\top \left[\boldsymbol{\lambda}^k - \beta(C\mathbf{z}^{k+1} + A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b}) \right].$$

 $(\mathbf{x}^{k+1}, \tilde{\mathbf{x}}^{k+1}) = \mathbf{xsub}(\mathbf{x}^k, \tilde{\mathbf{x}}^k, \mathbf{h}^k)$ with \mathbf{xsub} given in ALG.1.

$$\boldsymbol{\lambda}^{k+\frac{1}{2}} = \boldsymbol{\lambda}^k - \tau\beta(C\mathbf{z}^{k+1} + A\mathbf{x}^{k+1} + B\mathbf{y}^k - \mathbf{b}).$$

$$\mathbf{y}^{k+1} \in \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{y}) + \frac{\beta}{2} \left\| C\mathbf{z}^{k+1} + A\mathbf{x}^{k+1} + B\mathbf{y} - \mathbf{b} - \frac{\boldsymbol{\lambda}^{k+\frac{1}{2}}}{\beta} \right\|^2 + \frac{1}{2} \|\mathbf{y} - \mathbf{y}^k\|_L^2.$$

$$\boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^{k+\frac{1}{2}} - s\beta(A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} + C\mathbf{z}^{k+1} - \mathbf{b}).$$

end

Lemma A.4. Assume $C^\top A = \mathbf{0}$ and $\eta_k \in (0, 1/\nu)$. Then, the iterates generated by Algorithm 3 satisfy $\tilde{\mathbf{w}}^k \in \Omega$ and

$$\begin{aligned} & F(\mathbf{w}) - F(\tilde{\mathbf{w}}^k) + (\mathbf{w} - \tilde{\mathbf{w}}^k)^\top \mathcal{J}(\mathbf{w}) \\ & \geq \frac{1}{2} \left\{ \|\mathbf{v} - \mathbf{v}^{k+1}\|_{\tilde{Q}_k}^2 - \|\mathbf{v} - \mathbf{v}^k\|_{\tilde{Q}_k}^2 + \|\mathbf{v}^k - \tilde{\mathbf{v}}^k\|_{\tilde{G}_k}^2 \right\} + \zeta^k \end{aligned}$$

for any $\mathbf{w} \in \Omega$, where \tilde{Q}_k, \tilde{G}_k and ζ^k are given in Corollary 4.1 and (4.3), respectively. Moreover, we have

$$\begin{aligned} \|\mathbf{v}^k - \tilde{\mathbf{v}}^k\|_{\tilde{G}_k}^2 & \geq \|\mathbf{x}^k - \mathbf{x}^{k+1}\|_{\mathcal{D}_k}^2 + \omega_0 \|A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}\|^2 \\ & \quad + \omega_1 \left(\|A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b}\|^2 - \|A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b}\|^2 \right) \\ & \quad + \omega_2 \left(\|\mathbf{y}^k - \mathbf{y}^{k+1}\|_L^2 - \|\mathbf{y}^{k-1} - \mathbf{y}^k\|_L^2 \right), \end{aligned}$$

where $\omega_0, \omega_1, \omega_2 \geq 0$ are given in (4.16).

Proof. By the updates of \mathbf{h}^k and $\tilde{\boldsymbol{\lambda}}^k$ in Algorithm 3, it is easy to derive (4.1) as before. Then, according to the first-order optimality condition of \mathbf{z} -subproblem and the assumption that $C^\top A = \mathbf{0}$, we have

$$\mathbf{z}^{k+1} \in \mathcal{Z}, \quad l(\mathbf{z}) - l(\mathbf{z}^{k+1}) + \langle \mathbf{z} - \mathbf{z}^{k+1}, \mathbf{p}_z^k \rangle \geq 0, \quad \forall \mathbf{z} \in \mathcal{Z}, \quad (\text{A.7})$$

where

$$\begin{aligned} \mathbf{p}_z^k &= -C^\top \boldsymbol{\lambda}^k + \beta C^\top (C\mathbf{z}^{k+1} + A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b}) \\ &= -C^\top \tilde{\boldsymbol{\lambda}}^k - \beta C^\top A(\mathbf{x}^{k+1} - \mathbf{x}^k) = -C^\top \tilde{\boldsymbol{\lambda}}^k. \end{aligned}$$

Similarly, we have by the \mathbf{y} -update that

$$\mathbf{y}^{k+1} \in \mathcal{Y}, \quad g(\mathbf{y}) - g(\mathbf{y}^{k+1}) + \langle \mathbf{y} - \mathbf{y}^{k+1}, \mathbf{p}_y^k \rangle \geq 0, \quad \forall \mathbf{y} \in \mathcal{Y}, \quad (\text{A.8})$$

where

$$\begin{aligned} \mathbf{p}_y^k &= -B^\top \boldsymbol{\lambda}^{k+\frac{1}{2}} + \beta B^\top (\mathcal{K}\mathbf{w}^{k+1} - \mathbf{b}) + L(\mathbf{y}^{k+1} - \mathbf{y}^k) \\ &= -B^\top \boldsymbol{\lambda}^{k+\frac{1}{2}} + B^\top (\boldsymbol{\lambda}^k - \tilde{\boldsymbol{\lambda}}^k) + [L + \beta B^\top B] (\mathbf{y}^{k+1} - \mathbf{y}^k) \\ &= -B^\top \tilde{\boldsymbol{\lambda}}^k + \tau B^\top (\boldsymbol{\lambda}^k - \tilde{\boldsymbol{\lambda}}^k) + [L + \beta B^\top B] (\mathbf{y}^{k+1} - \mathbf{y}^k). \end{aligned}$$

