

A NON-MONOTONE ALTERNATING UPDATING METHOD FOR A CLASS OF MATRIX FACTORIZATION PROBLEMS

LEI YANG*, TING KEI PONG*, AND XIAOJUN CHEN*

Abstract. In this paper we consider a general matrix factorization model which covers a large class of existing models with many applications in areas such as machine learning and imaging sciences. To solve this possibly nonconvex, nonsmooth and non-Lipschitz problem, we develop a non-monotone alternating updating method based on a potential function. Our method essentially updates two blocks of variables in turn by inexactly minimizing this potential function, and updates another auxiliary block of variables using an explicit formula. The special structure of our potential function allows us to take advantage of efficient computational strategies for non-negative matrix factorization to perform the alternating minimization over the two blocks of variables. A suitable line search criterion is also incorporated to improve the numerical performance. Under some mild conditions, we show that the line search criterion is well defined, and establish that the sequence generated is bounded and any cluster point of the sequence is a stationary point. Finally, we conduct some numerical experiments using real datasets to compare our method with some existing efficient methods for non-negative matrix factorization and matrix completion. The numerical results show that our method can outperform these methods for these specific applications.

Key words. Matrix factorization; non-monotone line search; stationary point; alternating updating.

AMS subject classifications. 90C26, 90C30, 90C90, 65K05

1. Introduction. In this paper we consider a class of matrix factorization problems, which can be modeled as

$$(1.1) \quad \min_{X,Y} \mathcal{F}(X, Y) := \Psi(X) + \Phi(Y) + \frac{1}{2} \|\mathcal{A}(XY^\top) - \mathbf{b}\|^2,$$

where $X \in \mathbb{R}^{m \times r}$ and $Y \in \mathbb{R}^{n \times r}$ are decision variables with $r \leq \min\{m, n\}$, the functions $\Psi : \mathbb{R}^{m \times r} \rightarrow \mathbb{R} \cup \{\infty\}$ and $\Phi : \mathbb{R}^{n \times r} \rightarrow \mathbb{R} \cup \{\infty\}$ are proper closed but possibly *nonconvex*, *nonsmooth* and *non-Lipschitz*, $\mathbf{b} \in \mathbb{R}^q$ is a given vector and $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^q$ is a linear map with $q \leq mn$ and $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ (\mathcal{I}_q denotes the identity map from \mathbb{R}^q to \mathbb{R}^q). Model (1.1) covers many existing widely-studied models in many application areas such as machine learning [35] and imaging sciences [44]. In particular, $\Psi(X)$ and $\Phi(Y)$ can be various regularizers for inducing desired structures, and \mathcal{A} can be suitably chosen to model different scenarios. For example, when $\Psi(X)$ and $\Phi(Y)$ are chosen as the indicator functions (see the next section for notation and definitions) for $\mathcal{X} = \{X \in \mathbb{R}^{m \times r} : X \geq 0\}$ and $\mathcal{Y} = \{Y \in \mathbb{R}^{n \times r} : Y \geq 0\}$, respectively, and \mathcal{A} is the identity map, (1.1) reduces to the non-negative matrix factorization (NMF) problem, which has been widely used in data mining applications to provide interpretable decompositions of data. NMF was first introduced by Paatero and Tapper [25], and then popularized by Lee and Seung [17]. The basic task of NMF is to find two nonnegative matrices $X \in \mathbb{R}_+^{m \times r}$ and $Y \in \mathbb{R}_+^{n \times r}$ such that $M \approx XY^\top$ for a given nonnegative data matrix $M \in \mathbb{R}_+^{m \times n}$. We refer readers to [2, 9, 10, 18, 37] for more information on NMF and its variants. Another example of (1.1) arises in recent models of the matrix completion (MC) problem (see [30, 31, 32]), where $\Psi(X)$

*Department of Applied Mathematics, The Hong Kong Polytechnic University, Hung Hom, Kowloon, Hong Kong, P.R. China. (lei.yang@connect.polyu.hk, tk.pong@polyu.edu.hk, xiaojun.chen@polyu.edu.hk). The second author's work was supported in part by Hong Kong Research Grants Council PolyU153085/16p. The third author's work was supported in part by Hong Kong Research Grants Council PolyU153000/15p.

42 and $\Phi(Y)$ are chosen as the Schatten- p_1 quasi-norm and the Schatten- p_2 quasi-norm
 43 for suitable $p_1, p_2 > 0$, respectively, and \mathcal{A} is the sampling map. The MC problem
 44 aims to recover an unknown low rank matrix from a sample of its entries and arises in
 45 various applications (see, for example, [3, 22, 27, 33]). Many widely-studied models for
 46 MC are based on nuclear-norm minimization [5, 6, 26], or, more generally, Schatten- p
 47 ($0 < p \leq 1$) (quasi-)norm minimization [16, 23, 42]. Recently, models based on
 48 low-rank matrix factorization such as (1.1) have become popular because singular
 49 value decompositions or eigenvalue decompositions of huge ($m \times n$) matrices are not
 50 required for solving these models (see, for example, [15, 30, 31, 32, 34, 38]). More
 51 examples of (1.1) can be found in recent surveys [35, 44].

52 Problem (1.1) is in general nonconvex (even when Ψ, Φ are convex) and NP-hard¹.
 53 Therefore, in this paper, we focus on finding a *stationary point* of the objective \mathcal{F} in
 54 (1.1). Note that \mathcal{F} involves two blocks of variables. This kind of structure has been
 55 widely studied in the literature; see, for example, [1, 4, 13, 14, 40, 41, 43]. One popular
 56 class of methods for tackling this kind of problems is the alternating direction method
 57 of multipliers (ADMM) (see, for example, [41, 43]), in which each iteration consists of
 58 an alternating minimization of an augmented Lagrangian function that involves X, Y
 59 and some auxiliary variables, followed by updates of the associated multipliers. Ho-
 60 wever, the conditions presented in [41, 43] that guarantee convergence of the ADMM
 61 are too restrictive. Moreover, updating the auxiliary variables and the multipliers
 62 can be expensive for large-scale problems. Another class of methods for (1.1) is
 63 the alternating-minimization-based (or block-coordinate-descent-type) methods (see
 64 [1, 4, 8, 11, 20, 21, 40]), which alternately (exactly or inexactly) minimizes $\mathcal{F}(X, Y)$
 65 over each block of variables and converges under some mild conditions. When \mathcal{A} is not
 66 the identity map, the majorization technique can be used to simplify the subproblems.
 67 Some representative algorithms of this class are proximal alternating linearized mini-
 68 mization (PALM) [4], hierarchical alternating least squares (HALS) (for NMF only;
 69 see [8, 11, 20, 21]) and block coordinate descent (BCD) [40]. Comparing with ADMM,
 70 it was reported in [40] that BCD outperforms ADMM in both CPU time and solution
 71 quality for NMF.

72 PALM, HALS and BCD are currently the state-of-the-art algorithms for solving
 73 problems of the form (1.1). In this paper, we develop a new iterative method for (1.1),
 74 which, according to our numerical experiments in Section 6, outperforms HALS and
 75 BCD for NMF, and PALM for MC. Our method is based on the following potential
 76 function (specifically constructed for \mathcal{F} in (1.1)):

$$77 \quad (1.2) \quad \Theta_{\alpha, \beta}(X, Y, Z) := \Psi(X) + \Phi(Y) + \frac{\alpha}{2} \|XY^\top - Z\|_F^2 + \frac{\beta}{2} \|\mathcal{A}(Z) - \mathbf{b}\|^2,$$

78 where α and β are real numbers. Instead of alternately (exactly or inexactly) minimi-
 79 zing $\mathcal{F}(X, Y)$ or the augmented Lagrangian function, our method alternately updates
 80 X and Y by inexactly minimizing $\Theta_{\alpha, \beta}(X, Y, Z)$ over X and Y , and then updates Z
 81 by an *explicit formula*. Note that the coupled variables XY^\top is now separated from
 82 \mathcal{A} in our potential function. Thus, one can readily take advantage of efficient compu-
 83 tational strategies for NMF, such as those used in HALS (see the ‘‘hierarchical-prox’’
 84 updating strategy in Section 4), for inexactly minimizing $\Theta_{\alpha, \beta}(X, Y, Z)$ over X or Y .
 85 Furthermore, our method can be implemented for NMF and MC without explicitly
 86 forming the huge ($m \times n$) matrix Z (see (6.3) and (6.5)) in each iteration. This signifi-

¹Problem (1.1) is NP-hard because it contains NMF as a special case, which is NP-hard in general [36].

87 cantly reduces the computational cost per iteration. Finally, a suitable non-monotone
 88 line search criterion, which is motivated by recent studies on non-monotone algo-
 89 rithms (see, for example, [7, 12, 39]), is also incorporated to improve the numerical
 90 performance.

91 In the rest of this paper, we first present notation and preliminaries in Section 2.
 92 We then study the properties of our potential function $\Theta_{\alpha,\beta}$ in Section 3. Specifically,
 93 if $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and α, β are chosen such that $\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A} \succ 0$ and $\frac{1}{\alpha} + \frac{1}{\beta} = 1$, then the
 94 problem $\min_{X,Y,Z} \{\Theta_{\alpha,\beta}(X, Y, Z)\}$ is equivalent to (1.1) (see Theorem 3.2). Furthermore,
 95 under the weaker conditions that $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and $\frac{1}{\alpha} + \frac{1}{\beta} = 1$, we can show that (i) a
 96 stationary point of $\Theta_{\alpha,\beta}$ gives a stationary point of \mathcal{F} ; (ii) a stationary point of \mathcal{F} can
 97 be used to construct a stationary point of $\Theta_{\alpha,\beta}$ (see Theorem 3.3). Thus, one can find
 98 a stationary point of \mathcal{F} by finding a stationary point of $\Theta_{\alpha,\beta}$. In Section 4, we develop
 99 a non-monotone alternating updating method to find a stationary point of $\Theta_{\alpha,\beta}$, and
 100 hence of \mathcal{F} . The convergence analysis of our method is presented in Section 5. We
 101 show that our non-monotone line search criterion is well defined and any cluster point
 102 of the sequence generated by our method is a stationary point of \mathcal{F} under some mild
 103 conditions. Section 6 gives numerical experiments to evaluate the performance of our
 104 method for NMF and MC on real datasets. Our computational results illustrate the
 105 efficiency of our method. Finally, some concluding remarks are given in Section 7.

106 **2. Notation and preliminaries.** In this paper, for a vector $\mathbf{x} \in \mathbb{R}^m$, x_i de-
 107 notes its i -th entry, $\|\mathbf{x}\|$ denotes the Euclidean norm of \mathbf{x} and $\text{Diag}(\mathbf{x})$ denotes the
 108 diagonal matrix whose i -th diagonal element is x_i . For a matrix $X \in \mathbb{R}^{m \times n}$, x_{ij}
 109 denotes the ij -th entry of X , \mathbf{x}_j denotes the j -th column of X and $\text{tr}(X)$ deno-
 110 tes the trace of X . The Schatten- p (quasi-)norm ($0 < p < \infty$) of X is defined as
 111 $\|X\|_{S_p} = \left(\sum_{i=1}^{\min(m,n)} \varsigma_i^p(X) \right)^{\frac{1}{p}}$, where $\varsigma_i(X)$ is the i -th singular value of X . For
 112 $p = 2$, the Schatten-2 norm reduces to the Frobenius norm $\|X\|_F$, and for $p = 1$,
 113 the Schatten-1 norm reduces to the nuclear norm $\|X\|_*$. Moreover, the spectral
 114 norm is denoted by $\|X\|$, which is the largest singular value of X ; and the ℓ_1 -norm
 115 and ℓ_p -quasi-norm ($0 < p < 1$) of X are given by $\|X\|_1 := \sum_{i=1}^m \sum_{j=1}^n |x_{ij}|$ and
 116 $\|X\|_p := \left(\sum_{i=1}^m \sum_{j=1}^n |x_{ij}|^p \right)^{\frac{1}{p}}$, respectively. For two matrices X and Y of the same
 117 size, we denote their trace inner product by $\langle X, Y \rangle := \sum_{i=1}^m \sum_{j=1}^n x_{ij}y_{ij}$. We also
 118 use $X \leq Y$ (resp., $X \geq Y$) to denote $x_{ij} \leq y_{ij}$ (resp., $x_{ij} \geq y_{ij}$) for all (i, j) . Furt-
 119 hermore, for a linear map $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^q$, \mathcal{A}^* denotes the adjoint linear map and
 120 $\|\mathcal{A}\|$ denotes the induced operator norm of \mathcal{A} , i.e., $\|\mathcal{A}\| = \sup\{\|\mathcal{A}(X)\| : \|X\|_F \leq 1\}$.
 121 A linear self-map \mathcal{T} is said to be symmetric if $\mathcal{T} = \mathcal{T}^*$. For a symmetric linear self-
 122 map $\mathcal{T} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$, we say that \mathcal{T} is positive definite, denoted by $\mathcal{T} \succ 0$, if
 123 $\langle X, \mathcal{T}(X) \rangle > 0$ for all $X \neq 0$. The identity map from $\mathbb{R}^{m \times n}$ to $\mathbb{R}^{m \times n}$ is denoted by \mathcal{I}
 124 and the identity map from \mathbb{R}^q to \mathbb{R}^q is denoted by \mathcal{I}_q . Finally, for a nonempty closed
 125 set $\mathcal{C} \subseteq \mathbb{R}^{m \times n}$, its indicator function $\delta_{\mathcal{C}}$ is defined by

$$126 \quad \delta_{\mathcal{C}}(X) = \begin{cases} 0 & \text{if } X \in \mathcal{C}, \\ +\infty & \text{otherwise.} \end{cases}$$

127 For an extended-real-valued function $f : \mathbb{R}^{m \times n} \rightarrow [-\infty, \infty]$, we say that it is
 128 *proper* if $f(X) > -\infty$ for all $X \in \mathbb{R}^{m \times n}$ and its domain $\text{dom} f := \{X \in \mathbb{R}^{m \times n} :$
 129 $f(X) < \infty\}$ is nonempty. A function $f : \mathbb{R}^{m \times n} \rightarrow [-\infty, \infty]$ is level-bounded [28,
 130 Definition 1.8] if for every $\alpha \in \mathbb{R}$, the set $\{X \in \mathbb{R}^{m \times n} : f(X) \leq \alpha\}$ is bounded
 131 (possibly empty). For a proper function $f : \mathbb{R}^{m \times n} \rightarrow (-\infty, \infty]$, we use the notation

132 $Y \xrightarrow{f} X$ to denote $Y \rightarrow X$ (i.e., $\|Y - X\|_F \rightarrow 0$) and $f(Y) \rightarrow f(X)$. The (limiting)
 133 subdifferential [28, Definition 8.3] of f at $X \in \text{dom} f$ used in this paper, denoted by
 134 $\partial f(X)$, is defined as

$$135 \quad \partial f(X) := \left\{ D \in \mathbb{R}^{m \times n} : \exists X^k \xrightarrow{f} X \text{ and } D^k \rightarrow D \text{ with } D^k \in \widehat{\partial} f(X^k) \text{ for all } k \right\},$$

136 where $\widehat{\partial} f(\widetilde{Y})$ denotes the Fréchet subdifferential of f at $\widetilde{Y} \in \text{dom} f$, which is the set
 137 of all $D \in \mathbb{R}^{m \times n}$ satisfying

$$138 \quad \liminf_{Y \neq \widetilde{Y}, Y \rightarrow \widetilde{Y}} \frac{f(Y) - f(\widetilde{Y}) - \langle D, Y - \widetilde{Y} \rangle}{\|Y - \widetilde{Y}\|_F} \geq 0.$$

139 From the above definition, we can easily observe (see, for example, [28, Proposi-
 140 tion 8.7]) that

$$141 \quad (2.1) \quad \left\{ D \in \mathbb{R}^{m \times n} : \exists X^k \xrightarrow{f} X, D^k \rightarrow D, D^k \in \partial f(X^k) \right\} \subseteq \partial f(X).$$

142 When f is continuously differentiable or convex, the above subdifferential coincides
 143 with the classical concept of derivative or convex subdifferential of f ; see, for example,
 144 [28, Exercise 8.8] and [28, Proposition 8.12]. In this paper, we say that X^* is *stationary*
 145 *point* of f if $0 \in \partial f(X^*)$.

146 For a proper closed function $g : \mathbb{R}^m \rightarrow (-\infty, \infty]$, the proximal mapping $\text{Prox}_g : \mathbb{R}^m \rightarrow \mathbb{R}^m$
 147 of g is defined by $\text{Prox}_g(\mathbf{z}) := \underset{\mathbf{x} \in \mathbb{R}^m}{\text{Argmin}} \{g(\mathbf{x}) + \frac{1}{2}\|\mathbf{x} - \mathbf{z}\|^2\}$. For any $\nu > 0$,
 148 the matrix shrinkage operator $\mathcal{S}_\nu : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$ is defined by

$$149 \quad \mathcal{S}_\nu(X) := U \text{Diag}(\bar{\mathbf{s}}) V^\top \text{ with } \bar{s}_i = \begin{cases} s_i - \nu, & \text{if } s_i - \nu > 0, \\ 0, & \text{otherwise,} \end{cases}$$

150 where $U \in \mathbb{R}^{m \times t}$, $\mathbf{s} \in \mathbb{R}_+^t$ and $V \in \mathbb{R}^{n \times t}$ are given by the singular value decomposition
 151 of X , i.e., $X = U \text{Diag}(\mathbf{s}) V^\top$.