Besides, it follows from the updates of $\tilde{\boldsymbol{\lambda}}^k$ that

$$\left\langle \boldsymbol{\lambda} - \tilde{\boldsymbol{\lambda}}^k, \mathcal{K}\tilde{\mathbf{w}}^k - \mathbf{b} + \frac{1}{\beta} (\tilde{\boldsymbol{\lambda}}^k - \boldsymbol{\lambda}^k) - B(\tilde{\mathbf{y}}^k - \mathbf{y}^k) \right\rangle = 0, \quad \forall \boldsymbol{\lambda} \in \mathbb{R}^n. \quad (\text{A.9})$$

Combining the above inequalities (A.7), (A.8), (A.9) with (4.1), we can get

$$F(\mathbf{w}) - F(\tilde{\mathbf{w}}^k) + (\mathbf{w} - \tilde{\mathbf{w}}^k)^\top \mathcal{J}(\tilde{\mathbf{w}}^k) \geq (\mathbf{v} - \tilde{\mathbf{v}}^k)^\top Q_k(\mathbf{v}^k - \tilde{\mathbf{v}}^k) + \zeta^k, \quad (\text{A.10})$$

where ζ^k, Q_k are given by (4.3) and (4.5), respectively. Then, the rest proof will be similar to that of Corollary 4.1 and Lemma 4.4. \square

Based on the above Lemma A.4, the ergodic convergence of Algorithm 3 with a sub-linear convergence rate can be similarly established under the conditions of Theorem 4.2. Here, we omit the detailed proof. Note that, if the \mathbf{z} -subproblem is not easily solvable, one could also add a positive semidefinite proximal term to linearize it. However, the requirement $C^\top A = \mathbf{0}$ is quite strict in applications. In the next subsection we will propose a partially Jacobi update for the primal variables, for which $C^\top A = \mathbf{0}$ is not required.

A.2.2 Extension in partially Jacobi update

Now, let us consider Algorithm 4, where the block variables \mathbf{y} and \mathbf{z} are updated in a Jacobi fashion.

To establish the global convergence of Algorithm 4, we first have the following observations. Denoting

$$\tilde{\boldsymbol{\lambda}}^k = \boldsymbol{\lambda}^k - \beta (A\mathbf{x}^{k+1} + B\mathbf{y}^k + C\mathbf{z}^k - \mathbf{b}) \quad (\text{A.11})$$

Algorithm 4 Extension of SAS-ADMM in partially Jacobi update**Parameters:** $\beta > 0, (\tau, s) \in \Delta$, L_1 and L_2 satisfy (A.22).**Initialization:** $(\mathbf{x}^0, \mathbf{y}^0, \mathbf{z}^0, \boldsymbol{\lambda}^0) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z} \times \mathbb{R}^n := \Omega$, $\tilde{\mathbf{x}}^0 = \mathbf{x}^0$.For $k=0, 1, \dots$ Choose $m_k > 0$, $\eta_k > 0$ and \mathcal{M}_k such that $\mathcal{M}_k - \beta A^\top A \succeq \mathbf{0}$.

$$\mathbf{h}^k := -A^\top \left[\boldsymbol{\lambda}^k - \beta(A\mathbf{x}^k + B\mathbf{y}^k + C\mathbf{z}^k - \mathbf{b}) \right].$$

$$(\mathbf{x}^{k+1}, \tilde{\mathbf{x}}^{k+1}) = \mathbf{xsub}(\mathbf{x}^k, \tilde{\mathbf{x}}^k) \text{ with } \mathbf{xsub} \text{ given in ALG.1.}$$

$$\boldsymbol{\lambda}^{k+\frac{1}{2}} = \boldsymbol{\lambda}^k - \tau\beta(A\mathbf{x}^{k+1} + B\mathbf{y}^k + C\mathbf{z}^k - \mathbf{b}).$$

$$\mathbf{y}^{k+1} \in \arg\min_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{y}) + \frac{\beta}{2} \left\| A\mathbf{x}^{k+1} + B\mathbf{y} + C\mathbf{z}^k - \mathbf{b} - \frac{\boldsymbol{\lambda}^{k+\frac{1}{2}}}{\beta} \right\|^2 + \frac{1}{2} \|\mathbf{y} - \mathbf{y}^k\|_{L_1}^2.$$