152 We now present two propositions, which will be useful for developing our method
 153 in Section 4.

154 **PROPOSITION 2.1.** *Suppose that $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and $\alpha(\alpha + \beta) \neq 0$. Then, $\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A}$
 155 is invertible and its inverse is given by $\frac{1}{\alpha}\mathcal{I} - \frac{\beta}{\alpha(\alpha + \beta)}\mathcal{A}^*\mathcal{A}$.*

156 *Proof.* It is easy to check that $\frac{1}{\alpha}\mathcal{I} - \frac{\beta}{\alpha(\alpha + \beta)}\mathcal{A}^*\mathcal{A}$ is well defined since $\alpha(\alpha + \beta) \neq 0$,
 157 and that $(\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A}) \left(\frac{1}{\alpha}\mathcal{I} - \frac{\beta}{\alpha(\alpha + \beta)}\mathcal{A}^*\mathcal{A} \right) = \mathcal{I}$. This completes the proof. \square

158 **PROPOSITION 2.2.** *Let $\psi : \mathbb{R}^m \rightarrow (-\infty, \infty]$ and $\phi : \mathbb{R}^n \rightarrow (-\infty, \infty]$ be proper
 159 closed functions. Given $P, Q \in \mathbb{R}^{m \times n}$ and $\mathbf{a} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^m$ with $\|\mathbf{a}\| \neq 0$, $\|\mathbf{b}\| \neq 0$,
 160 the following statements hold.*

(i) *The problem $\min_{\mathbf{x} \in \mathbb{R}^m} \{ \psi(\mathbf{x}) + \frac{1}{2}\|\mathbf{x}\mathbf{a}^\top - P\|_F^2 \}$ is equivalent to*

$$\min_{\mathbf{x} \in \mathbb{R}^m} \left\{ \psi(\mathbf{x}) + \frac{\|\mathbf{a}\|^2}{2} \left\| \mathbf{x} - \frac{P\mathbf{a}}{\|\mathbf{a}\|^2} \right\|^2 \right\};$$

(ii) *The problem $\min_{\mathbf{y} \in \mathbb{R}^n} \{ \phi(\mathbf{y}) + \frac{1}{2}\|\mathbf{b}\mathbf{y}^\top - Q\|_F^2 \}$ is equivalent to*

$$\min_{\mathbf{y} \in \mathbb{R}^n} \left\{ \phi(\mathbf{y}) + \frac{\|\mathbf{b}\|^2}{2} \left\| \mathbf{y} - \frac{Q^\top \mathbf{b}}{\|\mathbf{b}\|^2} \right\|^2 \right\}.$$

161 *Proof.* Statement (i) can be easily proved by noticing that

$$162 \quad \begin{aligned} \|\mathbf{x}\mathbf{a}^\top - P\|_F^2 &= \|\mathbf{x}\mathbf{a}^\top\|_F^2 - 2\langle \mathbf{x}\mathbf{a}^\top, P \rangle + \|P\|_F^2 = \|\mathbf{a}\|^2\|\mathbf{x}\|^2 - 2\langle \mathbf{x}, P\mathbf{a} \rangle + \|P\|_F^2 \\ &= \|\mathbf{a}\|^2 \|\mathbf{x} - P\mathbf{a}/\|\mathbf{a}\|^2\|^2 - \|P\mathbf{a}\|^2/\|\mathbf{a}\|^2 + \|P\|_F^2. \end{aligned}$$

163 Then, statement (ii) can be easily proved by using statement (i) and $\|\mathbf{b}\mathbf{y}^\top - Q\|_F^2 =$
164 $\|\mathbf{y}\mathbf{b}^\top - Q^\top\|_F^2$. \square

165 Before ending this section, we discuss the first-order necessary conditions for (1.1).
166 First, from [28, Exercise 8.8] and [28, Proposition 10.5], we see that

$$167 \quad \partial\mathcal{F}(X, Y) = \begin{pmatrix} \partial\Psi(X) + \mathcal{A}^*(\mathcal{A}(XY^\top) - \mathbf{b})Y \\ \partial\Phi(Y) + (\mathcal{A}^*(\mathcal{A}(XY^\top) - \mathbf{b}))^\top X \end{pmatrix}.$$

168 Then, it follows from the generalized Fermat's rule [28, Theorem 10.1] that any local
169 minimizer (\bar{X}, \bar{Y}) of (1.1) satisfies $0 \in \partial\mathcal{F}(\bar{X}, \bar{Y})$, i.e.,

$$170 \quad (2.2) \quad \begin{cases} 0 \in \partial\Psi(\bar{X}) + \mathcal{A}^*(\mathcal{A}(\bar{X}\bar{Y}^\top) - \mathbf{b})\bar{Y}, \\ 0 \in \partial\Phi(\bar{Y}) + (\mathcal{A}^*(\mathcal{A}(\bar{X}\bar{Y}^\top) - \mathbf{b}))^\top \bar{X}, \end{cases}$$

171 which implies that (\bar{X}, \bar{Y}) is a stationary point of \mathcal{F} . In this paper, we focus on finding
172 a stationary point (X^*, Y^*) of \mathcal{F} , i.e., (X^*, Y^*) satisfies (2.2) in place of (\bar{X}, \bar{Y}) .

173 **3. The potential function for \mathcal{F} .** In this section, we analyze the relation
174 between \mathcal{F} and its potential function $\Theta_{\alpha, \beta}$ defined in (1.2). Intuitively, $\Theta_{\alpha, \beta}$ originates
175 from \mathcal{F} by separating the coupled variables XY^\top from the linear mapping \mathcal{A} via
176 introducing an auxiliary variable Z and penalizing $XY^\top = Z$. We will see later that
177 the stationary point of \mathcal{F} can be characterized by the stationary point of $\Theta_{\alpha, \beta}$. Before
178 proceeding, we prove the following technical lemma.

179 **LEMMA 3.1.** *Suppose that $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and $\frac{1}{\alpha} + \frac{1}{\beta} = 1$. Then, for any (X, Y, Z)*
180 *satisfying*

$$181 \quad (3.1) \quad Z = \left(\mathcal{I} - \frac{\beta}{\alpha+\beta}\mathcal{A}^*\mathcal{A}\right)(XY^\top) + \frac{\beta}{\alpha+\beta}\mathcal{A}^*(\mathbf{b}),$$

182 we have $\mathcal{F}(X, Y) = \Theta_{\alpha, \beta}(X, Y, Z)$.

183 *Proof.* First, from (3.1), we have

$$184 \quad (3.2) \quad XY^\top - Z = \frac{\beta}{\alpha+\beta}\mathcal{A}^*(\mathcal{A}(XY^\top) - \mathbf{b})$$

$$185 \quad \begin{aligned} \mathcal{A}(Z) - \mathbf{b} &= \mathcal{A}\left(XY^\top - \frac{\beta}{\alpha+\beta}\mathcal{A}^*\mathcal{A}(XY^\top) + \frac{\beta}{\alpha+\beta}\mathcal{A}^*(\mathbf{b})\right) - \mathbf{b} \\ 186 \quad (3.3) \quad &= \mathcal{A}(XY^\top) - \frac{\beta}{\alpha+\beta}\mathcal{A}\mathcal{A}^*\mathcal{A}(XY^\top) + \frac{\beta}{\alpha+\beta}\mathcal{A}\mathcal{A}^*(\mathbf{b}) - \mathbf{b} = \frac{\alpha}{\alpha+\beta}(\mathcal{A}(XY^\top) - \mathbf{b}), \end{aligned}$$

187 where the last equality follows from $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$. Then, we see that

$$188 \quad \begin{aligned} &\frac{\alpha}{2}\|XY^\top - Z\|_F^2 + \frac{\beta}{2}\|\mathcal{A}(Z) - \mathbf{b}\|^2 \\ &= \frac{\alpha}{2}\left\|\frac{\beta}{\alpha+\beta}\mathcal{A}^*(\mathcal{A}(XY^\top) - \mathbf{b})\right\|_F^2 + \frac{\beta}{2}\left\|\frac{\alpha}{\alpha+\beta}(\mathcal{A}(XY^\top) - \mathbf{b})\right\|^2 \\ &= \frac{\alpha\beta^2}{(\alpha+\beta)^2} \cdot \frac{1}{2}\|\mathcal{A}^*(\mathcal{A}(XY^\top) - \mathbf{b})\|_F^2 + \frac{\alpha^2\beta}{(\alpha+\beta)^2} \cdot \frac{1}{2}\|\mathcal{A}(XY^\top) - \mathbf{b}\|^2 \\ &= \frac{\alpha\beta^2}{(\alpha+\beta)^2} \cdot \frac{1}{2}\|\mathcal{A}(XY^\top) - \mathbf{b}\|^2 + \frac{\alpha^2\beta}{(\alpha+\beta)^2} \cdot \frac{1}{2}\|\mathcal{A}(XY^\top) - \mathbf{b}\|^2 \\ &= \frac{\alpha\beta}{\alpha+\beta} \cdot \frac{1}{2}\|\mathcal{A}(XY^\top) - \mathbf{b}\|^2, \end{aligned}$$

189 where the first equality follows from (3.2) and (3.3); and the third equality follows
 190 from $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$. This, together with $\frac{1}{\alpha} + \frac{1}{\beta} = 1$ and the definitions of \mathcal{F} and $\Theta_{\alpha,\beta}$
 191 completes the proof. \square

192 Based on the above lemma, we now establish the following property of $\Theta_{\alpha,\beta}$.

193 **THEOREM 3.2.** *Suppose that $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$. If α and β are chosen such that $\alpha\mathcal{I} +$
 194 $\beta\mathcal{A}^*\mathcal{A} \succ 0$ and $\frac{1}{\alpha} + \frac{1}{\beta} = 1$, then the problem $\min_{X,Y,Z} \{\Theta_{\alpha,\beta}(X, Y, Z)\}$ is equivalent to (1.1).*

195 *Proof.* First, it is easy to see from $\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A} \succ 0$ that the function $Z \mapsto$
 196 $\Theta_{\alpha,\beta}(X, Y, Z)$ is strongly convex. Thus, for any fixed X and Y , the optimal solution
 197 Z^* to the problem $\min_Z \{\Theta_{\alpha,\beta}(X, Y, Z)\}$ exists and is unique, and can be obtained
 198 explicitly. Indeed, from the optimality condition, we have

$$199 \quad \alpha(Z^* - XY^\top) + \beta\mathcal{A}^*(\mathcal{A}(Z^*) - \mathbf{b}) = 0.$$

200 Then, since $\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A}$ is invertible (as $\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A} \succ 0$), we see that

$$\begin{aligned} 201 \quad Z^* &= (\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A})^{-1} [\alpha XY^\top + \beta\mathcal{A}^*(\mathbf{b})] \\ &= \left[\frac{1}{\alpha}\mathcal{I} - \frac{\beta}{\alpha(\alpha+\beta)}\mathcal{A}^*\mathcal{A} \right] [\alpha XY^\top + \beta\mathcal{A}^*(\mathbf{b})] \\ &= \left(\mathcal{I} - \frac{\beta}{\alpha+\beta}\mathcal{A}^*\mathcal{A} \right) (XY^\top) + \left[\frac{\beta}{\alpha}\mathcal{A}^*(\mathbf{b}) - \frac{\beta^2}{\alpha(\alpha+\beta)}\mathcal{A}^*\mathcal{A}\mathcal{A}^*(\mathbf{b}) \right] \\ &= \left(\mathcal{I} - \frac{\beta}{\alpha+\beta}\mathcal{A}^*\mathcal{A} \right) (XY^\top) + \left[\frac{\beta}{\alpha} - \frac{\beta^2}{\alpha(\alpha+\beta)} \right] \mathcal{A}^*(\mathbf{b}) \\ &= \left(\mathcal{I} - \frac{\beta}{\alpha+\beta}\mathcal{A}^*\mathcal{A} \right) (XY^\top) + \frac{\beta}{\alpha+\beta}\mathcal{A}^*(\mathbf{b}), \end{aligned}$$

202 where the second equality follows from Proposition 2.1 and the fourth equality fol-
 203 lows from $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$. This, together with Lemma 3.1, implies that $\mathcal{F}(X, Y) =$
 204 $\Theta_{\alpha,\beta}(X, Y, Z^*)$. Then, we have that

$$\begin{aligned} 205 \quad \min_{X,Y,Z} \{\Theta_{\alpha,\beta}(X, Y, Z)\} &= \min_{X,Y} \left\{ \min_Z \{\Theta_{\alpha,\beta}(X, Y, Z)\} \right\} = \min_{X,Y} \{\Theta_{\alpha,\beta}(X, Y, Z^*)\} \\ &= \min_{X,Y} \{\mathcal{F}(X, Y)\}. \end{aligned}$$

206 This completes the proof. \square

207 **REMARK 3.1.** *From the proof of Lemma 3.1, we see that if Φ and Ψ are the indi-*
 208 *cator functions of some nonempty closed sets, then $\mathcal{F}(X, Y) = \left(\frac{1}{\alpha} + \frac{1}{\beta}\right) \Theta_{\alpha,\beta}(X, Y, Z)$*
 209 *holds with the special choice of Z in (3.1) whenever $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and $\frac{1}{\alpha} + \frac{1}{\beta} > 0$. Thus,*
 210 *the result in Theorem 3.2 remains valid whenever $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and α, β are chosen*
 211 *such that $\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A} \succ 0$ and $\frac{1}{\alpha} + \frac{1}{\beta} > 0$.*

212 We see from Theorem 3.2 that (1.1) is equivalent to minimizing $\Theta_{\alpha,\beta}$ with some
 213 suitable choices of α and β . On the other hand, we can also characterize the relation
 214 between the stationary points of \mathcal{F} and $\Theta_{\alpha,\beta}$ under weaker conditions on α and β .

215 **THEOREM 3.3.** *Suppose that $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and α, β are chosen such that $\frac{1}{\alpha} + \frac{1}{\beta} = 1$.
 216 Then, the following statements hold.*

- 217 (i) *If (X^*, Y^*, Z^*) is a stationary point of $\Theta_{\alpha,\beta}$, then (X^*, Y^*) is a stationary point*
 218 *of \mathcal{F} ;*
 219 (ii) *If (X^*, Y^*) is a stationary point of \mathcal{F} , then (X^*, Y^*, Z^*) is a stationary point of*
 220 *$\Theta_{\alpha,\beta}$, where Z^* is given by*

$$221 \quad (3.4) \quad Z^* = \left(\mathcal{I} - \frac{\beta}{\alpha+\beta}\mathcal{A}^*\mathcal{A} \right) (X^*(Y^*)^\top) + \frac{\beta}{\alpha+\beta}\mathcal{A}^*(\mathbf{b}).$$

Proof. First, if (X^*, Y^*, Z^*) is a stationary point of $\Theta_{\alpha, \beta}$, then we have $0 \in \partial\Theta_{\alpha, \beta}(X^*, Y^*, Z^*)$, i.e.,

$$\begin{aligned} (3.5a) \quad & \left\{ \begin{array}{l} 0 \in \partial\Psi(X^*) + \alpha(X^*(Y^*)^\top - Z^*)Y^*, \\ 0 \in \partial\Phi(Y^*) + \alpha(X^*(Y^*)^\top - Z^*)^\top X^*, \\ 0 = \alpha(Z^* - X^*(Y^*)^\top) + \beta\mathcal{A}^*(\mathcal{A}(Z^*) - \mathbf{b}). \end{array} \right. \\ (3.5b) \quad & \\ (3.5c) \quad & \end{aligned}$$

222 Since $\frac{1}{\alpha} + \frac{1}{\beta} = 1$, we have $\alpha(\alpha + \beta) \neq 0$ and hence $\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A}$ is invertible from Lemma
223 2.1. Then, using the same arguments in the proof of Theorem 3.2, we see from (3.5c)
224 that (X^*, Y^*, Z^*) satisfies (3.4). Moreover, using (3.4) and the same arguments in
225 (3.2) and (3.3), we have

$$226 (3.6) \quad X^*(Y^*)^\top - Z^* = \frac{\beta}{\alpha + \beta} \mathcal{A}^*(\mathcal{A}(X^*(Y^*)^\top) - \mathbf{b}),$$

$$227 (3.7) \quad \mathcal{A}(Z^*) - \mathbf{b} = \frac{\alpha}{\alpha + \beta} (\mathcal{A}(X^*(Y^*)^\top) - \mathbf{b}).$$

228 Thus, substituting (3.6) into (3.5a) and (3.5b), we see that

$$229 (3.8) \quad \left\{ \begin{array}{l} 0 \in \partial\Psi(X^*) + \frac{\alpha\beta}{\alpha + \beta} \mathcal{A}^*(\mathcal{A}(X^*(Y^*)^\top) - \mathbf{b})Y^*, \\ 0 \in \partial\Phi(Y^*) + \frac{\alpha\beta}{\alpha + \beta} (\mathcal{A}^*(\mathcal{A}(X^*(Y^*)^\top) - \mathbf{b}))^\top X^*. \end{array} \right.$$

230 This together with $\frac{1}{\alpha} + \frac{1}{\beta} = 1$ implies (X^*, Y^*) is a stationary point of \mathcal{F} . This proves
231 statement (i).

232 We now prove statement (ii). First, if (X^*, Y^*) is a stationary point of \mathcal{F} , then
233 invoking $\frac{1}{\alpha} + \frac{1}{\beta} = 1$ and (2.2), we have (3.8). Next, we consider (X^*, Y^*, Z^*) with Z^*
234 given by (3.4). Then, (X^*, Y^*, Z^*) satisfies (3.6) and (3.7). Thus, substituting (3.6)
235 into (3.8), we obtain (3.5a) and (3.5b). Moreover, we have from (3.6) and (3.7) that

$$236 (3.9) \quad \begin{aligned} & \alpha(Z^* - X^*(Y^*)^\top) + \beta\mathcal{A}^*(\mathcal{A}(Z^*) - \mathbf{b}) \\ & = -\frac{\alpha\beta}{\alpha + \beta} \mathcal{A}^* ((\mathcal{A}(X^*(Y^*)^\top) - \mathbf{b}) + \beta\mathcal{A}^* \left(\frac{\alpha}{\alpha + \beta} (\mathcal{A}(X^*(Y^*)^\top) - \mathbf{b}) \right)) = 0. \end{aligned}$$

237 This together with (3.5a) and (3.5b) implies that (X^*, Y^*, Z^*) is a stationary point
238 of $\Theta_{\alpha, \beta}$. This proves statement (ii). \square

239 **REMARK 3.2.** *From the proof of Theorem 3.3, one can see that if $\partial\Psi$ and $\partial\Phi$ are*
240 *cones, Theorem 3.3 remains valid under the weaker conditions that $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and*
241 *$\frac{1}{\alpha} + \frac{1}{\beta} > 0$.*

242 From Theorem 3.3, we see that a stationary point of \mathcal{F} can be obtained from a
243 stationary point of $\Theta_{\alpha, \beta}$ with a suitable choice of α and β , i.e., $\frac{1}{\alpha} + \frac{1}{\beta} = 1$. Since the
244 linear map \mathcal{A} is no longer associated with the coupled variables XY^\top in $\Theta_{\alpha, \beta}$, finding
245 a stationary point of $\Theta_{\alpha, \beta}$ is conceivably easier. Thus, one can consider finding a
246 stationary point of $\Theta_{\alpha, \beta}$ in order to find a stationary point of \mathcal{F} . Note that some
247 existing alternating-minimization-based methods (see, for example, [1, 40]) can be
248 used to find a stationary point of $\Theta_{\alpha, \beta}$, and hence of \mathcal{F} , under the conditions that
249 $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and α, β are chosen so that $\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A} \succ 0$ and $\frac{1}{\alpha} + \frac{1}{\beta} = 1$. These
250 conditions further imply that $\alpha > 1$ and $\beta = \frac{\alpha}{\alpha - 1} > 1$. However, as we will see from
251 our numerical results in Section 6, finding a stationary point of $\Theta_{\alpha, \beta}$ with $\alpha > 1$
252 can be slow. In view of this, in the next section, we develop a new non-monotone
253 alternating updating method for finding a stationary of $\Theta_{\alpha, \beta}$ (and hence of \mathcal{F}) under
254 the weaker conditions that $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and $\frac{1}{\alpha} + \frac{1}{\beta} = 1$. This allows more flexibilities
255 in choosing α and β .

256 **4. Non-monotone alternating updating method.** In this section, we con-
 257 sider a non-monotone alternating updating method (NAUM) for finding a stationary
 258 point of $\Theta_{\alpha,\beta}$ with $\frac{1}{\alpha} + \frac{1}{\beta} = 1$. Compared to existing alternating-minimization-based
 259 methods [1, 40] applied to $\Theta_{\alpha,\beta}$, which update X, Y, Z by alternately solving sub-
 260 problems related to $\Theta_{\alpha,\beta}$, NAUM updates Z by an *explicit formula* (see (4.5)) and
 261 updates X, Y by solving subproblems related to $\Theta_{\alpha,\beta}$ in a Gauss-Seidel manner. Be-
 262 fore presenting the complete algorithm, we first comment on the updates of X and
 263 Y .

264 Let (X^k, Y^k) denote the value of (X, Y) after the $(k-1)$ th iteration, and let (U, V)
 265 denote the candidate for (X^{k+1}, Y^{k+1}) at the k -th iteration (we will set (X^{k+1}, Y^{k+1})
 266 to be (U, V) if a line search criterion is satisfied; more details can be found in Algorithm
 267 1). For notational simplicity, we also define

$$268 \quad \mathcal{H}_\alpha(X, Y, Z) := \frac{\alpha}{2} \|XY^\top - Z\|_F^2$$

269 for any (X, Y, Z) . Then, at the k -th iteration, we first compute Z^k by (4.5) and, in
 270 the line search loop, we compute U in one of the following 3 ways for a given $\mu_k > 0$:

271 • **Proximal**

$$272 \quad (4.1a) \quad U \in \underset{X}{\operatorname{Argmin}} \Psi(X) + \mathcal{H}_\alpha(X, Y^k, Z^k) + \frac{\mu_k}{2} \|X - X^k\|_F^2.$$

273 • **Prox-linear**

$$274 \quad (4.1b) \quad U \in \underset{X}{\operatorname{Argmin}} \Psi(X) + \langle \nabla_X \mathcal{H}_\alpha(X^k, Y^k, Z^k), X - X^k \rangle + \frac{\mu_k}{2} \|X - X^k\|_F^2.$$

275 • **Hierarchical-prox** If Ψ is column-wise separable, i.e., $\Psi(X) = \sum_{i=1}^r \psi_i(\mathbf{x}_i)$ for
 $X = [\mathbf{x}_1, \dots, \mathbf{x}_r] \in \mathbb{R}^{m \times r}$, we can update U column-by-column. Specifically, for
 $i = 1, 2, \dots, r$, compute

$$276 \quad (4.1c) \quad \mathbf{u}_i \in \underset{\mathbf{x}_i}{\operatorname{Argmin}} \psi_i(\mathbf{x}_i) + \mathcal{H}_\alpha(\mathbf{u}_{j<i}, \mathbf{x}_i, \mathbf{x}_{j>i}^k, Y^k, Z^k) + \frac{\mu_k}{2} \|\mathbf{x}_i - \mathbf{x}_i^k\|^2,$$

277 where $\mathbf{u}_{j<i}$ denotes $(\mathbf{u}_1, \dots, \mathbf{u}_{i-1})$ and $\mathbf{x}_{j>i}^k$ denotes $(\mathbf{x}_{i+1}^k, \dots, \mathbf{x}_r^k)$.

278 After computing U , we compute V in one of the following 3 ways for a given $\sigma_k > 0$:

279 • **Proximal**

$$280 \quad (4.2a) \quad V \in \underset{Y}{\operatorname{Argmin}} \Phi(Y) + \mathcal{H}_\alpha(U, Y, Z^k) + \frac{\sigma_k}{2} \|Y - Y^k\|_F^2.$$

281 • **Prox-linear**

$$282 \quad (4.2b) \quad V \in \underset{Y}{\operatorname{Argmin}} \Phi(Y) + \langle \nabla_Y \mathcal{H}_\alpha(U, Y^k, Z^k), Y - Y^k \rangle + \frac{\sigma_k}{2} \|Y - Y^k\|_F^2.$$

283 • **Hierarchical-prox** If Φ is column-wise separable, i.e., $\Phi(Y) = \sum_{i=1}^r \phi_i(\mathbf{y}_i)$ for
 $Y = [\mathbf{y}_1, \dots, \mathbf{y}_r] \in \mathbb{R}^{n \times r}$, we can update V column-by-column. Specifically, for
 $i = 1, 2, \dots, r$, compute

$$284 \quad (4.2c) \quad \mathbf{v}_i \in \underset{\mathbf{y}_i}{\operatorname{Argmin}} \phi_i(\mathbf{y}_i) + \mathcal{H}_\alpha(U, \mathbf{v}_{j<i}, \mathbf{y}_i, \mathbf{y}_{j>i}^k, Z^k) + \frac{\sigma_k}{2} \|\mathbf{y}_i - \mathbf{y}_i^k\|^2,$$

285 where $\mathbf{v}_{j<i}$ denotes $(\mathbf{v}_1, \dots, \mathbf{v}_{i-1})$ and $\mathbf{y}_{j>i}^k$ denotes $(\mathbf{y}_{i+1}^k, \dots, \mathbf{y}_r^k)$.

286 For notational simplicity, we further let

$$287 \quad (4.3) \quad \rho := \left\| \mathcal{I} - \frac{\beta}{\alpha + \beta} \mathcal{A}^* \mathcal{A} \right\|^2$$

288 and let $\gamma \geq 0$ be a nonnegative number satisfying

$$289 \quad (4.4) \quad (\alpha + \gamma) \mathcal{I} + \beta \mathcal{A}^* \mathcal{A} \succeq 0.$$

290 **REMARK 4.1 (Comments on “hierarchical-prox”).** *The hierarchical-prox up-*
 291 *dating scheme requires the column-wise separability of Ψ or Φ . This is satisfied for*
 292 *many common regularizers, for example, $\|\cdot\|_F^2$, $\|\cdot\|_1$, $\|\cdot\|_p^p$ ($0 < p < 1$), and the*
 293 *indicator function of the nonnegativity (or box) constraint.*

294 **REMARK 4.2 (Comments on ρ and γ).** *Since $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$, we see that the ei-*
 295 *genvalues of $\mathcal{A}^* \mathcal{A}$ are either 0 or 1. Then, the eigenvalues of $\mathcal{I} - \frac{\beta}{\alpha + \beta} \mathcal{A}^* \mathcal{A}$ must*
 296 *be either 1 or $\frac{\alpha}{\alpha + \beta}$, and hence $\rho = \max\{1, \alpha^2/(\alpha + \beta)^2\}$. Similarly, the eigenva-*
 297 *lues of $-(\alpha \mathcal{I} + \beta \mathcal{A}^* \mathcal{A})$ are either $-\alpha$ or $-(\alpha + \beta)$. Then, (4.4) is satisfied whenever*
 298 *$\gamma \geq \max\{0, -\alpha, -(\alpha + \beta)\}$.*

299 Now, we are ready to present NAUM as Algorithm 1.

Algorithm 1 NAUM for finding a stationary point of \mathcal{F}

Input: (X^0, Y^0) , α and β such that $\frac{1}{\alpha} + \frac{1}{\beta} = 1$, ρ as in (4.3), $\gamma \geq 0$ satisfying (4.4),
 $\tau > 1$, $c > 0$, $\mu^{\min} > 0$, $\sigma^{\max} > \sigma^{\min} > 0$, and an integer $N \geq 0$. Set $k = 0$.

while a termination criterion is not met, **do**

Step 1. Compute Z^k by

$$(4.5) \quad Z^k = \left(\mathcal{I} - \frac{\beta}{\alpha + \beta} \mathcal{A}^* \mathcal{A} \right) (X^k (Y^k)^\top) + \frac{\beta}{\alpha + \beta} \mathcal{A}^* (\mathbf{b}).$$

Step 2. Choose $\mu_k^0 \geq \mu^{\min}$ and $\sigma_k^0 \in [\sigma^{\min}, \sigma^{\max}]$ arbitrarily. Set $\tilde{\mu}_k = \mu_k^0$,
 $\sigma_k = \sigma_k^0$ and $\mu_k^{\max} = (\alpha + 2\gamma\rho) \|Y^k\|^2 + c$.

(2a) Set $\mu_k \leftarrow \min\{\tilde{\mu}_k, \mu_k^{\max}\}$. Compute U by either (4.1a), (4.1b) or
 (4.1c).

(2b) Compute V by either (4.2a), (4.2b) or (4.2c).

(2c) If

$$(4.6) \quad \mathcal{F}(U, V) - \max_{[k-N]_+ \leq i \leq k} \mathcal{F}(X^i, Y^i) \leq -\frac{c}{2} (\|U - X^k\|_F^2 + \|V - Y^k\|_F^2),$$

 then go to **Step 3**.

(2d) If $\mu_k = \mu_k^{\max}$, set $\sigma_k^{\max} = (\alpha + 2\gamma\rho) \|U\|^2 + c$, $\sigma_k \leftarrow \min\{\tau\sigma_k, \sigma_k^{\max}\}$
 and then, go to step **(2b)**; otherwise, set $\tilde{\mu}_k \leftarrow \tau\mu_k$ and $\sigma_k \leftarrow \tau\sigma_k$
 and then, go to step **(2a)**.

Step 3. Set $X^{k+1} \leftarrow U$, $Y^{k+1} \leftarrow V$, $\bar{\mu}_k \leftarrow \mu_k$, $\bar{\sigma}_k \leftarrow \sigma_k$, $k \leftarrow k + 1$ and go to

Step 1.

end while

Output: (X^k, Y^k)

300 In Algorithm 1, the update for Z^k is given explicitly. This is motivated by the
 301 condition on Z at a stationary point of $\Theta_{\alpha, \beta}$; see (3.5c). In fact, following the
 302 same arguments in (3.9), we see that (3.5c) always holds at (X^k, Y^k, Z^k) with Z^k

303 given in (4.5) when $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and $\frac{1}{\alpha} + \frac{1}{\beta} = 1$. If, in addition, $\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A} \succ$
 304 0 holds, one can show that Z^k is actually the optimal solution to the problem
 305 $\min_Z \{\Theta_{\alpha,\beta}(X^k, Y^k, Z)\}$. In this case, our NAUM with $N = 0$ in (4.6) can be viewed
 306 as an alternating-minimization-based method (see, for example, [1, 40]) applied to
 307 the problem $\min_{X,Y,Z} \{\Theta_{\alpha,\beta}(X, Y, Z)\}$. However, if $\alpha\mathcal{I} + \beta\mathcal{A}^*\mathcal{A} \not\succeq 0$,² then the corre-
 308 sponding $\inf_Z \{\Theta_{\alpha,\beta}(X^k, Y^k, Z)\} = -\infty$ for all k , and Z^k is only a stationary point
 309 of $Z \mapsto \Theta_{\alpha,\beta}(X^k, Y^k, Z)$. In this case, the function value of $\Theta_{\alpha,\beta}$ may increase after
 310 updating Z by (4.5). Fortunately, as we shall see later in (5.8) and (5.9), as long as
 311 $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and $\frac{1}{\alpha} + \frac{1}{\beta} = 1$, we still have $\Theta_{\alpha,\beta}(X^{k+1}, Y^{k+1}, Z^k) < \Theta_{\alpha,\beta}(X^k, Y^k, Z^k)$
 312 by updating X^{k+1} and Y^{k+1} with properly chosen parameters μ_k and σ_k . Thus, if
 313 the possible increase in $\Theta_{\alpha,\beta}$ induced by the Z -update is not too large, one can still
 314 ensure $\Theta_{\alpha,\beta}(X^{k+1}, Y^{k+1}, Z^{k+1}) < \Theta_{\alpha,\beta}(X^k, Y^k, Z^k)$. Moreover, it can be seen from
 315 Lemma 3.1 and (4.5) that $\mathcal{F}(X^k, Y^k) = \Theta_{\alpha,\beta}(X^k, Y^k, Z^k)$ and hence the decrease of
 316 $\Theta_{\alpha,\beta}$ translates to that of \mathcal{F} (see Lemma 5.1 below). In view of this, $\Theta_{\alpha,\beta}$ is a valid
 317 potential function for minimizing \mathcal{F} as long as $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$ and $\frac{1}{\alpha} + \frac{1}{\beta} = 1$, even when
 318 $\beta < 0$ or $\alpha < 0$. Allowing negative α or β makes our NAUM (even with $N = 0$ in
 319 (4.6)) different from the classical alternating minimization schemes.

320 Our NAUM also allows U and V to be updated in three different ways respectively,
 321 and hence there are 9 possible combinations. Thus, one can choose suitable updating
 322 schemes to fit different applications. In particular, if Ψ or Φ are column-wise separable,
 323 taking advantage of the structure of $\Theta_{\alpha,\beta}$ and the fact that XY^\top can be written as
 324 $\sum_{i=1}^r \mathbf{x}_i \mathbf{y}_i^\top$ with $X = [\mathbf{x}_1, \dots, \mathbf{x}_r] \in \mathbb{R}^{m \times r}$ and $Y = [\mathbf{y}_1, \dots, \mathbf{y}_r] \in \mathbb{R}^{n \times r}$, one can
 325 update X or Y column-wise even when $\mathcal{A} \neq \mathcal{I}$. The motivation for updating X
 326 (or Y) column-wise rather than updating the whole X (or Y) is that the resulting
 327 subproblems (4.1c) (or (4.2c)) can be reduced to the computation of the proximal
 328 mapping of ψ_i (or ϕ_i), which is easy for many commonly used ψ_i (or ϕ_i). Indeed,
 329 from (4.1c) and (4.2c), \mathbf{u}_i and \mathbf{v}_i are given by

$$330 \quad (4.7) \quad \begin{cases} \mathbf{u}_i \in \underset{\mathbf{x}_i}{\text{Argmin}} \left\{ \psi_i(\mathbf{x}_i) + \frac{\alpha}{2} \|\mathbf{x}_i(\mathbf{y}_i^k)^\top - P_i^k\|_F^2 + \frac{\mu_k}{2} \|\mathbf{x}_i - \mathbf{x}_i^k\|^2 \right\}, \\ \mathbf{v}_i \in \underset{\mathbf{y}_i}{\text{Argmin}} \left\{ \phi_i(\mathbf{y}_i) + \frac{\alpha}{2} \|\mathbf{u}_i \mathbf{y}_i^\top - Q_i^k\|_F^2 + \frac{\sigma_k}{2} \|\mathbf{y}_i - \mathbf{y}_i^k\|^2 \right\}, \end{cases}$$

331 where P_i^k and Q_i^k are defined by

$$332 \quad (4.8) \quad \begin{aligned} P_i^k &:= Z^k - \sum_{j=1}^{i-1} \mathbf{u}_j(\mathbf{y}_j^k)^\top - \sum_{j=i+1}^r \mathbf{x}_j^k(\mathbf{y}_j^k)^\top, \\ Q_i^k &:= Z^k - \sum_{j=1}^{i-1} \mathbf{u}_j \mathbf{v}_j^\top - \sum_{j=i+1}^r \mathbf{u}_j(\mathbf{y}_j^k)^\top. \end{aligned}$$

333 Then, from Proposition 2.2, we can reformulate the subproblems in (4.7) and obtain
 334 the corresponding solutions by computing the proximal mappings of ψ_i and ϕ_i , which
 335 can be computed efficiently when ψ_i and ϕ_i are some common regularizers used in
 336 the literature. In particular, when $\psi_i(\cdot)$ and $\phi_i(\cdot)$ are $\|\cdot\|_1$, $\|\cdot\|_2^2$ or the indicator
 337 function of the box constraint, these subproblems have closed-form solutions. This
 338 updating strategy has also been used for NMF; see, for example, [8, 20, 21]. However,
 339 the methods used in [8, 20, 21] can only be applied for some specific problems with
 340 $\mathcal{A} = \mathcal{I}$, while NAUM can be applied for more general problems with $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$.

341 Our NAUM adapts a non-monotone line search criterion (see Step 2 in Algorithm
 342 1) to improve the numerical performance. This is motivated by recent studies on

²This may happen when $0 < \alpha < 1$ so that $\beta = \alpha(\alpha - 1)^{-1} < 0$, or $0 < \beta < 1$ so that $\alpha = \beta(\beta - 1)^{-1} < 0$.

343 non-monotone algorithms with promising performances; see, for example, [7, 12, 39].
 344 However, different from the non-monotone line search criteria used there, NAUM only
 345 includes (U, V) in the line search loop and checks the stopping criterion (4.6) after
 346 updating a pair of (U, V) , rather than checking (4.6) immediately once U or V is
 347 updated. Thus, we do not need to compute the function value after updating each
 348 block of variable. This may reduce the cost of the line search and make NAUM more
 349 practical, especially when computing the function value is relatively expensive.