$$\mathbf{z}^{k+1} \in \arg\min_{\mathbf{z} \in \mathcal{Z}} l(\mathbf{z}) + \frac{\beta}{2} \left\| A\mathbf{x}^{k+1} + B\mathbf{y}^k + C\mathbf{z} - \mathbf{b} - \frac{\boldsymbol{\lambda}^{k+\frac{1}{2}}}{\beta} \right\|^2 + \frac{1}{2} \|\mathbf{z} - \mathbf{z}^k\|_{L_2}^2.$$

$$\boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^{k+\frac{1}{2}} - s\beta(A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} + C\mathbf{z}^{k+1} - \mathbf{b}).$$

end

and using the first-order optimality condition of the \mathbf{y} -subproblem, we have

$$\mathbf{y}^{k+1} \in \mathcal{Y}, \quad g(\mathbf{y}) - g(\mathbf{y}^{k+1}) + \langle \mathbf{y} - \mathbf{y}^{k+1}, \mathbf{p}_y^k \rangle \geq 0, \quad \forall \mathbf{y} \in \mathcal{Y}, \quad (\text{A.12})$$

where

$$\begin{aligned} \mathbf{p}_y^k &= -B^\top \boldsymbol{\lambda}^{k+\frac{1}{2}} + \beta B^\top (A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} + C\mathbf{z}^k - \mathbf{b}) + L_1 (\mathbf{y}^{k+1} - \mathbf{y}^k) \\ &= -B^\top \boldsymbol{\lambda}^{k+\frac{1}{2}} + \beta B^\top (A\mathbf{x}^{k+1} + B\mathbf{y}^k + C\mathbf{z}^k - \mathbf{b}) + (L_1 + \beta B^\top B) (\mathbf{y}^{k+1} - \mathbf{y}^k) \\ &= -B^\top \tilde{\boldsymbol{\lambda}}^k + \tau B^\top (\boldsymbol{\lambda}^k - \tilde{\boldsymbol{\lambda}}^k) + (L_1 + \beta B^\top B) (\mathbf{y}^{k+1} - \mathbf{y}^k), \end{aligned}$$

and we use the relationship

$$\boldsymbol{\lambda}^{k+\frac{1}{2}} = \boldsymbol{\lambda}^k - \tau(\boldsymbol{\lambda}^k - \tilde{\boldsymbol{\lambda}}^k). \quad (\text{A.13})$$

Combining (A.12) and the definition of \mathbf{p}_y^k , we have

$$\begin{aligned} g(\mathbf{y}) - g(\mathbf{y}^{k+1}) + \langle \mathbf{y} - \mathbf{y}^{k+1}, -B^\top \boldsymbol{\lambda}^{k+\frac{1}{2}} + \beta B^\top (\mathcal{K}\mathbf{w}^{k+1} - \mathbf{b}) \\ - \beta B^\top C(\mathbf{z}^{k+1} - \mathbf{z}^k) + L_1 (\mathbf{y}^{k+1} - \mathbf{y}^k) \rangle \geq 0. \end{aligned} \quad (\text{A.14})$$

Similarly, by the first-order optimality condition of the \mathbf{z} -subproblem, we have

$$\begin{aligned} l(\mathbf{z}) - l(\mathbf{z}^{k+1}) + \langle \mathbf{z} - \mathbf{z}^{k+1}, -B^\top \boldsymbol{\lambda}^{k+\frac{1}{2}} + \beta C^\top (\mathcal{K}\mathbf{w}^{k+1} - \mathbf{b}) \\ - \beta C^\top B(\mathbf{y}^{k+1} - \mathbf{y}^k) + L_2 (\mathbf{z}^{k+1} - \mathbf{z}^k) \rangle \geq 0. \end{aligned} \quad (\text{A.15})$$

Adding the above two inequalities (A.14) and (A.15), we can see $(\mathbf{y}^{k+1}, \mathbf{z}^{k+1})$ satisfies the first-order optimality condition, hence is a solution, of the following problem

$$(\mathbf{y}^{k+1}, \mathbf{z}^{k+1}) \in \arg \min_{\mathbf{y} \in \mathcal{Y}, \mathbf{z} \in \mathcal{Z}} \left\{ \begin{array}{l} g(\mathbf{y}) + l(\mathbf{z}) + \frac{1}{2} \|(\mathbf{y} - \mathbf{y}^k, \mathbf{z} - \mathbf{z}^k)\|_{\bar{L}}^2 \\ + \frac{\beta}{2} \left\| A\mathbf{x}^{k+1} + B\mathbf{y} + C\mathbf{z} - \mathbf{b} - \frac{\boldsymbol{\lambda}^{k+\frac{1}{2}}}{\beta} \right\|^2 \end{array} \right\}, \quad (\text{A.16})$$

where

$$\bar{L} = \begin{bmatrix} L_1 & -\beta B^\top C \\ -\beta C^\top B & L_2 \end{bmatrix}. \quad (\text{A.17})$$

Hence, by considering (\mathbf{y}, \mathbf{z}) as one block variable, Algorithm 4 is essentially a particular version of Algorithm 1 for solving a 2-block problem with L and B being replaced by \bar{L} and (B, C) , respectively.

From the above observations, we can directly establish the following properties of Algorithm 4.