350 Before moving to the convergence analysis of NAUM, we would like to point out
 351 an interesting connection between NAUM and the low-rank matrix fitting algorithm,
 352 LMaFit [38], for solving the following matrix completion model without regularizers:

$$\min_{X, Y} \frac{1}{2} \|\mathcal{P}_\Omega(XY^\top - M)\|_F^2,$$

354 where Ω is the index set of the known entries of M , and $\mathcal{P}_\Omega(Z)$ keeps the entries of
 355 Z in Ω and sets the remaining ones to zero. If we apply our NAUM with (4.1a) and
 356 (4.2a), then at the k -th iteration, the iterates Z^k , X^{k+1} and Y^{k+1} are given by

$$\begin{aligned}
 Z^k &= \left(\mathcal{I} - \frac{\beta}{\alpha+\beta}\mathcal{P}_\Omega\right) X^k(Y^k)^\top + \frac{\beta}{\alpha+\beta}\mathcal{P}_\Omega(M), \\
 X^{k+1} &= (\bar{\mu}_k X^k + \alpha Z^k Y^k) (\bar{\mu}_k I + \alpha(Y^k)^\top Y^k)^{-1}, \\
 Y^{k+1} &= (\bar{\sigma}_k Y^k + \alpha(Z^k)^\top X^{k+1}) (\bar{\sigma}_k I + \alpha(X^{k+1})^\top X^{k+1})^{-1}.
 \end{aligned}$$

358 One can verify that the sequence $\{(Z^k, X^{k+1}, Y^{k+1})\}$ above can be equivalently ge-
 359 nerated by the following scheme with $\tilde{Z}^0 = \mathcal{P}_\Omega(M) + \mathcal{P}_{\Omega^c}(X^0(Y^0)^\top)$:

$$\begin{aligned}
 Z^k &= \frac{\beta}{\alpha+\beta}\tilde{Z}^k + \left(1 - \frac{\beta}{\alpha+\beta}\right) X^k(Y^k)^\top, \\
 X^{k+1} &= (\bar{\mu}_k X^k + \alpha Z^k Y^k) (\bar{\mu}_k I + \alpha(Y^k)^\top Y^k)^{-1}, \\
 Y^{k+1} &= (\bar{\sigma}_k Y^k + \alpha(Z^k)^\top X^{k+1}) (\bar{\sigma}_k I + \alpha(X^{k+1})^\top X^{k+1})^{-1}, \\
 \tilde{Z}^{k+1} &= \mathcal{P}_\Omega(M) + \mathcal{P}_{\Omega^c}(X^{k+1}(Y^{k+1})^\top),
 \end{aligned}$$

361 where Ω^c is the complement set of Ω . Surprisingly, when $\bar{\mu}_k = \bar{\sigma}_k = 0$, this scheme
 362 is exactly the SOR(successive over-relaxation)-like scheme used in LMaFit (see [38,
 363 Eq.(2.11)]) with $\omega := \frac{\beta}{\alpha+\beta}$ being an over-relaxation weight. With this connection,
 364 our NAUM, in some sense, can be viewed as an SOR-based algorithm. Moreover, just
 365 like the classical SOR for solving a system of linear equations, LMaFit with $\omega > 1$
 366 also appears to be more efficient from the extensive numerical experiments reported
 367 in [38]. Then, it is natural to consider $\frac{\beta}{\alpha+\beta} > 1$ and hence $\frac{1}{\alpha} > 1$ (since $\frac{1}{\alpha} + \frac{1}{\beta} = 1$) in
 368 NAUM. This also gives some insights for the necessity of allowing more flexibilities in
 369 choosing α and β , and the promising performance of NAUM with a relatively small
 370 $\alpha \in (0, 1)$ as we shall see in Section 6.

371 **5. Convergence analysis of NAUM.** In this section, we discuss the conver-
 372 gence properties of Algorithm 1. First, we present the first-order optimality conditions
 373 for the three different updating schemes in (2a) of Algorithm 1 as follows:

374 • **Proximal**

$$375 \quad (5.1a) \quad 0 \in \partial\Psi(U) + \alpha(U(Y^k)^\top - Z^k)Y^k + \mu_k(U - X^k).$$

376 • **Prox-linear**

$$377 \quad (5.1b) \quad 0 \in \partial\Psi(U) + \alpha(X^k(Y^k)^\top - Z^k)Y^k + \mu_k(U - X^k).$$

378 • **Hierarchical-prox** For $i = 1, 2, \dots, r$,

$$379 \quad (5.1c) \quad 0 \in \partial\psi_i(\mathbf{u}_i) + \alpha \left(\sum_{j=1}^i \mathbf{u}_j (\mathbf{y}_j^k)^\top + \sum_{j=i+1}^r \mathbf{x}_j^k (\mathbf{y}_j^k)^\top - Z^k \right) \mathbf{y}_i^k + \mu_k (\mathbf{u}_i - \mathbf{x}_i^k).$$

380 Similarly, the first-order optimality conditions for the three different updating schemes
381 in **(2b)** of Algorithm 1 are

382 • **Proximal**

$$383 \quad (5.2a) \quad 0 \in \partial\Phi(V) + \alpha (UV^\top - Z^k)^\top U + \sigma_k(V - Y^k).$$

384 • **Prox-linear**

$$385 \quad (5.2b) \quad 0 \in \partial\Phi(V) + \alpha (U(Y^k)^\top - Z^k)^\top U + \sigma_k(V - Y^k).$$

386 • **Hierarchical-prox** For $i = 1, 2, \dots, r$,

$$387 \quad (5.2c) \quad 0 \in \partial\phi_i(\mathbf{v}_i) + \alpha \left(\sum_{j=1}^i \mathbf{u}_j \mathbf{v}_j^\top + \sum_{j=i+1}^r \mathbf{u}_j (\mathbf{y}_j^k)^\top - Z^k \right)^\top \mathbf{u}_i + \sigma_k(\mathbf{v}_i - \mathbf{y}_i^k).$$

388 We also need to make the following assumptions.

389 ASSUMPTION 5.1.

390 **(a1)** Ψ, Φ are proper, closed, level-bounded functions and continuous on their domains respectively;

392 **(a2)** $\mathcal{A}\mathcal{A}^* = \mathcal{I}_q$;

393 **(a3)** $\frac{1}{\alpha} + \frac{1}{\beta} = 1$.

394 REMARK 5.1. (i) From **(a1)**, one can see from [28, Theorem 1.9] that $\inf \Psi$ and
395 $\inf \Phi$ are finite, i.e., Ψ and Φ are bounded from below. In particular, the iterates
396 (4.1a), (4.1b), (4.1c), (4.2a), (4.2b) and (4.2c) are well defined; (ii) The continuity
397 assumption in **(a1)** holds for many common regularizers, for example, ℓ_1 -norm, nuclear
398 norm and the indicator function of a nonempty closed set; (iii) **(a2)** is satisfied
399 for some common linear maps, for example, the identity map and the sampling map.

400 We start our convergence analysis by proving the following auxiliary lemma.

401 LEMMA 5.1 (**Sufficient descent of \mathcal{F}**). Suppose that Assumption 5.1 holds.
402 Let (X^k, Y^k) be generated by Algorithm 1 at the k -th iteration, and (U, V) be the
403 candidate for (X^{k+1}, Y^{k+1}) generated by steps **(2a)** and **(2b)**. Then, for any integer
404 $k \geq 0$, we have

$$405 \quad (5.3) \quad \begin{aligned} & \mathcal{F}(U, V) - \mathcal{F}(X^k, Y^k) \\ & \leq -\frac{\mu_k - (\alpha + 2\gamma\rho)\|Y^k\|^2}{2} \|U - X^k\|_F^2 - \frac{\sigma_k - (\alpha + 2\gamma\rho)\|U\|^2}{2} \|V - Y^k\|_F^2. \end{aligned}$$

406 *Proof.* First, from Lemma 3.1 and (4.5), we see that $\mathcal{F}(X^k, Y^k) = \Theta_{\alpha, \beta}(X^k, Y^k,$
407 $Z^k)$. For any (U, V) , let

$$408 \quad (5.4) \quad W = \left(\mathcal{I} - \frac{\beta}{\alpha + \beta} \mathcal{A}^* \mathcal{A} \right) (UV^\top) + \frac{\beta}{\alpha + \beta} \mathcal{A}^* (\mathbf{b}).$$

409 Then, from Lemma 3.1, we have $\mathcal{F}(U, V) = \Theta_{\alpha, \beta}(U, V, W)$. Thus, to establish (5.3),
410 we only need to consider the difference $\Theta_{\alpha, \beta}(U, V, W) - \Theta_{\alpha, \beta}(X^k, Y^k, Z^k)$.

411 We start by noting that

$$412 \quad (5.5) \quad \begin{aligned} \mathcal{A}^* \mathcal{A}(W) &= \left(\mathcal{A}^* \mathcal{A} - \frac{\beta}{\alpha + \beta} \mathcal{A}^* (\mathcal{A}\mathcal{A}^*) \mathcal{A} \right) (UV^\top) + \frac{\beta}{\alpha + \beta} \mathcal{A}^* (\mathcal{A}\mathcal{A}^*) (\mathbf{b}) \\ &= \frac{\alpha}{\alpha + \beta} \mathcal{A}^* \mathcal{A} (UV^\top) + \frac{\beta}{\alpha + \beta} \mathcal{A}^* (\mathbf{b}), \end{aligned}$$

413 where the last equality follows from **(a2)** in Assumption 5.1. Then, we obtain that

$$414 \quad \begin{aligned} \nabla_Z \Theta_{\alpha,\beta}(U, V, W) &= \alpha(W - UV^\top) + \beta \mathcal{A}^* \mathcal{A}(W) - \beta \mathcal{A}^*(\mathbf{b}) \\ &= \alpha \left[-\frac{\beta}{\alpha+\beta} \mathcal{A}^* \mathcal{A}(UV^\top) + \frac{\beta}{\alpha+\beta} \mathcal{A}^*(\mathbf{b}) \right] + \beta \left[\frac{\alpha}{\alpha+\beta} \mathcal{A}^* \mathcal{A}(UV^\top) + \frac{\beta}{\alpha+\beta} \mathcal{A}^*(\mathbf{b}) \right] - \beta \mathcal{A}^*(\mathbf{b}) = 0, \end{aligned}$$

415 where the second equality follows from (5.4) and (5.5). Moreover, since γ is chosen
416 such that $(\alpha + \gamma)\mathcal{I} + \beta \mathcal{A}^* \mathcal{A} \succeq 0$ (see (4.4)), we see that, for any $k \geq 0$, the function
417 $Z \mapsto \Theta_{\alpha,\beta}(U, V, Z) + \frac{\gamma}{2} \|Z - Z^k\|_F^2$ is convex and hence

$$418 \quad \begin{aligned} &\Theta_{\alpha,\beta}(U, V, Z^k) + \underbrace{\frac{\gamma}{2} \|Z^k - Z^k\|_F^2}_{=0} \\ &\geq \Theta_{\alpha,\beta}(U, V, W) + \frac{\gamma}{2} \|W - Z^k\|_F^2 + \underbrace{\langle \nabla_Z \Theta_{\alpha,\beta}(U, V, W), W - Z^k \rangle}_{=0} + \gamma(W - Z^k), Z^k - W, \end{aligned}$$

419 which implies that

$$420 \quad (5.6) \quad \Theta_{\alpha,\beta}(U, V, W) - \Theta_{\alpha,\beta}(U, V, Z^k) \leq \frac{\gamma}{2} \|W - Z^k\|_F^2.$$

421 Then, substituting (4.5) and (5.4) into (5.6), we obtain

$$422 \quad (5.7) \quad \begin{aligned} \Theta_{\alpha,\beta}(U, V, W) - \Theta_{\alpha,\beta}(U, V, Z^k) &\leq \frac{\gamma}{2} \left\| \left(\mathcal{I} - \frac{\beta}{\alpha+\beta} \mathcal{A}^* \mathcal{A} \right) (UV^\top - X^k(Y^k)^\top) \right\|_F^2 \\ &\leq \frac{\gamma}{2} \left\| \mathcal{I} - \frac{\beta}{\alpha+\beta} \mathcal{A}^* \mathcal{A} \right\|^2 \cdot \|UV^\top - X^k(Y^k)^\top\|_F^2 \\ &= \frac{\gamma\rho}{2} \|U(V - Y^k)^\top + (U - X^k)(Y^k)^\top\|_F^2 \\ &\leq \frac{\gamma\rho}{2} \left(\|U(V - Y^k)^\top\|_F + \|(U - X^k)(Y^k)^\top\|_F \right)^2 \\ &\stackrel{(i)}{\leq} \frac{\gamma\rho}{2} \left(\|U\| \|V - Y^k\|_F + \|Y^k\| \|U - X^k\|_F \right)^2 \\ &\stackrel{(ii)}{\leq} \gamma\rho \left(\|U\|^2 \|V - Y^k\|_F^2 + \|Y^k\|^2 \|U - X^k\|_F^2 \right), \end{aligned}$$

423 where the equality follows from the definition of ρ in (4.3); (i) follows from the relation
424 $\|AB\|_F \leq \|A\| \|B\|_F$; and (ii) follows from the relation $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$.

425 Next, we claim that

$$426 \quad (5.8) \quad \Theta_{\alpha,\beta}(U, V, Z^k) - \Theta_{\alpha,\beta}(U, Y^k, Z^k) \leq \frac{\alpha\|U\|^2 - \sigma_k}{2} \|V - Y^k\|_F^2,$$

$$427 \quad (5.9) \quad \Theta_{\alpha,\beta}(U, Y^k, Z^k) - \Theta_{\alpha,\beta}(X^k, Y^k, Z^k) \leq \frac{\alpha\|Y^k\|^2 - \mu_k}{2} \|U - X^k\|_F^2.$$

428 Below, we will only prove (5.8). The proof for (5.9) can be done in a similar way.

429 To prove (5.8), we consider the following three cases.

430 • Proximal: In this case, we have

$$431 \quad \begin{aligned} &\Theta_{\alpha,\beta}(U, V, Z^k) - \Theta_{\alpha,\beta}(U, Y^k, Z^k) = \Phi(V) + \mathcal{H}_\alpha(U, V, Z^k) - \Phi(Y^k) - \mathcal{H}_\alpha(U, Y^k, Z^k) \\ &= \left[\Phi(V) + \mathcal{H}_\alpha(U, V, Z^k) + \frac{\sigma_k}{2} \|V - Y^k\|_F^2 \right] - \left[\Phi(Y^k) + \mathcal{H}_\alpha(U, Y^k, Z^k) \right] - \frac{\sigma_k}{2} \|V - Y^k\|_F^2 \\ &\leq -\frac{\sigma_k}{2} \|V - Y^k\|_F^2, \end{aligned}$$

432 where the inequality follows from the definition of V as a minimizer of (4.2a).

433 This implies (5.8).

434 • Prox-linear: In this case, we have

$$\begin{aligned}
& \Theta_{\alpha,\beta}(U, V, Z^k) - \Theta_{\alpha,\beta}(U, Y^k, Z^k) = \Phi(V) + \mathcal{H}_\alpha(U, V, Z^k) - \Phi(Y^k) - \mathcal{H}_\alpha(U, Y^k, Z^k) \\
& \leq \Phi(V) + \mathcal{H}_\alpha(U, Y^k, Z^k) + \langle \nabla_Y \mathcal{H}_\alpha(U, Y^k, Z^k), V - Y^k \rangle + \frac{\alpha \|U\|^2}{2} \|V - Y^k\|_F^2 \\
& \quad - \Phi(Y^k) - \mathcal{H}_\alpha(U, Y^k, Z^k) \\
& = \Phi(V) + \langle \nabla_Y \mathcal{H}_\alpha(U, Y^k, Z^k), V - Y^k \rangle + \frac{\sigma_k}{2} \|V - Y^k\|_F^2 - \Phi(Y^k) + \frac{\alpha \|U\|^2 - \sigma_k}{2} \|V - Y^k\|_F^2 \\
& \leq \frac{\alpha \|U\|^2 - \sigma_k}{2} \|V - Y^k\|_F^2,
\end{aligned}$$

436 where the first inequality follows from the fact that $Y \mapsto \nabla_Y \mathcal{H}_\alpha(X, Y, Z)$ is Lip-
437 schitz with modulus $\alpha \|X\|^2$ and the last inequality follows from the definition
438 of V as a minimizer of (4.2b).

439 • Hierarchical-prox: In this case, for any $1 \leq i \leq r$, we have

$$\begin{aligned}
& \Theta_{\alpha,\beta}(U, \mathbf{v}_{j < i}, \mathbf{v}_i, \mathbf{y}_{j > i}^k, Z^k) - \Theta_{\alpha,\beta}(U, \mathbf{v}_{j < i}, \mathbf{y}_i^k, \mathbf{y}_{j > i}^k, Z^k) \\
& = \phi_i(\mathbf{v}_i) + \mathcal{H}_\alpha(U, \mathbf{v}_{j < i}, \mathbf{v}_i, \mathbf{y}_{j > i}^k, Z^k) - \phi_i(\mathbf{y}_i^k) - \mathcal{H}_\alpha(U, \mathbf{v}_{j < i}, \mathbf{y}_i^k, \mathbf{y}_{j > i}^k, Z^k) \\
& = \left[\phi_i(\mathbf{v}_i) + \mathcal{H}_\alpha(U, \mathbf{v}_{j < i}, \mathbf{v}_i, \mathbf{y}_{j > i}^k, Z^k) + \frac{\sigma_k}{2} \|\mathbf{v}_i - \mathbf{y}_i^k\|^2 \right] - \frac{\sigma_k}{2} \|\mathbf{v}_i - \mathbf{y}_i^k\|^2 \\
& \quad - \left[\phi_i(\mathbf{y}_i^k) + \mathcal{H}_\alpha(U, \mathbf{v}_{j < i}, \mathbf{y}_i^k, \mathbf{y}_{j > i}^k, Z^k) \right] \\
& \leq -\frac{\sigma_k}{2} \|\mathbf{v}_i - \mathbf{y}_i^k\|^2,
\end{aligned}$$

441 where the inequality follows from the definition of \mathbf{v}_i as a minimizer of (4.2c).
442 Then, summing the above relation from $i = r$ to $i = 1$ and simplifying the
443 resulting inequality, we obtain (5.8).

444 The inequality (5.9) can be obtained via a similar argument.

445 Now, summing (5.7), (5.8) and (5.9), and using $\mathcal{F}(U, V) = \Theta_{\alpha,\beta}(U, V, W)$ and
446 $\mathcal{F}(X^k, Y^k) = \Theta_{\alpha,\beta}(X^k, Y^k, Z^k)$, we obtain (5.3). This completes the proof. \square

447 From Lemma 5.1, we see that the sufficient descent of $\mathcal{F}(X, Y)$ can be guaranteed
448 as long as μ_k and σ_k are sufficiently large. Thus, based on this lemma, we can show
449 in the following proposition that our non-monotone line search criterion in Algorithm
450 1 is well defined.

451 **PROPOSITION 5.2 (Well-definedness of the non-monotone line search cri-**
452 **terion).** *Suppose that Assumption 5.1 holds and Algorithm 1 is applied. Then,*
453 *for each $k \geq 0$, the line search criterion (4.6) is satisfied after finitely many inner*
454 *iterations.*

455 *Proof.* We prove this proposition by contradiction. Assume that there exists
456 a $k \geq 0$ such that the line search criterion (4.6) cannot be satisfied after finitely
457 many inner iterations. Note from (2a) and (2d) in Step 2 of Algorithm 1 that
458 $\mu_k \leq \mu_k^{\max} = (\alpha + 2\gamma\rho) \|Y^k\|^2 + c$ and hence $\mu_k = \mu_k^{\max}$ must be satisfied after finitely
459 many inner iterations. Let n_k denote the number of inner iterations when $\mu_k = \mu_k^{\max}$
460 is satisfied for the first time. If $\mu_k^0 \geq \mu_k^{\max}$, then $n_k = 1$; otherwise, we have

$$461 \quad \mu^{\min} \tau^{n_k - 2} \leq \mu_k^0 \tau^{n_k - 2} < \mu_k^{\max},$$

462 which implies that

$$463 \quad (5.10) \quad n_k \leq \left\lceil \frac{\log(\mu_k^{\max}) - \log(\mu^{\min})}{\log \tau} + 2 \right\rceil.$$

464 Then, from **(2d)** in Step 2 of Algorithm 1, we have $U \equiv U_{\mu_k^{\max}}$ and $\sigma_k^{\max} = (\alpha +$
 465 $2\gamma\rho)\|U_{\mu_k^{\max}}\|^2 + c$ after at most $n_k + 1$ inner iterations, where $U_{\mu_k^{\max}}$ is computed by
 466 (4.1a), (4.1b) or (4.1c) with $\mu_k = \mu_k^{\max}$. Moreover, we see that $\sigma_k = \sigma_k^{\max}$ must be
 467 satisfied after finitely many inner iterations. Similarly, let \hat{n}_k denote the number of
 468 inner iterations when $\sigma_k = \sigma_k^{\max}$ is satisfied for the *first* time. If $\sigma_k^0 > \sigma_k^{\max}$, then
 469 $\hat{n}_k = n_k$; if $\sigma_k^0 = \sigma_k^{\max}$, then $\hat{n}_k = 0$; otherwise, we have

$$470 \quad \sigma^{\min} \tau^{\hat{n}_k - 1} \leq \sigma_k^0 \tau^{\hat{n}_k - 1} < \sigma_k^{\max},$$

471 which implies that

$$472 \quad \hat{n}_k \leq \left\lceil \frac{\log(\sigma_k^{\max}) - \log(\sigma^{\min})}{\log \tau} + 1 \right\rceil.$$

473 Thus, after at most $\max\{n_k, \hat{n}_k\} + 1$ inner iterations, we must have $V \equiv V_{\sigma_k^{\max}}$, where
 474 $V_{\sigma_k^{\max}}$ is computed by (4.2a), (4.2b) or (4.2c) with $\sigma_k = \sigma_k^{\max}$. Therefore, after at
 475 most $\max\{n_k, \hat{n}_k\} + 1$ inner iterations, we have

$$\begin{aligned}
 & \mathcal{F}(U_{\mu_k^{\max}}, V_{\sigma_k^{\max}}) - \mathcal{F}(X^k, Y^k) \\
 476 \quad & \leq -\frac{\mu_k^{\max} - (\alpha + 2\gamma\rho)\|Y^k\|^2}{2} \|U_{\mu_k^{\max}} - X^k\|_F^2 - \frac{\sigma_k^{\max} - (\alpha + 2\gamma\rho)\|U_{\mu_k^{\max}}\|^2}{2} \|V_{\sigma_k^{\max}} - Y^k\|_F^2 \\
 & = -\frac{c}{2} (\|U_{\mu_k^{\max}} - X^k\|_F^2 + \|V_{\sigma_k^{\max}} - Y^k\|_F^2),
 \end{aligned}$$

477 where the inequality follows from (5.3) and the equality follows from $\mu_k^{\max} = (\alpha +$
 478 $2\gamma\rho)\|Y^k\|^2 + c$ and $\sigma_k^{\max} = (\alpha + 2\gamma\rho)\|U_{\mu_k^{\max}}\|^2 + c$. This together with

$$479 \quad \mathcal{F}(X^k, Y^k) \leq \max_{[k-N]_+ \leq i \leq k} \mathcal{F}(X^i, Y^i)$$

480 implies that (4.6) must be satisfied after at most $\max\{n_k, \hat{n}_k\} + 1$ inner iterations,
 481 which leads to a contradiction. \square

482 Now, we are ready to prove our main convergence result, which characterizes a
 483 cluster point of the sequence generated by Algorithm 1. Our proof of statement (ii)
 484 in the following theorem is similar to that of [39, Lemma 4]. However, the arguments
 485 involved are more intricate since we have two blocks of variables in our line search
 486 loop.

487 **THEOREM 5.3.** *Suppose that Assumption 5.1 holds. Let $\{(X^k, Y^k)\}$ be the se-*
 488 *quence generated by Algorithm 1. Then,*

- 489 (i) **(boundedness of sequence)** $\{(X^k, Y^k)\}$, $\{\bar{\mu}_k\}$ and $\{\bar{\sigma}_k\}$ are bounded;
- 490 (ii) **(diminishing successive changes)** $\lim_{k \rightarrow \infty} \|X^{k+1} - X^k\|_F + \|Y^{k+1} - Y^k\|_F = 0$;
- 491 (iii) **(global subsequential convergence)** any cluster point (X^*, Y^*) of $\{(X^k, Y^k)\}$
 492 is a stationary point of \mathcal{F} .

493 *Proof.* Statement (i). We first show that

$$494 \quad (5.11) \quad \mathcal{F}(X^k, Y^k) \leq \mathcal{F}(X^0, Y^0)$$

495 for all $k \geq 1$. We will prove it by induction. Indeed, for $k = 1$, it follows from
 496 Proposition 5.2 that

$$497 \quad \mathcal{F}(X^1, Y^1) - \mathcal{F}(X^0, Y^0) \leq -\frac{c}{2} (\|X^1 - X^0\|_F^2 + \|Y^1 - Y^0\|_F^2) \leq 0$$

498 is satisfied after finitely many inner iterations. Hence, (5.11) holds for $k = 1$. We now
 499 suppose that (5.11) holds for all $k \leq K$ for some integer $K \geq 1$. Then, we only need
 500 to show that (5.11) also holds for $k = K + 1$. For $k = K + 1$, we have

$$501 \quad \begin{aligned} \mathcal{F}(X^{K+1}, Y^{K+1}) - \mathcal{F}(X^0, Y^0) &\leq \mathcal{F}(X^{K+1}, Y^{K+1}) - \max_{[K-N]_+ \leq i \leq K} \mathcal{F}(X^i, Y^i) \\ &\leq -\frac{c}{2} (\|X^{K+1} - X^K\|_F^2 + \|Y^{K+1} - Y^K\|_F^2) \leq 0, \end{aligned}$$

502 where the first inequality follows from the induction hypothesis and the second ine-
 503 quality follows from (4.6). Hence, (5.11) holds for $k = K + 1$. This completes the
 504 induction. Then, from (5.11), we have that for any $k \geq 0$,

$$505 \quad \mathcal{F}(X^0, Y^0) \geq \mathcal{F}(X^k, Y^k) = \Psi(X^k) + \Phi(Y^k) + \frac{1}{2} \|\mathcal{A}(X^k(Y^k)^\top) - \mathbf{b}\|^2,$$

506 which, together with (a1) in Assumption 5.1, implies that the sequences $\{X^k\}$, $\{Y^k\}$
 507 and $\{\|\mathcal{A}(X^k(Y^k)^\top) - \mathbf{b}\|\}$ are bounded. Moreover, from Step 2 and Step 3 in Algo-
 508 rithm 1, it is easy to see $\bar{\mu}_k \leq \mu_k^{\max} = (\alpha + 2\gamma\rho)\|Y^k\|^2 + c$ for all k . Since $\{Y^k\}$ is
 509 bounded, the sequences $\{\mu_k^{\max}\}$ and $\{\bar{\mu}_k\}$ are bounded. Next, we prove the bounded-
 510 ness of $\{\bar{\sigma}_k\}$. Indeed, at the k -th iteration, there are three possibilities:

- 511 • $\bar{\mu}_k < \mu_k^{\max}$: In this case, we have $\bar{\sigma}_k \leq \sigma_k^0 \tau^{\tilde{n}_k} \leq \sigma^{\max} \tau^{\tilde{n}_k}$, where \tilde{n}_k denotes
 512 the number of inner iterations for the line search at the k -th iteration and $\tilde{n}_k \leq$
 513 $\max\left\{1, \left\lfloor \frac{\log(\mu_k^{\max}) - \log(\mu^{\min})}{\log \tau} + 2 \right\rfloor\right\}$ (see (5.10) and the discussions preceding it).
- 514 • $\bar{\mu}_k = \mu_k^{\max}$ and $\bar{\sigma}_k > \sigma_k^{\max}$: In this case, we have $\bar{\sigma}_k \leq \sigma_k^0 \tau^{\tilde{n}_k} \leq \sigma^{\max} \tau^{\tilde{n}_k}$, where
 515 $\tilde{n}_k \leq \max\left\{1, \left\lfloor \frac{\log(\mu_k^{\max}) - \log(\mu^{\min})}{\log \tau} + 2 \right\rfloor\right\}$.
- 516 • Otherwise, we have $\bar{\sigma}_k \leq \sigma_k^{\max} = (\alpha + 2\gamma\rho)\|X^{k+1}\|^2 + c$.

517 Note that $\{\tilde{n}_k\}$ is bounded as $\{\mu_k^{\max}\}$ is bounded. Thus, $\{\bar{\sigma}_k\}$ is bounded as the
 518 sequences $\{X^k\}$ and $\{\tilde{n}_k\}$ are bounded. This proves statement (i).

519 *Statement (ii).* We first claim that any cluster point of $\{(X^k, Y^k)\}$ is in $\text{dom}\mathcal{F}$.
 520 Since $\{(X^k, Y^k)\}$ is bounded from statement (i), there exists at least one cluster
 521 point. Suppose that (X^*, Y^*) is a cluster point of $\{(X^k, Y^k)\}$ and let $\{(X^{k_i}, Y^{k_i})\}$ be
 522 a convergent subsequence such that $\lim_{i \rightarrow \infty} (X^{k_i}, Y^{k_i}) = (X^*, Y^*)$. Then, from the lower
 523 semicontinuity of \mathcal{F} (since Ψ, Φ are closed by (a1) in Assumption 5.1) and (5.11), we
 524 have

$$525 \quad \mathcal{F}(X^*, Y^*) \leq \liminf_{i \rightarrow \infty} \mathcal{F}(X^{k_i}, Y^{k_i}) \leq \mathcal{F}(X^0, Y^0),$$

526 which implies that $\mathcal{F}(X^*, Y^*)$ is finite and hence $(X^*, Y^*) \in \text{dom}\mathcal{F}$.

527 For notational simplicity, from now on, we let $\Delta_{X^k} := X^{k+1} - X^k$, $\Delta_{Y^k} :=$
 528 $Y^{k+1} - Y^k$, $\Delta_{Z^k} := Z^{k+1} - Z^k$ and

$$529 \quad (5.12) \quad \ell(k) = \arg \max_i \{ \mathcal{F}(X^i, Y^i) : i = [k-N]_+, \dots, k \}.$$

530 Then, the line search criterion (4.6) can be rewritten as

$$531 \quad (5.13) \quad \mathcal{F}(X^{k+1}, Y^{k+1}) - \mathcal{F}(X^{\ell(k)}, Y^{\ell(k)}) \leq -\frac{c}{2} (\|\Delta_{X^k}\|_F^2 + \|\Delta_{Y^k}\|_F^2) \leq 0.$$

532 Observe that

$$\begin{aligned}
 & \mathcal{F}(X^{\ell(k+1)}, Y^{\ell(k+1)}) \\
 &= \max_{[k+1-N]_+ \leq i \leq k+1} \mathcal{F}(X^i, Y^i) = \max \left\{ \mathcal{F}(X^{k+1}, Y^{k+1}), \max_{[k+1-N]_+ \leq i \leq k} \mathcal{F}(X^i, Y^i) \right\} \\
 533 & \stackrel{(i)}{\leq} \max \left\{ \mathcal{F}(X^{\ell(k)}, Y^{\ell(k)}), \max_{[k+1-N]_+ \leq i \leq k} \mathcal{F}(X^i, Y^i) \right\} \\
 & \leq \max \left\{ \mathcal{F}(X^{\ell(k)}, Y^{\ell(k)}), \max_{[k-N]_+ \leq i \leq k} \mathcal{F}(X^i, Y^i) \right\} \\
 & \stackrel{(ii)}{=} \max \left\{ \mathcal{F}(X^{\ell(k)}, Y^{\ell(k)}), \mathcal{F}(X^{\ell(k)}, Y^{\ell(k)}) \right\} = \mathcal{F}(X^{\ell(k)}, Y^{\ell(k)}),
 \end{aligned}$$

534 where (i) follows from (5.13) and (ii) follows from (5.12). Therefore, the sequence
 535 $\{\mathcal{F}(X^{\ell(k)}, Y^{\ell(k)})\}$ is non-increasing. Since $\mathcal{F}(X^{\ell(k)}, Y^{\ell(k)})$ is also bounded from below
 536 (due to **(a1)** in Assumption 5.1), we conclude that there exists a number $\tilde{\mathcal{F}}$ such that

$$537 \quad (5.14) \quad \lim_{k \rightarrow \infty} \mathcal{F}(X^{\ell(k)}, Y^{\ell(k)}) = \tilde{\mathcal{F}}.$$

We next prove by induction that for all $j \geq 1$,

$$538 \quad (5.15a) \quad \left\{ \begin{array}{l} \lim_{k \rightarrow \infty} \Delta_{X^{\ell(k)-j}} = \lim_{k \rightarrow \infty} \Delta_{Y^{\ell(k)-j}} = 0, \\ 539 \quad (5.15b) \quad \lim_{k \rightarrow \infty} \mathcal{F}(X^{\ell(k)-j}, Y^{\ell(k)-j}) = \tilde{\mathcal{F}}. \end{array} \right.$$

538 We first prove (5.15a) and (5.15b) for $j = 1$. Applying (5.13) with k replaced by
 539 $\ell(k) - 1$, we obtain

$$540 \quad \mathcal{F}(X^{\ell(k)}, Y^{\ell(k)}) - \mathcal{F}(X^{\ell(\ell(k)-1)}, Y^{\ell(\ell(k)-1)}) \leq -\frac{c}{2} (\|\Delta_{X^{\ell(k)-1}}\|_F^2 + \|\Delta_{Y^{\ell(k)-1}}\|_F^2),$$

541 which, together with (5.14), implies that

$$542 \quad (5.16) \quad \lim_{k \rightarrow \infty} \Delta_{X^{\ell(k)-1}} = \lim_{k \rightarrow \infty} \Delta_{Y^{\ell(k)-1}} = 0.$$

543 Then, from (5.14) and (5.16), we have

$$\begin{aligned}
 544 \quad \tilde{\mathcal{F}} &= \lim_{k \rightarrow \infty} \mathcal{F}(X^{\ell(k)}, Y^{\ell(k)}) = \lim_{k \rightarrow \infty} \mathcal{F}(X^{\ell(k)-1} + \Delta_{X^{\ell(k)-1}}, Y^{\ell(k)-1} + \Delta_{Y^{\ell(k)-1}}) \\
 &= \lim_{k \rightarrow \infty} \mathcal{F}(X^{\ell(k)-1}, Y^{\ell(k)-1}),
 \end{aligned}$$

545 where the last equality follows because $\{(X^k, Y^k)\}$ is bounded, any cluster point of
 546 $\{(X^k, Y^k)\}$ is in $\text{dom}\mathcal{F}$ and \mathcal{F} is uniformly continuous on any compact subset of
 547 $\text{dom}\mathcal{F}$ under **(a1)** in Assumption 5.1. Thus, (5.15a) and (5.15b) hold for $j = 1$.

548 We next suppose that (5.15a) and (5.15b) hold for $j = J$ for some $J \geq 1$. It
 549 remains to show that they also hold for $j = J + 1$. Indeed, from (5.13) with k
 550 replaced by $\ell(k) - J - 1$ (here, without loss of generality, we assume that k is large
 551 enough such that $\ell(k) - J - 1$ is nonnegative), we have

$$552 \quad \mathcal{F}(X^{\ell(k)-J}, Y^{\ell(k)-J}) - \mathcal{F}(X^{\ell(\ell(k)-J-1)}, Y^{\ell(\ell(k)-J-1)}) \leq -\frac{c}{2} (\|\Delta_{X^{\ell(k)-J-1}}\|_F^2 + \|\Delta_{Y^{\ell(k)-J-1}}\|_F^2),$$

553 which implies that

$$554 \quad \|\Delta_{X^{\ell(k)-J-1}}\|_F^2 + \|\Delta_{Y^{\ell(k)-J-1}}\|_F^2 \leq \frac{2}{c} (\mathcal{F}(X^{\ell(\ell(k)-J-1)}, Y^{\ell(\ell(k)-J-1)}) - \mathcal{F}(X^{\ell(k)-J}, Y^{\ell(k)-J})).$$

555 This together with (5.14) and the induction hypothesis implies that

$$556 \quad \lim_{k \rightarrow \infty} \Delta_{X^{\ell(k)-(J+1)}} = \lim_{k \rightarrow \infty} \Delta_{Y^{\ell(k)-(J+1)}} = 0.$$

557 Thus, (5.15a) holds for $j = J + 1$. From this, we further have

$$558 \quad \begin{aligned} \lim_{k \rightarrow \infty} \mathcal{F}(X^{\ell(k)-(J+1)}, Y^{\ell(k)-(J+1)}) &= \lim_{k \rightarrow \infty} \mathcal{F}(X^{\ell(k)-J} - \Delta_{X^{\ell(k)-(J+1)}}, Y^{\ell(k)-J} - \Delta_{Y^{\ell(k)-(J+1)}}) \\ &= \lim_{k \rightarrow \infty} \mathcal{F}(X^{\ell(k)-J}, Y^{\ell(k)-J}) = \tilde{\mathcal{F}}, \end{aligned}$$

559 where the second equality follows because $\{(X^k, Y^k)\}$ is bounded, any cluster point
560 of $\{(X^k, Y^k)\}$ is in $\text{dom}\mathcal{F}$ and \mathcal{F} is uniformly continuous on any compact subset of
561 $\text{dom}\mathcal{F}$ under (a1) in Assumption 5.1. Hence, (5.15b) also holds for $j = J + 1$. This
562 completes the induction.