Lemma A.5. *The iterates generated by Algorithm 4 satisfy*

$$F(\mathbf{w}) - F(\tilde{\mathbf{w}}^k) + \langle \mathbf{w} - \tilde{\mathbf{w}}^k, \mathcal{J}(\mathbf{w}) \rangle \geq (\mathbf{w} - \tilde{\mathbf{w}}^k)^\top Q_k (\mathbf{w}^k - \tilde{\mathbf{w}}^k) + \zeta^k$$

for any $\mathbf{w} \in \Omega$, where ζ^k is given by (4.3),

$$Q_k = \begin{bmatrix} \mathcal{D}_k & & & \\ & L_1 + \beta B^\top B & & -\tau B^\top \\ & -B & L_2 + \beta C^\top C & -\tau C^\top \\ & & -C & \frac{1}{\beta} \mathbf{I} \end{bmatrix}. \quad (\text{A.18})$$

Proof. Notice that

$$\bar{L} + \beta(B, C)^\top (B, C) = \begin{bmatrix} L_1 + \beta B^\top B & \\ & L_2 + \beta C^\top C \end{bmatrix}.$$

So, replacing L and B in Lemma 4.2 by \bar{L} and (B, C) , respectively, this lemma directly follows from Lemma 4.2. \square

Similarly, identifying L and B in (4.11) with \bar{L} and (B, C) , respectively, it follows from Corollary 4.1 that

$$\begin{aligned} & F(\mathbf{w}) - F(\tilde{\mathbf{w}}^k) + \langle \mathbf{w} - \tilde{\mathbf{w}}^k, \mathcal{J}(\mathbf{w}) \rangle \\ & \geq \frac{1}{2} \left\{ \left\| \mathbf{w} - \mathbf{w}^{k+1} \right\|_{\tilde{Q}_k}^2 - \left\| \mathbf{w} - \mathbf{w}^k \right\|_{\tilde{Q}_k}^2 + \left\| \mathbf{w}^k - \tilde{\mathbf{w}}^k \right\|_{\tilde{G}_k}^2 \right\} + \zeta^k, \end{aligned} \quad (\text{A.19})$$

where ζ^k is given by (4.3) and

$$\tilde{Q}_k = \begin{bmatrix} \mathcal{D}_k & & & \\ L_1 + (1 - \frac{\tau s}{\tau+s})\beta B^\top B & -\frac{\tau s}{\tau+s}\beta B^\top C & -\frac{\tau}{\tau+s}B^\top & \\ -\frac{\tau s}{\tau+s}\beta C^\top B & L_2 + (1 - \frac{\tau s}{\tau+s})\beta C^\top C & -\frac{\tau}{\tau+s}C^\top & \\ -\frac{\tau}{\tau+s}B & -\frac{\tau}{\tau+s}C & \frac{1}{(\tau+s)\beta}I & \end{bmatrix}, \quad (\text{A.20})$$

$$\tilde{G}_k = \begin{bmatrix} \mathcal{D}_k & & \\ L_1 + (1-s)\beta B^\top B & -s\beta B^\top C & (s-1)B^\top & \\ -s\beta C^\top B & L_2 + (1-s)\beta C^\top C & (s-1)C^\top & \\ (s-1)B & (s-1)C & \frac{2-\tau-s}{\beta}I & \end{bmatrix}. \quad (\text{A.21})$$

Then, we have the following estimate on a lower bound of $\|w^k - \tilde{w}^k\|_{\tilde{G}_k}$.

Lemma A.6. Suppose there exist $\gamma_1 > 0$ and $\gamma_2 > 0$ with $\gamma_1\gamma_2 \geq 1$ such that

$$L_1 \succeq \gamma_1 \beta B^\top B \quad \text{and} \quad L_2 \succeq \gamma_2 \beta C^\top C. \quad (\text{A.22})$$

Then, for any $(\tau, s) \in \Delta$ defined in (3.1), we have \tilde{Q}_k defined in (A.20) is positive semidefinite and

$$\begin{aligned} \|w^k - \tilde{w}^k\|_{\tilde{G}_k}^2 &\geq \|x^k - x^{k+1}\|_{\mathcal{D}_k}^2 + \omega_0 \|\mathcal{K}w^{k+1} - b\|^2 \\ &\quad + \omega_1 \left(\|\mathcal{K}w^{k+1} - b\|^2 - \|\mathcal{K}w^k - b\|^2 \right) \\ &\quad + \omega_2 \left(\left\| \begin{pmatrix} y^k - y^{k+1} \\ z^k - z^{k+1} \end{pmatrix} \right\|_{\bar{L}}^2 - \left\| \begin{pmatrix} y^k - y^{k-1} \\ z^k - z^{k-1} \end{pmatrix} \right\|_{\bar{L}}^2 \right), \end{aligned} \quad (\text{A.23})$$

where $\omega_0, \omega_1, \omega_2 \geq 0$ is defined in (4.16) and \bar{L} is defined in (A.17).

Proof. First, since $L_1 \succeq \gamma_1 \beta B^\top B$ and $L_2 \succeq \gamma_2 \beta C^\top C$, it follows from $\gamma_1 > 0, \gamma_2 > 0$ and $\gamma_1\gamma_2 \geq 1$ that

$$\bar{L} = \begin{bmatrix} L_1 & -\beta B^\top C \\ -\beta C^\top B & L_2 \end{bmatrix} \succeq \beta \begin{bmatrix} \gamma_1 B^\top B & -B^\top C \\ -C^\top B & \gamma_2 C^\top C \end{bmatrix} \succeq 0. \quad (\text{A.24})$$

By Lemma 4.3, we have \tilde{Q}_k defined in (A.20) is positive semidefinite if

$$\bar{L} \succeq (\tau-1)\beta(B, C)^\top(B, C). \quad (\text{A.25})$$

Since $\tau \leq 1$ for any $(\tau, s) \in \Delta$, we have $0 \succeq (\tau-1)\beta(B, C)^\top(B, C)$. Therefore, we have from (A.24) that (A.25) holds automatically and therefore, \tilde{Q}_k defined in (A.20) is positive semidefinite. Furthermore, it follows from Lemma 4.4 that (A.23) holds as long as $\bar{L} \succeq 0$ which is verified by (A.24). \square

Now, defining $\mathbf{w}_T := \frac{1}{T} \sum_{k=\kappa}^{\kappa+T} \tilde{\mathbf{w}}^k$ for some integers $T > 0$ and $\kappa > 0$, under the same conditions in Theorem 4.1, by Lemma A.6 and similar to the proof of Theorem 4.1, we can obtain

$$\begin{aligned} & \mathbb{E} \left[F(\mathbf{w}_T) - F(\mathbf{w}) + (\mathbf{w}_T - \mathbf{w})^\top \mathcal{J}(\mathbf{w}) \right] \\ & \leq \frac{1}{2T} \left\{ \sigma^2 \sum_{k=\kappa}^{\kappa+T} \eta_k m_k + \frac{4}{m_\kappa(m_\kappa+1)\eta_\kappa} \|\mathbf{x} - \check{\mathbf{x}}^\kappa\|_{\mathcal{H}}^2 + \|\mathbf{w} - \mathbf{w}^\kappa\|_{\tilde{\mathcal{Q}}_\kappa}^2 \right. \\ & \quad \left. + \omega_1 \|\mathcal{K}\mathbf{w}^\kappa - \mathbf{b}\|^2 + \omega_2 \left\| \begin{pmatrix} \mathbf{y}^\kappa - \mathbf{y}^{\kappa-1} \\ \mathbf{z}^\kappa - \mathbf{z}^{\kappa-1} \end{pmatrix} \right\|_{\tilde{\mathcal{L}}}^2 \right\}, \end{aligned}$$

where $\omega_1 \geq 0$ and $\omega_2 \geq 0$ given in (4.16). So, by the choice of the parameters (η_k, m_k) chosen in Theorem 4.2, we can obtain

$$|\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)]| = E_\varrho(T) = \mathbb{E}[\|A\mathbf{x}_T + B\mathbf{y}_T + C\mathbf{z}_T - \mathbf{b}\|],$$

where $E_\varrho(T) = \mathcal{O}(1/T)$ for the parameter $\varrho > 1$ and $E_\varrho(T) = \mathcal{O}(T^{-1} \log T)$ for $\varrho = 1$.

Remark A.1. Observing from the above analysis, Algorithm 4 could be in fact generalized to Algorithm 5 for solving the multi-block separable convex optimization:

$$\begin{aligned} \min \quad & F(\mathbf{w}) := f(\mathbf{x}) + \sum_{i=1}^q g_i(\mathbf{y}_i) \\ \text{s.t.} \quad & \mathcal{K}\mathbf{w} := A\mathbf{x} + \sum_{i=1}^q B_i \mathbf{y}_i = \mathbf{b}, \\ & \mathbf{x} \in \mathcal{X}, \quad \mathbf{y}_i \in \mathcal{Y}_i, \quad i = 1, 2, \dots, q, \end{aligned} \tag{A.26}$$

where f has the same definition as in (1.1), $g_i: \mathcal{Y}_i \rightarrow \mathbb{R} \cup \{+\infty\}$ is a convex but possibly nonsmooth function, $B_i \in \mathbb{R}^{n \times n_i}$ and $\mathcal{Y}_i \subset \mathbb{R}^{n_i}$ is a closed convex subset.