563 We are now ready to prove the main result in this statement. Indeed, from (5.12),
564 we can see $k - N \leq \ell(k) \leq k$ (without loss of generality, we assume that k is large
565 enough such that $k \geq N$). Thus, for any k , we must have $k - N - 1 = \ell(k) - j_k$ for
566 $1 \leq j_k \leq N + 1$. Then, we have

$$567 \quad \begin{aligned} \|\Delta_{X^{k-N-1}}\|_F &= \|\Delta_{X^{\ell(k)-j_k}}\|_F \leq \max_{1 \leq j \leq N+1} \|\Delta_{X^{\ell(k)-j}\|_F, \\ \|\Delta_{Y^{k-N-1}}\|_F &= \|\Delta_{Y^{\ell(k)-j_k}}\|_F \leq \max_{1 \leq j \leq N+1} \|\Delta_{Y^{\ell(k)-j}\|_F. \end{aligned}$$

568 This together with (5.15a) implies that

$$569 \quad \begin{aligned} \lim_{k \rightarrow \infty} \Delta_{X^k} &= \lim_{k \rightarrow \infty} \Delta_{X^{k-N-1}} = 0, \\ \lim_{k \rightarrow \infty} \Delta_{Y^k} &= \lim_{k \rightarrow \infty} \Delta_{Y^{k-N-1}} = 0. \end{aligned}$$

570 This proves the statement (ii).

571 *Statement (iii).* Again, let (X^*, Y^*) be a cluster point of $\{(X^k, Y^k)\}$ and let
572 $\{(X^{k_i}, Y^{k_i})\}$ be a convergent subsequence such that $\lim_{i \rightarrow \infty} (X^{k_i}, Y^{k_i}) = (X^*, Y^*)$. Recall
573 that $(X^*, Y^*) \in \text{dom}\mathcal{F}$. On the other hand, it is easy to see from (4.5) that $\lim_{i \rightarrow \infty} Z^{k_i} =$
574 Z^* , where Z^* is given by (3.4). Thus, it can be shown as in (3.9) that

$$575 \quad (5.17) \quad \alpha(Z^* - X^*(Y^*)^\top) + \beta \mathcal{A}^*(\mathcal{A}(Z^*) - \mathbf{b}) = 0.$$

We next show that

$$(5.18a) \quad \begin{cases} 0 \in \partial\Psi(X^*) + \alpha(X^*(Y^*)^\top - Z^*)Y^*, \\ (5.18b) \quad 0 \in \partial\Phi(Y^*) + \alpha(X^*(Y^*)^\top - Z^*)^\top X^*. \end{cases}$$

576 We start by showing (5.18a) in the following cases:

- 577 • **Proximal&Prox-linear:** In these two cases, passing to the limit along $\{(X^{k_i},$
578 $Y^{k_i})\}$ in (5.1a) or (5.1b) with X^{k_i+1} in place of U and $\bar{\mu}_{k_i}$ in place of μ_k , and
579 invoking (a1) in Assumption 5.1, statements (i), (ii), $(X^*, Y^*) \in \text{dom}\mathcal{F}$ and (2.1),
580 we obtain (5.18a).
- 581 • **Hierarchical-prox:** In this case, passing to the limit along $\{(X^{k_i}, Y^{k_i})\}$ in (5.1c)
582 with X^{k_i+1} in place of U and $\bar{\mu}_{k_i}$ in place of μ_k , and invoking (a1) in Assumption
583 5.1, statements (i), (ii), $(X^*, Y^*) \in \text{dom}\mathcal{F}$ and (2.1), we have

$$584 \quad 0 \in \partial\psi_i(\mathbf{x}_i^*) + \alpha(X^*(Y^*)^\top - Z^*)\mathbf{y}_i^*$$

585 for any $i = 1, 2, \dots, r$. Then, stacking them up, we obtain (5.18a).

586 Similarly, we can obtain (5.18b). Thus, combining (5.17), (5.18a) and (5.18b), we see
 587 that (X^*, Y^*, Z^*) is a stationary point of $\Theta_{\alpha, \beta}$, which further implies (X^*, Y^*) is a
 588 stationary point of \mathcal{F} from Theorem 3.3. This proves statement (iii). \square

589 **REMARK 5.2 (Comment on (a3) in Assumption 5.1).** *If Φ and Ψ are the*
 590 *indicator functions of some nonempty closed sets, Theorem 5.3 can remain valid under*
 591 *the weaker condition on α and β that $\frac{1}{\alpha} + \frac{1}{\beta} > 0$ with a slight modification in (4.6)*
 592 *of Algorithm 1. Indeed, when Φ and Ψ are the indicator functions, one can see from*
 593 *Remark 3.1 and the proofs of Lemma 5.1 and Proposition 5.2 that if $\frac{1}{\alpha} + \frac{1}{\beta} > 0$, then*

$$594 \quad \begin{aligned} \mathcal{F}(U, V) - \mathcal{F}(X^k, Y^k) &= \left(\frac{1}{\alpha} + \frac{1}{\beta}\right) (\Theta_{\alpha, \beta}(U, V, W) - \Theta_{\alpha, \beta}(X^k, Y^k, Z^k)) \\ &\leq -\left(\frac{1}{\alpha} + \frac{1}{\beta}\right) \left(\frac{\mu_k - (\alpha + 2\gamma\rho)\|Y^k\|^2}{2} \cdot \|U - X^k\|_F^2 + \frac{\sigma_k - (\alpha + 2\gamma\rho)\|U\|^2}{2} \cdot \|V - Y^k\|_F^2 \right), \end{aligned}$$

595 *and the line search criterion is well defined with c replaced by $\left(\frac{1}{\alpha} + \frac{1}{\beta}\right)c$. Moreover,*
 596 *recalling [28, Exercise 8.14], we see that $\partial\Psi$ and $\partial\Phi$ are normal cones. Thus, following*
 597 *Remark 3.2 and the similar augments in Theorem 5.3, we can obtain the same results*
 598 *when $\frac{1}{\alpha} + \frac{1}{\beta} > 0$ with c replaced by $\left(\frac{1}{\alpha} + \frac{1}{\beta}\right)c$ in (4.6) of Algorithm 1.*

599 **REMARK 5.3 (Comments on updating μ_k^{\max} and σ_k^{\max}).** *In Algorithm 1,*
 600 *we need to evaluate $\mu_k^{\max} = (\alpha + 2\gamma\rho)\|Y^k\|^2 + c$ and $\sigma_k^{\max} = (\alpha + 2\gamma\rho)\|U\|^2 + c$ in*
 601 *each iteration. However, computing the spectral norms of Y^k and U might be costly,*
 602 *especially when r is large. Hence, in our experiments, instead of computing $\|Y^k\|^2$*
 603 *and $\|U\|^2$, we compute $\|Y^k\|_F^2$ and $\|U\|_F^2$, and update μ_k^{\max} and σ_k^{\max} by $\mu_k^{\max} =$
 604 $(\alpha + 2\gamma\rho)\|Y^k\|_F^2 + c$ and $\sigma_k^{\max} = (\alpha + 2\gamma\rho)\|U\|_F^2 + c$ instead. Since $\|Y^k\| \leq \|Y^k\|_F$
 605 and $\|U\| \leq \|U\|_F$, it follows from (5.3) that*

$$606 \quad \mathcal{F}(U, V) - \mathcal{F}(X^k, Y^k) \leq -\frac{\mu_k - (\alpha + 2\gamma\rho)\|Y^k\|_F^2}{2} \|U - X^k\|_F^2 - \frac{\sigma_k - (\alpha + 2\gamma\rho)\|U\|_F^2}{2} \|V - Y^k\|_F^2.$$

607 *Then, one can show that Proposition 5.2 and Theorem 5.3 remain valid. In addition,*
 608 *we compute the quantities $\|U\|_F^2$ and $\|Y^k\|_F^2$ by $\text{tr}(U^\top U)$ and $\text{tr}((Y^k)^\top Y^k)$, respecti-*
 609 *vely. For some cases, the matrices $U^\top U$ and $(Y^k)^\top Y^k$ can be used repeatedly in*
 610 *updating the variables and evaluating the objective value and successive changes to*
 611 *reduce the cost of line search; see a concrete example in Section 6.1.*

612 **6. Numerical experiments.** In this section, we conduct numerical experiments
 613 to test our algorithm for NMF and MC on real datasets. All experiments are run in
 614 MATLAB R2015b on a 64-bit PC with an Intel Core i7-4790 CPU (3.60 GHz) and
 615 32 GB of RAM equipped with Windows 10 OS.

616 **6.1. Non-negative matrix factorization.** We first consider NMF

$$617 \quad (6.1) \quad \min_{X, Y} \frac{1}{2} \|XY^\top - M\|_F^2 \quad \text{s.t.} \quad X \geq 0, \quad Y \geq 0,$$

618 where $X \in \mathbb{R}^{m \times r}$ and $Y \in \mathbb{R}^{n \times r}$ are decision variables. Note that the feasible set of
 619 (6.1) is unbounded. We hence focus on the following model:

$$620 \quad (6.2) \quad \min_{X, Y} \frac{1}{2} \|XY^\top - M\|_F^2 \quad \text{s.t.} \quad 0 \leq X \leq X^{\max}, \quad 0 \leq Y \leq Y^{\max},$$

621 where $X^{\max} \geq 0$ and $Y^{\max} \geq 0$ are upper bound matrices. One can show that, when
 622 X_{ij}^{\max} and Y_{ij}^{\max} are sufficiently large³, solving (6.2) gives a solution of (6.1). In our

³The estimations of X_{ij}^{\max} and Y_{ij}^{\max} have been discussed in [9, Page 67].

623 experiments, for simplicity, we set $X_{ij}^{\max} = 10^{16}$ and $Y_{ij}^{\max} = 10^{16}$ for all (i, j) . Now,
 624 we see that (6.2) corresponds to (1.1) with $\Psi(X) = \delta_{\mathcal{X}}(X)$, $\Phi(Y) = \delta_{\mathcal{Y}}(Y)$ and $\mathcal{A} = \mathcal{I}$,
 625 where $\mathcal{X} = \{X \in \mathbb{R}^{m \times r} : 0 \leq X \leq X^{\max}\}$ and $\mathcal{Y} = \{Y \in \mathbb{R}^{n \times r} : 0 \leq Y \leq Y^{\max}\}$.
 626 We apply NAUM to solving (6.2), and use (4.1c) and (4.2c) to update U and V . The
 627 specific updates of Z^k , \mathbf{u}_i and \mathbf{v}_i are

$$\begin{aligned} Z^k &= \frac{\alpha}{\alpha+\beta} X^k (Y^k)^\top + \frac{\beta}{\alpha+\beta} M, \\ \mathbf{u}_i &= \max \left\{ 0, \min \left\{ \mathbf{x}_i^{\max}, \frac{\alpha P_i^k \mathbf{y}_i^k + \mu_k \mathbf{x}_i^k}{\alpha \|\mathbf{y}_i^k\|^2 + \mu_k} \right\} \right\}, \quad i = 1, 2, \dots, r, \\ \mathbf{v}_i &= \max \left\{ 0, \min \left\{ \mathbf{y}_i^{\max}, \frac{\alpha (Q_i^k)^\top \mathbf{u}_i + \sigma_k \mathbf{y}_i^k}{\alpha \|\mathbf{u}_i\|^2 + \sigma_k} \right\} \right\}, \quad i = 1, 2, \dots, r, \end{aligned}$$

629 where P_i^k and Q_i^k are defined in (4.8). Note that here it is not necessary to update
 630 Z^k explicitly. Indeed, we can directly compute $P_i^k \mathbf{y}_i^k$ and $(Q_i^k)^\top \mathbf{u}_i$ by substituting
 631 Z^k as below:

$$(6.3) \quad \begin{aligned} P_i^k \mathbf{y}_i^k &= \frac{\alpha}{\alpha+\beta} X^k (Y^k)^\top \mathbf{y}_i^k + \frac{\beta}{\alpha+\beta} M \mathbf{y}_i^k - \sum_{j=1}^{i-1} \mathbf{u}_j (\mathbf{y}_j^k)^\top \mathbf{y}_i^k - \sum_{j=i+1}^r \mathbf{x}_j^k (\mathbf{y}_j^k)^\top \mathbf{y}_i^k, \\ (Q_i^k)^\top \mathbf{u}_i &= \frac{\alpha}{\alpha+\beta} Y^k (X^k)^\top \mathbf{u}_i + \frac{\beta}{\alpha+\beta} M^\top \mathbf{u}_i - \sum_{j=1}^{i-1} \mathbf{v}_j \mathbf{u}_j^\top \mathbf{u}_i - \sum_{j=i+1}^r \mathbf{y}_j^k \mathbf{u}_j^\top \mathbf{u}_i. \end{aligned}$$

633 When computing $X^k (Y^k)^\top \mathbf{y}_i^k$ and $Y^k (X^k)^\top \mathbf{u}_i$ in the above, we first compute $(Y^k)^\top \mathbf{y}_i^k$
 634 and $(X^k)^\top \mathbf{u}_i$ to avoid forming the huge $(m \times n)$ matrix $X^k (Y^k)^\top$. Moreover, the
 635 matrices $(X^k)^\top U$, $U^\top U$, $(Y^k)^\top Y^k$ and $M^\top U$ that have been computed in (6.3) can
 636 be used again to evaluate the successive changes and the objective value as follows:

$$\begin{aligned} \|U - X^k\|_F^2 &= \text{tr}(U^\top U) - 2\text{tr}((X^k)^\top U) + \text{tr}((X^k)^\top X^k), \\ \|V - Y^k\|_F^2 &= \text{tr}(V^\top V) - 2\text{tr}((Y^k)^\top V) + \text{tr}((Y^k)^\top Y^k), \\ \|UV^\top - M\|_F^2 &= \text{tr}((U^\top U)(V^\top V)) - 2\text{tr}((M^\top U)V^\top) + \|M\|_F^2. \end{aligned}$$

638 In the above relations, $(X^k)^\top X^k$ and $(Y^k)^\top Y^k$ can be obtained from $U^\top U$ and $V^\top V$
 639 in the previous iteration, respectively, and $\|M\|_F^2$ can be computed in advance. Ad-
 640 ditionally, as we discussed in Remark 5.3, $\text{tr}((Y^k)^\top Y^k)$ and $\text{tr}(U^\top U)$ can also be
 641 used in computing μ_k^{\max} and σ_k^{\max} , respectively. These techniques were also used in
 642 many popular algorithms for NMF to reduce the computational cost (see, for example,
 643 [2, 9, 10, 18, 37]).

644 The experiments are conducted on the face datasets (dense matrices) and the
 645 text datasets (sparse matrices). For face datasets, we use CBCL⁴, ORL⁵ [29] and the
 646 extended Yale Face Database B (e-YaleB)⁶ [19] for our test. CBCL contains 2429
 647 images of faces with 19×19 pixels, ORL contains 400 images of faces with 112×92
 648 pixels, and e-YaleB contains 2414 images of faces with 168×192 pixels. In our
 649 experiments, for each face dataset, each image is vectorized and stacked as a column
 650 of a data matrix M of size $m \times n$. For text datasets, we use three datasets from the
 651 CLUTO toolkit⁷. The specific values of m and n for each dataset and the values of r
 652 used for our tests are summarized in Table 1.

653 The parameters in NAUM are set as follows: $\mu^{\min} = \bar{\mu}_{-1} = 1$, $\sigma^{\min} = \bar{\sigma}_{-1} =$
 654 1 , $\sigma^{\max} = 10^6$, $\tau = 4$, $c = 10^{-4}$, $N = 3$, $\mu_k^0 = \max\{0.1\bar{\mu}_{k-1}, \mu^{\min}\}$ and $\sigma_k^0 =$

⁴Available in <http://cbcl.mit.edu/cbcl/software-datasets/FaceData2.html>.

⁵Available in <http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html>.

⁶Available in <http://vision.ucsd.edu/~iskwak/ExtYaleDatabase/ExtYaleB.html>.

⁷Available in <http://glaros.dtc.umn.edu/gkhome/cluto/cluto/download>.

TABLE 1
Real data sets

Face Datasets (dense matrices)					Text Datasets (sparse matrices)				
Data	Pixels	m	n	r	Data	Sparsity	m	n	r
CBCL	19×19	361	2429	30, 60	classic	99.92%	7094	41681	10, 20
ORL	112×92	10304	400	30, 60	sports	99.14%	8580	14870	10, 20
e-YaleB	168×192	32256	2414	30, 60	ohscal	99.47%	11162	11465	10, 20

655 $\min \{ \max \{ 0.1\bar{\sigma}_{k-1}, \sigma^{\min} \}, \sigma^{\max} \}$ for any $k \geq 0$. Moreover, we set $\beta = \frac{\alpha}{\alpha-1}$, $\gamma =$
 656 $\max \{ 0, -\alpha, -(\alpha + \beta) \}$ and $\rho = \max \{ 1, \alpha^2 / (\alpha + \beta)^2 \}$ for some given α .