The convergence of Algorithm 5 can be analogously established with proper modifications on the convergence proof of Algorithm 4. Here, we only give a very brief explanation. Denote $g(\mathbf{y}) = \sum_{i=1}^q g_i(\mathbf{y}_i)$, $B = (B_1, \dots, B_q)$, $\mathbf{y} = (\mathbf{y}_1; \dots; \mathbf{y}_q)$, $\mathbf{y}^k = (\mathbf{y}_1^k; \dots; \mathbf{y}_q^k)$ and $\mathcal{Y} = \mathcal{Y}_1 \times \dots \times \mathcal{Y}_q$. Then, by the first-order optimality condition of \mathbf{y}_i -subproblem, we have $\mathbf{y}_i^{k+1} \in \mathcal{Y}_i$ and

$$\begin{aligned} & g_i(\mathbf{y}_i) - g_i(\mathbf{y}_i^{k+1}) + \left\langle \mathbf{y}_i - \mathbf{y}_i^{k+1}, -B_i^\top \boldsymbol{\lambda}^{k+\frac{1}{2}} + \beta B_i^\top (\mathcal{K}\mathbf{w}^{k+1} - \mathbf{b}) - \right. \\ & \quad \left. \beta \sum_{l \neq i, l=1}^q B_l^\top B_l (\mathbf{y}_l^{k+1} - \mathbf{y}_l^k) + L_i (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) \right\rangle \geq 0, \quad \forall \mathbf{y}_i \in \mathcal{Y}_i. \end{aligned}$$

After adding the above inequality from $i = 1$ to q , we can see \mathbf{y}^{k+1} satisfies the first-order optimality condition, hence is a solution, of the following problem:

$$\mathbf{y}^{k+1} \in \arg \min_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{y}) + \frac{1}{2} \|\mathbf{y} - \mathbf{y}^k\|_{\tilde{\mathcal{L}}}^2 + \frac{\beta}{2} \left\| A\mathbf{x}^{k+1} + B\mathbf{y} - \mathbf{b} - \frac{\boldsymbol{\lambda}^{k+\frac{1}{2}}}{\beta} \right\|^2,$$

Algorithm 5 Multi-block extension of SAS-ADMM in partially Jacobi update

Parameters: $\beta > 0, (\tau, s) \in \Delta$ and $L_i \succeq (q-1)\beta B_i^\top B_i$ for all $i=1, \dots, q$.

Initialization: $(\mathbf{x}^0, \mathbf{y}^0, \boldsymbol{\lambda}^0) \in \mathcal{X} \times \mathcal{Y} \times \mathbb{R}^n$, $\check{\mathbf{x}}^0 = \mathbf{x}^0$.

For $k=0, 1, \dots$

Choose $m_k > 0$, $\eta_k > 0$ and \mathcal{M}_k such that $\mathcal{M}_k - \beta A^\top A \succeq \mathbf{0}$.

$\mathbf{h}^k := -A^\top \left[\boldsymbol{\lambda}^k - \beta(A\mathbf{x}^k + B\mathbf{y}^k - \mathbf{b}) \right]$.

$(\mathbf{x}^{k+1}, \check{\mathbf{x}}^{k+1}) = \mathbf{xsub}(\mathbf{x}^k, \check{\mathbf{x}}^k)$ with \mathbf{xsub} given in ALG.1.

$\boldsymbol{\lambda}^{k+\frac{1}{2}} = \boldsymbol{\lambda}^k - \tau\beta(A\mathbf{x}^{k+1} + B\mathbf{y}^k - \mathbf{b})$.

For $i=1, 2, \dots, q$,

$\mathbf{y}_i^{k+1} \in \arg \min_{\mathbf{y}_i \in \mathcal{Y}_i} g_i(\mathbf{y}_i) + \frac{\beta}{2} \left\| A\mathbf{x}^{k+1} + B_i\mathbf{y}_i + \sum_{l \neq i, l=1}^q B_l\mathbf{y}_l^k - \mathbf{b} - \frac{\boldsymbol{\lambda}^{k+\frac{1}{2}}}{\beta} \right\|^2 + \frac{1}{2} \|\mathbf{y}_i - \mathbf{y}_i^k\|_{L_i}^2$.

end

$\boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^{k+\frac{1}{2}} - s\beta(A\mathbf{x}^{k+1} + B\mathbf{y}^{k+1} - \mathbf{b})$.

end

where

$$\tilde{L} = \begin{bmatrix} L_1 & -\beta B_1^\top B_2 & \cdots & -\beta B_1^\top B_q \\ -\beta B_2^\top B_1 & L_2 & \cdots & -\beta B_2^\top B_q \\ \vdots & \vdots & \ddots & \vdots \\ -\beta B_q^\top B_1 & -\beta B_q^\top B_2 & \cdots & L_q \end{bmatrix}. \quad (\text{A.27})$$

So, by a similar analysis to Algorithm 4, the inequality (A.19) holds with

$$\tilde{Q}_k = \left[\begin{array}{c|cccc|c} \mathcal{D}_k & & & & & \\ \hline & L_1 + (1 - \frac{\tau s}{\tau+s})\beta B_1^\top B_1 & -\frac{\tau s}{\tau+s}\beta B_1^\top B_2 & \cdots & -\frac{\tau s}{\tau+s}\beta B_1^\top B_q & -\frac{\tau}{\tau+s}B_1^\top \\ & -\frac{\tau s}{\tau+s}\beta B_2^\top B_1 & L_2 + (1 - \frac{\tau s}{\tau+s})\beta B_2^\top B_2 & \cdots & -\frac{\tau s}{\tau+s}\beta B_2^\top B_q & -\frac{\tau}{\tau+s}B_2^\top \\ & \vdots & \vdots & \ddots & \vdots & \vdots \\ & -\frac{\tau s}{\tau+s}\beta B_q^\top B_1 & -\frac{\tau s}{\tau+s}\beta B_q^\top B_2 & \cdots & L_q + (1 - \frac{\tau s}{\tau+s})\beta B_q^\top B_q & -\frac{\tau}{\tau+s}B_q^\top \\ \hline & -\frac{\tau}{\tau+s}B_1 & -\frac{\tau}{\tau+s}B_2 & \cdots & -\frac{\tau}{\tau+s}B_q & \frac{1}{(\tau+s)\beta}\mathbf{I} \end{array} \right],$$

$$\tilde{G}_k = \left[\begin{array}{c|cccc|c} \mathcal{D}_k & & & & & \\ \hline & L_1 + (1-s)\beta B_1^\top B_1 & -s\beta B_1^\top B_2 & \cdots & -s\beta B_1^\top B_q & (s-1)B_1^\top \\ & -s\beta B_2^\top B_1 & L_2 + (1-s)\beta B_2^\top B_2 & \cdots & -s\beta B_2^\top B_q & (s-1)B_2^\top \\ & \vdots & \vdots & \ddots & \vdots & \vdots \\ & -s\beta B_q^\top B_1 & -s\beta B_q^\top B_2 & \cdots & L_q + (1-s)\beta B_q^\top B_q & (s-1)B_q^\top \\ \hline & (s-1)B_1 & (s-1)B_2 & \cdots & (s-1)B_q & \frac{2-\tau-s}{\beta}\mathbf{I} \end{array} \right].$$

If $L_i \succeq (q-1)\beta B_i^\top B_i$ for $i=1, \dots, q$, then for any $(\tau, s) \in \Delta$ defined by (3.1), the above matrix \tilde{Q}_k is positive semidefinite and

$$\begin{aligned} \|\mathbf{w}^k - \tilde{\mathbf{w}}^k\|_{\tilde{G}_k}^2 &\geq \|\mathbf{x}^k - \mathbf{x}^{k+1}\|_{\mathcal{D}_k}^2 + \omega_0 \|\mathcal{K}\mathbf{w}^{k+1} - \mathbf{b}\|^2 \\ &\quad + \omega_1 \left(\|\mathcal{K}\mathbf{w}^{k+1} - \mathbf{b}\|^2 - \|\mathcal{K}\mathbf{w}^k - \mathbf{b}\|^2 \right) \\ &\quad + \omega_2 \left(\|\mathbf{y}^{k+1} - \mathbf{y}^k\|_{\tilde{L}}^2 - \|\mathbf{y}^k - \mathbf{y}^{k-1}\|_{\tilde{L}}^2 \right), \end{aligned}$$

where $\omega_0, \omega_1, \omega_2 \geq 0$ is defined in (4.16) and \tilde{L} is defined in (A.27). The above discussions imply that Algorithm 5 has the same convergence properties as Algorithm 4 and can be also considered as a stochastic extension of the deterministic GS-ADMM [1] for solving the grouped multi-block separable convex optimization problem.

References

- [1] J. Bai, J. Li, F. Xu and H. Zhang, Generalized symmetric ADMM for separable convex optimization, *Comput. Optim. Appl.*, 70 (2018), 129-170.
- [2] J. Bai, X. Chang, J. Li and F. Xu, Convergence revisit on generalized symmetric ADMM, *Optimization*, 70 (2021), 149-168.
- [3] J. Bai, W. Hager and H. Zhang, Accelerated stochastic ADMM for separable convex optimization, *Comput. Optim. Appl.*, (2022), DOI:10.1007/s10589-021-00338-8.
- [4] J. Bai, Y. Ma, H. Sun and M. Zhang, Iteration complexity analysis of a partial LQP-based alternating direction method of multipliers, *Appl. Numer. Math.*, 165 (2021), 500-518.
- [5] X. Cai, D. Han and X. Yuan, On the convergence of the direct extension of ADMM for three-block separable convex minimization models with one strongly convex function, *Comput. Optim. Appl.*, 66 (2017), 39-73.
- [6] X. Chang, J. Bai, D. Song and S. Liu, Linearized symmetric multi-block ADMM with indefinite proximal regularization and optimal proximal parameter, *Calcolo*, 57 (2020), 1-36.
- [7] Y. Dai, D. Han, X. Yuan, and W. Zhang, A sequential updating scheme of the lagrange multiplier for separable convex programming, *Math. Comput.*, 86 (2017), 315-343.
- [8] J. Eckstein and D. Bertsekas, On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators, *Math. Program.*, 55 (1992), 293-318.
- [9] X. Fang, B. He, H. Liu and X. Yuan, Generalized alternating direction method of multipliers: New theoretical insights and applications, *Math. Prog. Comp.*, 7 (2015), 149-187.
- [10] J. Friedman, T. Hastie and R. Tibshirani, Sparse inverse covariance estimation with the graphical lasso, *Biostatistics*, 9 (2008), 432-441.
- [11] R. Glowinski and A. Marrocco, Approximation par éléments finis d'ordre un et résolution, par pénalisation-dualité d'une classe de problèmes de Dirichlet non linéaires, *Rev. Fr. Autom. Inform. Rech. Opér. Anal. Numér.*, 2 (1975), 41-76.
- [12] D. Gabay and B. Mercier, A dual algorithm for the solution of nonlinear variational problems via finite element approximations, *Comput. Math. Appl.*, 2 (1976), 17-40.
- [13] T. Goldstein, B. Donoghue, S. Setzer and R. Baraniuk, Fast alternating direction optimization methods, *SIAM J. Imaging Sci.*, 7 (2014), 1588-1623.