657 We then compare the performances of NAUM with different α . In our compari-
 658 sons, we initialize NAUM with different α at the same random initialization (X^0, Y^0) ⁸
 659 and terminate them if one of the following stopping criteria is satisfied:

- 660 • $\frac{|\mathcal{F}_{\text{nmf}}^k - \mathcal{F}_{\text{nmf}}^{k-1}|}{\mathcal{F}_{\text{nmf}}^k + 1} \leq 10^{-4}$ holds for 3 consecutive iterations;
- 661 • $\frac{\|X^k - X^{k-1}\|_F + \|Y^k - Y^{k-1}\|_F}{\|X^k\|_F + \|Y^k\|_F + 1} \leq 10^{-4}$ holds,

662 where $\mathcal{F}_{\text{nmf}}^k := \frac{1}{2} \|X^k (Y^k)^\top - M\|_F^2$ denotes the objective value at (X^k, Y^k) . Table 2
 663 presents the results of NAUM with different α for two face datasets (CBCL and ORL)
 664 and $r = 30, 60$. In the table, “iter” denotes the number of iterations; “relerr” denotes
 665 the relative error $\frac{\|X^*(Y^*)^\top - M\|_F}{\|M\|_F}$, where (X^*, Y^*) is a terminating point obtained by
 666 each NUAM in a trial; “time” denotes the computational time (in seconds). All the
 667 results presented are the average of 10 independent trials. From Table 2, we can see
 668 that NAUM with a relatively small α (e.g., 0.6 and 0.8) has better numerical perfor-
 669 mance. However, α cannot be too small. Observe that NAUM with $\alpha = 0.5, 0.4, 0.2$
 670 are not competitive and, surprisingly, $\alpha = 0.5$ leads to the worst performance. In
 671 view of this, we do not choose $\alpha < 0.6$ in our following experiments for NMF.

TABLE 2
Comparisons of NAUM with different α

α	iter	relerr	time	α	iter	relerr	time
CBCL, $r = 30$				CBCL, $r = 60$			
2.0	488	1.0519e-01	1.72	2.0	626	7.4388e-02	4.94
1.1	381	1.0448e-01	1.35	1.1	555	7.3477e-02	4.38
0.8	315	1.0426e-01	1.09	0.8	511	7.2986e-02	4.09
0.6	268	1.0406e-01	0.94	0.6	419	7.2998e-02	3.32
0.5	833	1.0593e-01	4.74	0.5	1372	7.5864e-02	19.49
0.4	440	1.0489e-01	3.05	0.4	599	7.4568e-02	10.02
0.2	556	1.0674e-01	4.18	0.2	782	7.7654e-02	14.30
ORL, $r = 30$				ORL, $r = 60$			
2.0	232	1.6673e-01	3.45	2.0	277	1.4078e-01	7.92
1.1	188	1.6619e-01	2.78	1.1	210	1.4042e-01	6.04
0.8	158	1.6603e-01	2.33	0.8	182	1.4017e-01	5.20
0.6	132	1.6578e-01	2.01	0.6	156	1.3996e-01	4.44
0.5	652	1.7216e-01	15.79	0.5	695	1.4583e-01	32.91
0.4	280	1.6615e-01	7.55	0.4	353	1.4061e-01	19.17
0.2	307	1.6753e-01	8.71	0.2	358	1.4272e-01	20.77

672 We next compare NAUM with two existing efficient algorithms⁹ for NMF: the

⁸We use the Matlab commands: $X0 = \max(0, \text{randn}(m, r)); Y0 = \max(0, \text{randn}(n, r)); X0 = X0 / \text{norm}(X0, 'fro') * \text{sqrt}(\text{norm}(M, 'fro'));$ $Y0 = Y0 / \text{norm}(Y0, 'fro') * \text{sqrt}(\text{norm}(M, 'fro'));$

⁹Most existing algorithms are directly developed for (6.1). However, they need the assumption that the sequence generated is bounded in their convergence analysis. Although this assumption is uncheckable and may fail, these algorithms always work well in practice. Thus, we directly use these algorithms in our comparisons, rather than modifying them for (6.2).

673 hierarchical alternating least squares (HALS) method¹⁰ (see, for example, [8, 9, 10,
674 11, 20, 21]) and the block coordinate descent method for NMF (BCD-NMF¹¹) (see
675 Algorithm 2 in Section 3.2 in [40]).

676 To better evaluate the performances of different algorithms, we follow [11] to use
677 an evolution of the objective function value. To define this evolution, we first define

$$678 \quad e(k) := \frac{\mathcal{F}^k - \mathcal{F}_{\min}}{\mathcal{F}^0 - \mathcal{F}_{\min}},$$

679 where \mathcal{F}^k denotes the objective function value obtained by an algorithm at (X^k, Y^k)
680 and \mathcal{F}_{\min} denotes the minimum of the objective function values obtained among *all*
681 algorithms across *all* initializations. We also use $\mathcal{T}(k)$ to denote the total computa-
682 tional time after completing the k -th iteration of an algorithm. Thus, $\mathcal{T}(0) = 0$ and
683 $\mathcal{T}(k)$ is non-decreasing with respect to k . Then, the evolution of the function value
684 obtained from a particular algorithm with respect to time t is defined as

$$685 \quad E(t) := \min \{e(k) : k \in \{i : \mathcal{T}(i) \leq t\}\}.$$

686 One can see that $0 \leq E(t) \leq 1$ (since $0 \leq e(k) \leq 1$ for all k) and $E(t)$ is non-increasing
687 with respect to t . $E(t)$ can be considered as a normalized measure of the reduction of
688 the function value with respect to time. For a given matrix M and a positive integer
689 r , one can take the average of $E(t)$ over several independent trials with different
690 initializations, and plot the average $E(t)$ within time t for a given algorithm.

691 In our experiments, we initialize all the algorithms at the same random initial
692 point (X^0, Y^0) and terminate them *only* by the maximum running time T^{\max} . The
693 specific values of T^{\max} are given in Fig. 1 and Fig. 2. Additionally, we use the default
694 settings for BCD-NMF. For NAUM, we choose $\alpha = 0.6, 0.8, 1.1, 2$. We then plot the
695 average $E(t)$ for each algorithm within time T^{\max} .

696 Fig. 1 and Fig. 2 show the average $E(t)$ of 30 independent trials for NMF on face
697 datasets and text datasets, respectively. From the results, we can see that NAUM
698 with $\alpha = 0.6$ performs best in most cases, and NAUM with $\alpha = 0.6$ or 0.8 always
699 performs better than NAUM with $\alpha > 1$. This shows that choosing α and β under the
700 weaker condition $\frac{1}{\alpha} + \frac{1}{\beta} = 1$ (hence α can be small than 1) can improve the numerical
701 performance of NAUM.

702 **6.2. Matrix completion.** We next consider a recent model for MC:

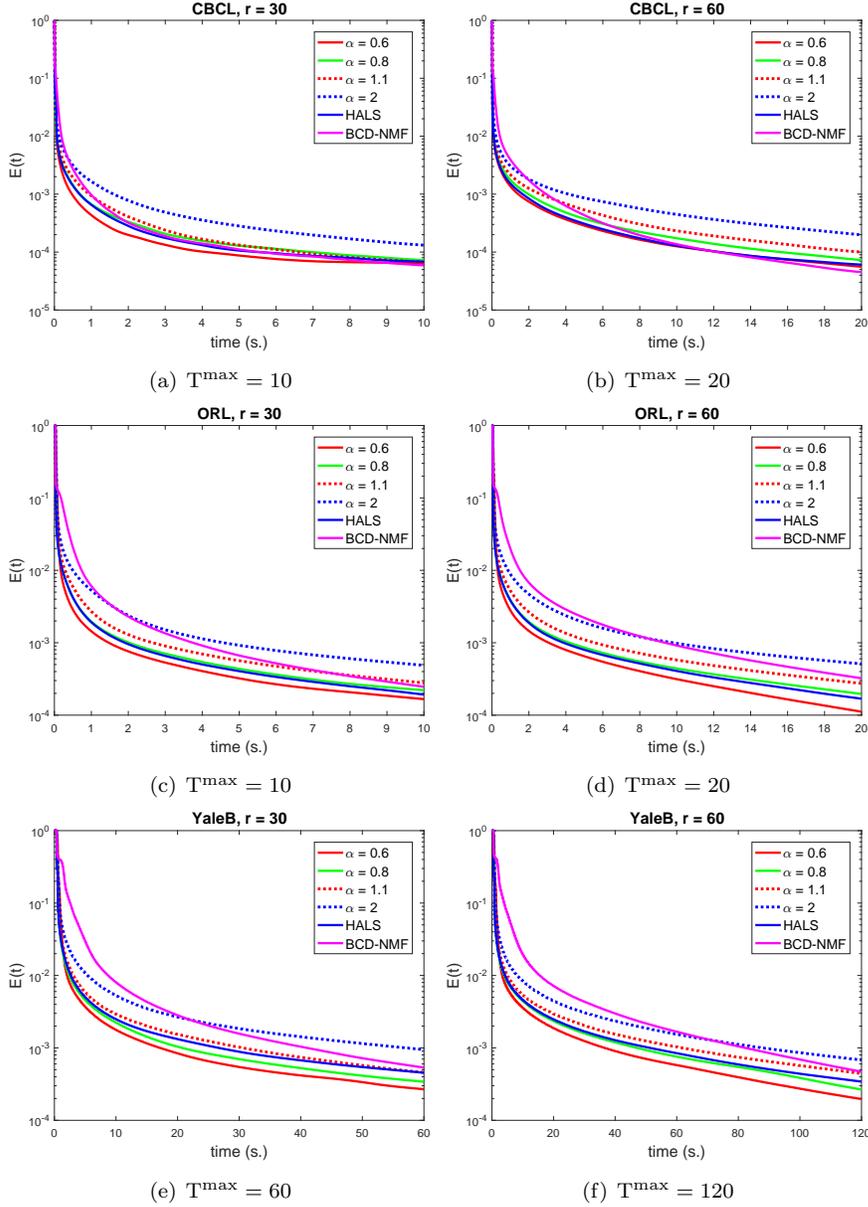
$$703 \quad (6.4) \quad \min_{X, Y} \frac{\eta}{2} \|X\|_* + \frac{\eta}{2} \|Y\|_* + \frac{1}{2} \|\mathcal{P}_{\Omega}(XY^{\top} - M)\|_F^2,$$

704 where $\eta > 0$ is a penalty parameter, Ω is the index set of the known entries of M ,
705 and $\mathcal{P}_{\Omega}(Z)$ keeps the entries of Z in Ω and sets the remaining ones to zero. This
706 model was first considered in [30, 31] and was shown to be equivalent to Schatten- $\frac{1}{2}$
707 quasi-norm minimization. Encouraging numerical performance of this model has also

¹⁰HALS for (6.1) is given by

$$\begin{aligned} \mathbf{x}_i^{k+1} &= \max \left\{ 0, \frac{M \mathbf{y}_i^k - \sum_{j=1}^{i-1} \mathbf{x}_j^{k+1} (\mathbf{y}_j^k)^{\top} \mathbf{y}_i^k - \sum_{j=i+1}^r \mathbf{x}_j^k (\mathbf{y}_j^k)^{\top} \mathbf{y}_i^k}{\|\mathbf{y}_i^k\|^2} \right\}, \quad i = 1, \dots, r, \\ \mathbf{y}_i^{k+1} &= \max \left\{ 0, \frac{M^{\top} \mathbf{x}_i^{k+1} - \sum_{j=1}^{i-1} \mathbf{y}_j^{k+1} (\mathbf{x}_j^{k+1})^{\top} \mathbf{x}_i^{k+1} - \sum_{j=i+1}^r \mathbf{y}_j^k (\mathbf{x}_j^{k+1})^{\top} \mathbf{x}_i^{k+1}}{\|\mathbf{x}_i^{k+1}\|^2} \right\}, \quad i = 1, \dots, r. \end{aligned}$$

¹¹Available at <http://www.math.ucla.edu/~wotaoyin/papers/bcu/nmf/index.html>.


 FIG. 1. Average $E(t)$ of 30 independent trials for NMF on face datasets.

708 been reported in [30, 31]. Note that (6.4) corresponds to (1.1) with $\Psi(X) = \frac{\eta}{2}\|X\|_*$,
 709 $\Phi(Y) = \frac{\eta}{2}\|Y\|_*$ and $\mathcal{A} = \mathcal{P}_\Omega$. Thus, we can apply NAUM with (4.1b) and (4.2b) to

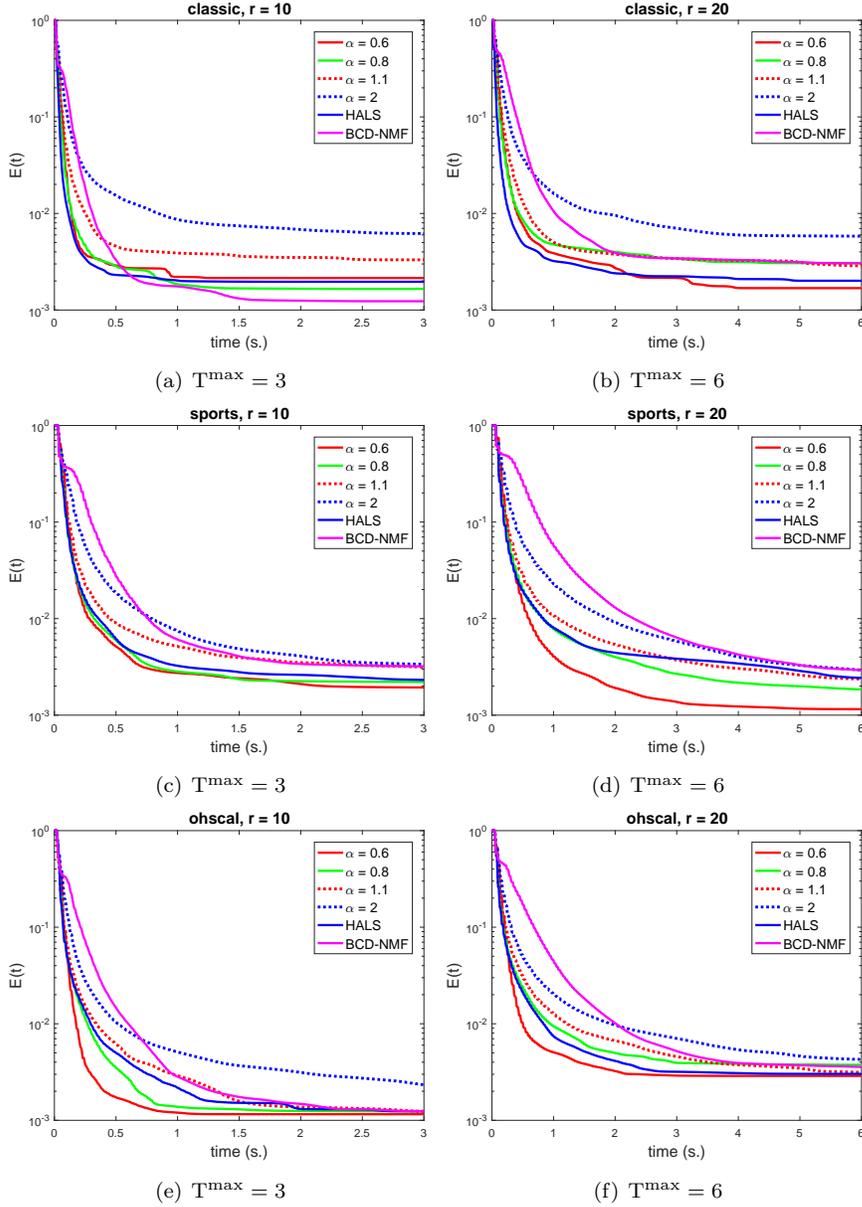


FIG. 2. Average $E(t)$ of 30 independent trials for NMF on text datasets.

710 solving (6.4). The updates of Z^k , U and V are

$$\begin{aligned}
 Z^k &= X^k (Y^k)^\top + \frac{\beta}{\alpha + \beta} \mathcal{P}_\Omega (M - X^k (Y^k)^\top), \\
 U &= \mathcal{S}_{\eta/(2\mu_k)} \left(X^k - \frac{\alpha}{\mu_k} (X^k (Y^k)^\top - Z^k) Y^k \right), \\
 V &= \mathcal{S}_{\eta/(2\sigma_k)} \left(Y^k - \frac{\alpha}{\sigma_k} (U (Y^k)^\top - Z^k)^\top U \right).
 \end{aligned}$$

711

712 Substituting Z^k into U and V and using $\frac{1}{\alpha} + \frac{1}{\beta} = 1$ gives

$$713 \quad (6.5) \quad \begin{aligned} U &= \mathcal{S}_{\eta/(2\mu_k)} \left(X^k - \frac{1}{\mu_k} [\mathcal{P}_\Omega(X^k(Y^k)^\top - M)] Y^k \right), \\ V &= \mathcal{S}_{\eta/(2\sigma_k)} \left(Y^k - \frac{\alpha}{\sigma_k} Y^k (U - X^k)^\top U - \frac{1}{\sigma_k} [\mathcal{P}_\Omega(X^k(Y^k)^\top - M)]^\top U \right). \end{aligned}$$

714 Thus, similar to NAUM for NMF, we do not need to update Z^k explicitly for MC.

715 We compare NAUM with proximal alternating linearized minimization (PALM),
716 which was proposed in [4] and was used to solve (6.4) in [30, 31]. For ease of future
717 reference, we recall that the PALM for solving (6.4) is given by

$$718 \quad \begin{aligned} X^{k+1} &= \mathcal{S}_{\frac{\eta}{2\|Y^k\|^2}} \left(X^k - \frac{1}{\|Y^k\|^2} [\mathcal{P}_\Omega(X^k(Y^k)^\top - M)] Y^k \right), \\ Y^{k+1} &= \mathcal{S}_{\frac{\eta}{2\|X^{k+1}\|^2}} \left(Y^k - \frac{1}{\|X^{k+1}\|^2} [\mathcal{P}_\Omega(X^{k+1}(Y^k)^\top - M)]^\top X^{k+1} \right). \end{aligned}$$

719 For NAUM, we use the same parameter settings as in Section 6.1, but choose $\alpha =$
720 0.4, 0.6, 1.1. All the algorithms are initialized at the same random initialization
721 (X^0, Y^0) ¹² and terminated if one of the following stopping criteria is satisfied:

- 722 • $\frac{|\mathcal{F}_{\text{mc}}^k - \mathcal{F}_{\text{mc}}^{k-1}|}{\mathcal{F}_{\text{mc}}^k} \leq 10^{-4}$ holds for 3 consecutive iterations;
- 723 • $\frac{\|X^k - X^{k-1}\|_F + \|Y^k - Y^{k-1}\|_F}{\|X^k\|_F + \|Y^k\|_F} \leq 10^{-4}$ holds;
- 724 • the running time is more than 300 seconds,

725 where $\mathcal{F}_{\text{mc}}^k := \frac{\eta}{2}\|X^k\|_* + \frac{\eta}{2}\|Y^k\|_* + \frac{1}{2}\|\mathcal{P}_\Omega(X^k(Y^k)^\top - M)\|_F^2$ denotes the objective
726 function value obtained by each algorithm at (X^k, Y^k) .