- [14] Y. Gu, B. Jiang and D. Han, A semi-proximal-based strictly contractive Peaceman-Rachford splitting method, arXiv:1506.02221, (2015), 1-20.
- [15] G. Gu, B. He and J. Yang, Inexact alternating-direction-based contraction methods for separable linearly constrained convex optimization, *J. Optim. Theory Appl.*, 163 (2014), 105-129.
- [16] W. Hager and H. Zhang, Inexact alternating direction multiplier methods for separable convex optimization, *Comput. Optim. Appl.*, 73 (2019), 201-235.
- [17] D. Han, A hybrid entropic proximal decomposition method with self-adaptive strategy for solving variational inequality problems, *Comput. Math. Appl.*, 55 (2008), 101-115.
- [18] B. He, H. Liu, Z. Wang and X. Yuan, A strictly contractive Peaceman-Rachford splitting method for convex programming, *SIAM J. Optim.*, 24 (2014), 1011-1040.
- [19] B. He, F. Ma and X. Yuan, Convergence study on the symmetric version of ADMM with larger step sizes, *SIAM J. Imaging Sci.*, 9 (2016), 1467-1501.
- [20] M. Hong and Z. Luo, On the linear convergence of alternating direction method of multipliers, *Math. Program.*, 162 (2017), 165-199.
- [21] F. Jiang, Z. Wu and X. Cai, Generalized ADMM with optimal indefinite proximal term for linearly constrained convex optimization, *J. Indust. Manag. Optim.*, 16 (2020), 835-856.
- [22] Z. Jia, X. Gao, X. Cai and D. Han, The convergence rate analysis of the symmetric ADMM for the nonconvex separable optimization problems, *J. Indust. Manag. Optim.*, 17 (2021), 1943-1971.
- [23] T. Lin, S. Ma and S. Zhang, On the global linear convergence of the ADMM with multi-block variables, *SIAM J. Optim.*, 25 (2015), 1478-1497.
- [24] Y. Liu, F. Shang, and J. Cheng, Accelerated variance reduced stochastic ADMM, in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI17*, AAAI Press, 2017, 2287-2293.
- [25] G. Luo and Q. Yang, A fast symmetric alternating direction method of multipliers, *Numer. Math. Theor. Meth. Appl.*, 13 (2020), 200-219.
- [26] H. Ouyang, N. He, L. Tran and A. Gray, Stochastic alternating direction method of multipliers, *Proc. 30th Int. Conf. Mach. Learn.*, (2013), 80-88.
- [27] D. Peaceman and H. Rachford JR., The numerical solution of parabolic and elliptic differential equations, *J. Soc. Indust. Appl. Math.*, 3 (1955), 28-41.
- [28] H. Sun, M. Sun and Y. Wang, Proximal ADMM with larger step size for two-block separable convex programming and its application to the correlation matrices calibrating problems, *J. Nonlinear Sci. Appl.*, 10 (2017), 5038-5051.
- [29] Z. Wu and M. Li, An LQP-based symmetric alternating direction method of multipliers with larger step sizes, *J. Oper. Res. Soc. China*, 7 (2019), 365-383.
- [30] Y. Xiao, L. Chen and D. Li, A generalized alternating direction method of multipliers with semi-proximal terms for convex composite conic programming, *Math. Prog. Comp.*, 10 (2018), 533-555.
- [31] M. Xu and T. Wu, A class of linearized proximal alternating direction methods, *J. Optim. Theory Appl.*, 151 (2011), 321-337.
- [32] Z. Yang and Z. Yan, Fast linearized alternating direction method of multipliers for the augmented l_1 -regularized problem, *SIViP*, 9 (2015), 1601-1612.
- [33] X. Yuan, S. Zeng and J. Zhang, Discerning the linear convergence of ADMM for sturcuted convex optimization through the lens of cariationsl analysis, *J. Mach. Learn. Res.*, 21 (2020), 1-74.
- [34] W. Yang and D. Han, Linear convergence of the alternating direction method of multipliers for a class of convex optimization problems, *SIAM J. Numer. Anal.*, 54 (2016), 625-640.

- [35] S. Zhao, W. Li and Z. Zhou, Scalable stochastic alternating direction method of multipliers, arXiv:1502.03529, (2015), 1-24.
- [36] N. Zhang, J. Wu, L. Zhang, A linearly convergent majorized ADMM with indefinite proximal terms for convex composite programming and its applications, Math. Comput., 89 (2020), 1867-1894.