727 Table 3 presents the numerical results of different algorithms for different pro-
728 blems, where two face datasets (CBCL and ORL) are used as our test matrices M
729 and a subset Ω of entries is sampled uniformly at random. In the table, sr denotes
730 the sampling ratio, i.e., a subset Ω of (rounded) $mn * sr$ entries is sampled; r denotes
731 the rank used for test; “iter” denotes the number of iterations; “Normalized fval”
732 denotes the normalized function value $\frac{\mathcal{F}(X^*, Y^*) - \mathcal{F}_{\text{min}}}{\mathcal{F}_{\text{max}} - \mathcal{F}_{\text{min}}}$, where (X^*, Y^*) is obtained
733 by each algorithm, $\mathcal{F}(X^*, Y^*)$ is the function value at (X^*, Y^*) for each algorithm
734 and \mathcal{F}_{max} (resp. \mathcal{F}_{min}) denotes the maximum (resp. minimum) of the terminating
735 function values obtained from *all* algorithms in *a* trial (one random initialization and
736 Ω); “RecErr” denotes the recovery error $\frac{\|X^*(Y^*)^\top - M\|_F}{\|M\|_F}$. All the results presented are
737 the average of 10 independent trials.

738 From Table 3, we can see that NAUM with $\alpha = 0.4$ gives the smallest function
739 values and the smallest recovery error within least CPU time in most cases. Moreover,
740 NAUM with $\alpha = 0.6$ also performs better than NAUM with $\alpha = 1.1$ and PALM with
741 respect to the function value and the recovery error in most cases. This again shows
742 that a flexible choice of α and β can lead to better numerical performances and the
743 choice of $\alpha = 0.4$ performs best for MC from our experiments.

744 **7. Concluding remarks.** In this paper, we consider a class of matrix facto-
745 rization problems involving two blocks of variables. To solve this kind of possibly
746 nonconvex, nonsmooth and non-Lipschitz problems, we introduce a specially con-
747 structed potential function $\Theta_{\alpha, \beta}$ defined in (1.2) which contains one auxiliary block
748 of variables. We then develop a non-monotone alternating updating method with a
749 suitable line search criterion based on this potential function. Unlike other existing

¹²We use the Matlab commands: $X0 = \text{randn}(m, r)$; $Y0 = \text{randn}(n, r)$;

TABLE 3
Numerical results for MC on face datasets

η	data	sr	r	$\alpha = 0.4$	$\alpha = 0.6$	$\alpha = 1.1$	PALM	$\alpha = 0.4$	$\alpha = 0.6$	$\alpha = 1.1$	PALM
				iter				Normalized fval			
5	CBCL	0.5	30	780	1189	3320	3306	1.13e-01	7.50e-02	4.52e-01	1
		0.5	60	921	1218	3850	4654	3.24e-02	5.10e-02	3.85e-01	1
		0.2	30	1174	2366	4767	3573	8.01e-03	2.21e-01	6.87e-01	9.60e-01
		0.2	60	1577	1919	5360	5037	1.03e-02	8.95e-02	8.08e-01	8.86e-01
	ORL	0.5	30	1218	1243	1241	1468	0	2.94e-01	5.06e-01	1
		0.5	60	1049	1051	1051	1327	0	1	4.00e-01	7.73e-01
		0.2	30	2074	325	385	2691	2.59e-03	7.01e-01	1	1.31e-01
		0.2	60	1551	1551	356	2222	0	3.82e-01	1	2.12e-01
10	CBCL	0.5	30	457	654	1793	1935	2.20e-02	1.29e-01	3.60e-01	9.81e-01
		0.5	60	514	594	1950	2559	2.65e-01	1.15e-01	3.79e-01	8.71e-01
		0.2	30	627	1313	2513	2116	1.91e-02	3.75e-02	8.35e-01	7.79e-01
		0.2	60	866	1095	2713	2889	2.07e-02	2.89e-02	9.22e-01	4.86e-01
	ORL	0.5	30	1003	1186	1192	1402	3.30e-02	1.47e-01	4.30e-01	1
		0.5	60	975	1009	1012	1276	0	8.58e-01	6.11e-01	9.99e-01
		0.2	30	1409	364	411	2646	0	7.16e-01	1	8.10e-02
		0.2	60	1241	1504	376	2185	4.05e-06	3.97e-02	1	2.21e-01
				CPU time				RecErr			
5	CBCL	0.5	30	35.56	54.14	151.23	119.05	1.05e-01	1.05e-01	1.06e-01	1.08e-01
		0.5	60	57.66	76.09	240.19	206.47	8.81e-02	9.02e-02	9.04e-02	8.99e-02
		0.2	30	34.04	68.57	137.97	75.56	1.37e-01	1.37e-01	1.38e-01	1.43e-01
		0.2	60	72.01	87.82	245.21	147.08	1.34e-01	1.35e-01	1.35e-01	1.36e-01
	ORL	0.5	30	294.20	300	300	300	1.72e-01	1.84e-01	2.01e-01	2.12e-01
		0.5	60	300	300	300	300	1.66e-01	2.11e-01	2.05e-01	2.11e-01
		0.2	30	300	47.35	55.86	300	2.08e-01	3.04e-01	3.81e-01	2.24e-01
		0.2	60	300	300	69.21	300	2.16e-01	2.35e-01	3.49e-01	2.61e-01
10	CBCL	0.5	30	21.01	30.12	82.45	70.32	1.16e-01	1.19e-01	1.18e-01	1.17e-01
		0.5	60	32.40	37.38	122.51	113.80	1.09e-01	1.11e-01	1.14e-01	1.11e-01
		0.2	30	18.15	38.01	72.84	44.62	1.60e-01	1.61e-01	1.62e-01	1.60e-01
		0.2	60	39.13	49.37	123.74	83.52	1.57e-01	1.57e-01	1.58e-01	1.56e-01
	ORL	0.5	30	252.15	300	300	300	1.71e-01	1.77e-01	1.95e-01	2.08e-01
		0.5	60	289.57	300	300	300	1.53e-01	2.01e-01	2.03e-01	2.09e-01
		0.2	30	207.22	53.08	60.54	300	1.95e-01	3.06e-01	3.83e-01	2.14e-01
		0.2	60	243.45	295.60	74.09	300	1.87e-01	1.95e-01	3.60e-01	2.36e-01

750 methods such as those based on alternating minimization, our method essentially up-
751 dates the two blocks of variables alternately by solving subproblems related to $\Theta_{\alpha,\beta}$
752 and then updates the auxiliary block of variables by an explicit formula (see (4.5)).
753 Using the special structure of $\Theta_{\alpha,\beta}$, we demonstrate how some efficient computational
754 strategies for NMF can be used to solve the associated subproblems in our method.
755 Moreover, under some mild conditions, we establish that the sequence generated by
756 our method is bounded and any cluster point of the sequence gives a stationary point
757 of our problem. Finally, we conduct some numerical experiments for NMF and MC
758 on real datasets to illustrate the efficiency of our method.

759 Note that the parameter α (and $\beta = \alpha/(\alpha - 1)$) plays a significant role in our
760 NAUM. Although it has been observed in our experiments that a relatively small α
761 (e.g., 0.6, 0.8) can improve the numerical performance of NAUM, how to choose an
762 optimal α is still unknown. In view of the recent work [24] on adaptively choosing the
763 extrapolation parameter in FISTA for solving a class of possibly nonconvex problems,
764 it may be possible to derive a strategy to adaptively update α in our NAUM. This is
765 a possible future research topic.

766 **Acknowledgments.** The authors are grateful to Nicolas Gillis for his helpful
767 comments. The authors are also grateful to the editor and the anonymous referees
768 for their valuable suggestions and comments, which helped improve this paper.

- 770 [1] H. ATTOUCH, J. BOLTE, AND B. SVAITER, *Convergence of descent methods for semi-algebraic*
771 *and tame problems: proximal algorithms, forward-backward splitting, and regularized*
772 *Gauss-Seidel methods*, Mathematical Programming, 137 (2013), pp. 91–129.
- 773 [2] M. BERRY, M. BROWNE, A. LANGVILLE, V. PAUCA, AND R. PLEMMONS, *Algorithms and ap-*
774 *lications for approximate nonnegative matrix factorization*, Computational Statistics &
775 Data Analysis, 52 (2007), pp. 155–173.
- 776 [3] P. BISWAS, T. LIANG, K. TOH, Y. YE, AND T. WANG, *Semidefinite programming approaches*
777 *for sensor network localization with noisy distance measurements*, IEEE Transactions on
778 Automation Science and Engineering, 3 (2006), pp. 360–371.
- 779 [4] J. BOLTE, S. SABACH, AND M. TEOUBILLE, *Proximal alternating linearized minimization for*
780 *nonconvex and nonsmooth problems*, Mathematical Programming, 146 (2014), pp. 459–494.
- 781 [5] E. CANDÈS AND B. RECHT, *Exact matrix completion via convex optimization*, Foundations of
782 Computational Mathematics, 9 (2009), pp. 717–772.
- 783 [6] E. CANDÈS AND T. TAO, *The power of convex relaxation: Near-optimal matrix completion*,
784 IEEE Transactions on Information Theory, 56 (2010), pp. 2053–2080.
- 785 [7] X. CHEN, Z. LU, AND T. PONG, *Penalty methods for a class of non-Lipschitz optimization*
786 *problems*, SIAM Journal on Optimization, 26 (2016), pp. 1465–1492.
- 787 [8] A. CICHOCKI, R. ZDUNEK, AND S. AMARI, *Hierarchical ALS algorithms for nonnegative ma-*
788 *trix and 3D tensor factorization*, in International Conference on Independent Component
789 Analysis and Signal Separation, Springer, 2007, pp. 169–176.
- 790 [9] N. GILLIS, *Nonnegative Matrix Factorization: Complexity, Algorithms and Applications*, PhD
791 thesis, Université catholique de Louvain, 2011.
- 792 [10] N. GILLIS, *The why and how of nonnegative matrix factorization*, in Regularization, Optimiza-
- 793 tion, Kernels, and Support Vector Machines, Chapman and Hall/CRC, 2014, pp. 257–291.
- 794 [11] N. GILLIS AND F. GLINEUR, *Accelerated multiplicative updates and hierarchical ALS algorithms*
795 *for nonnegative matrix factorization*, Neural Computation, 24 (2012), pp. 1085–1105.
- 796 [12] P. GONG, C. ZHANG, Z. LU, J. HUANG, AND J. YE, *A general iterative shrinkage and thresh-*
797 *olding algorithm for non-convex regularized optimization problems*, in Proceedings of the
798 International Conference on Machine Learning, vol. 28, 2013, pp. 37–45.
- 799 [13] M. HONG, *Decomposing linearly constrained nonconvex problems by a proximal primal dual*
800 *approach: Algorithms, convergence and applications*, arXiv preprint arXiv:1604.00543,
801 (2016).
- 802 [14] B. JIANG, T. LIN, S. MA, AND S. ZHANG, *Structured nonconvex and nonsmooth optimization:*
803 *Algorithms and iteration complexity analysis*, arXiv preprint arXiv: 1605.02408, (2016).
- 804 [15] R. KESHAVAN, A. MONTANARI, AND S. OH, *Matrix completion from a few entries*, IEEE Tran-
- 805 sactions on Information Theory, 56 (2010), pp. 2980–2998.
- 806 [16] M.-J. LAI, Y. XU, AND W. YIN, *Improved iteratively reweighted least squares for unconstrained*
807 *smoothed ℓ_p minimization*, SIAM Journal on Numerical Analysis, 51 (2013), pp. 927–957.
- 808 [17] D. LEE AND H. SEUNG, *Learning the parts of objects by nonnegative matrix factorization*,
809 Nature, 401 (1999), pp. 788–791.
- 810 [18] D. LEE AND H. SEUNG, *Algorithms for non-negative matrix factorization*, in Proceedings of
811 NIPS, 2001, pp. 556–562.
- 812 [19] K. LEE, J. HO, AND D. KRIEGMAN, *Acquiring linear subspaces for face recognition under*
813 *variable lighting*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 27
814 (2005), pp. 684–698.
- 815 [20] L. LI AND Y. ZHANG, *FastNMF: Highly efficient monotonic fixed-point nonnegative matrix*
816 *factorization algorithm with good applicability*, Journal of Electronic Imaging, 18 (2009),
817 p. 033004.
- 818 [21] J. LIU, J. LIU, P. WONKA, AND J. YE, *Sparse non-negative tensor factorization using column-*
819 *wise coordinate descent*, Pattern Recognition, 45 (2012), pp. 649–656.
- 820 [22] Z. LIU AND L. VANDENBERGHE, *Interior-point method for nuclear norm approximation with*
821 *application to system identification*, SIAM Journal on Matrix Analysis and Applications,
822 31 (2009), pp. 1235–1256.
- 823 [23] K. MOHAN AND M. FAZEL, *Iterative reweighted algorithms for matrix rank minimization*, Jour-
- 824 nal of Machine Learning Research, 13 (2012), pp. 3441–3473.
- 825 [24] P. OCHS AND T. POCK, *Adaptive FISTA*, arXiv preprint arXiv:1711.04343, (2017).
- 826 [25] P. PAATERO AND U. TAPPER, *Positive matrix factorization: A non-negative factor model with*
827 *optimal utilization of error estimates of data values*, Environmetrics, 5 (1994), pp. 111–126.
- 828 [26] B. RECHT, M. FAZEL, AND P. PARRILO, *Guaranteed minimum-rank solutions of linear matrix*
829 *equations via nuclear norm minimization*, SIAM Review, 52 (2010), pp. 471–501.
- 830 [27] J. RENNIE AND N. SREBRO, *Fast maximum margin matrix factorization for collaborative pre-*
831 *diction*, in Proceedings of the 22nd international conference on Machine learning, ACM,

- 832 2005, pp. 713–719.
- 833 [28] R. ROCKAFELLAR AND R.-B. WETS, *Variational Analysis*, Springer, 1998.
- 834 [29] F. SAMARIA AND A. HARTER, *Parameterisation of a stochastic model for human face identifi-*
835 *cation*, in Proceedings of the Second IEEE Workshop on Applications of Computer Vision,
836 IEEE, 1994, pp. 138–142.
- 837 [30] F. SHANG, Y. LIU, AND J. CHENG, *Scalable algorithms for tractable Schatten quasi-norm mi-*
838 *nimization*, in Proceedings of the 30th AAAI Conference on Artificial Intelligence, 2016,
839 pp. 2016–2022.
- 840 [31] F. SHANG, Y. LIU, AND J. CHENG, *Tractable and scalable Schatten quasi-norm approximations*
841 *for rank minimization*, in Proceedings of the 19th International Conference on Artificial
842 Intelligence and Statistics, 2016, pp. 620–629.
- 843 [32] F. SHANG, Y. LIU, AND J. CHENG, *Unified scalable equivalent formulations for Schatten quasi-*
844 *norms*, arXiv preprint arXiv: 1606.00668, (2016).
- 845 [33] N. SREBRO, *Learning with Matrix Factorizations*, PhD thesis, Massachusetts Institute of
846 Technology, 2004.
- 847 [34] R. SUN AND Z.-Q. LUO, *Guaranteed matrix completion via non-convex factorization*, IEEE
848 Transactions on Information Theory, 62 (2016), pp. 6535–6579.
- 849 [35] M. UDELL, C. HORN, R. ZADEH, AND S. BOYD, *Generalized low rank models*, Foundations and
850 Trends in Machine Learning, 9 (2016), pp. 1–118.
- 851 [36] S. VAVASIS, *On the complexity of nonnegative matrix factorization*, SIAM Journal on Optimi-
852 zation, 20 (2009), pp. 1364–1377.
- 853 [37] Y. WANG AND Y. ZHANG, *Nonnegative matrix factorization: A comprehensive review*, IEEE
854 Transactions on Knowledge and Data Engineering, 25 (2013), pp. 1336–1353.
- 855 [38] Z. WEN, W. YIN, AND Y. ZHANG, *Solving a low-rank factorization model for matrix com-*
856 *pletion by a non-linear successive over-relaxation algorithm*, Mathematical Programming
857 Computation, 4 (2012), pp. 333–361.
- 858 [39] S. WRIGHT, R. NOWAK, AND M. FIGUEIREDO, *Sparse reconstruction by separable approxima-*
859 *tion*, IEEE Transactions on Signal Processing, 57 (2009), pp. 2479–2493.
- 860 [40] Y. XU AND W. YIN, *A block coordinate descent method for regularized multiconvex optimization*
861 *with applications to nonnegative tensor factorization and completion*, SIAM Journal on
862 Imaging Sciences, 6 (2013), pp. 1758–1789.
- 863 [41] Y. XU, W. YIN, Z. WEN, AND Y. ZHANG, *An alternating direction algorithm for matrix com-*
864 *pletion with nonnegative factors*, Frontiers of Mathematics in China, 7 (2012), pp. 365–384.
- 865 [42] M. ZHANG, Z.-H. HUANG, AND Y. ZHANG, *Restricted p -isometry properties of nonconvex matrix*
866 *recovery*, IEEE Transactions on Information Theory, 59 (2013), pp. 4316–4323.
- 867 [43] Y. ZHANG, *An alternating direction algorithm for nonnegative matrix factorization*, Preprint,
868 (2010).
- 869 [44] X. ZHOU, C. YANG, H. ZHAO, AND W. YU, *Low-rank modeling and its applications in image*
870 *analysis*, ACM Computing Survey, 47 (2015), p. 36.