

1 **ORTHOGONAL NONNEGATIVE MATRIX FACTORIZATION VIA**
2 **MINIMIZATION OVER THE NULL SPACE**

3 XIAOJUN CHEN*, WEN LI†, AND QILUN LUO‡

4 **Abstract.** This paper gives a necessary and sufficient condition for a nonnegative matrix that
5 has an orthogonal nonnegative matrix factorization (ONMF) via characterization of the null space.
6 We propose an optimization model to minimize the Frobenius norm of the product of a given nonneg-
7 ative matrix and a variable matrix subject to the constraints defined by the necessary and sufficient
8 condition. Moreover, we present an augmented Lagrangian algorithm for solving this minimization
9 model, and prove the global convergence to a stationary point. Two factor matrices for the ONMF
10 of the given matrix can be easily obtained by the outputs of the algorithm. Preliminary numerical
11 results using synthetic and real-world data with applications in clustering show that our approach
12 outperforms some existing ONMF methods regarding accuracy and robustness.

13 **Key words.** Orthogonal nonnegative matrix factorization, null space, augmented Lagrangian
14 method, clustering.

15 **MSC codes.** 90C26, 90C30, 65F30

16 **1. Introduction.** Given a nonnegative matrix $A \in \mathbb{R}_+^{m \times n}$ with $m \geq n$ and a
17 positive integer $r < n$, the minimization problem

18 (1.1) $\min \|A - BC^\top\|_F^2, \quad \text{subject to } B \in \mathbb{R}_+^{m \times r}, C \in \mathbb{R}_+^{n \times r}, C^\top C = I$

19 is a popular mathematical model for Orthogonal Nonnegative Matrix Factorization
20 (ONMF), which aims to approximate A as the product of two matrices B, C^\top in the
21 feasible set of (1.1). ONMF has been successfully applied in various areas including
22 clustering, image science, hyperspectral unmixing [12, 14, 19, 20]. However, solving
23 (1.1) is generally NP-hard due to the combinatorial features. In particular, for any
24 feasible point, each row of C has at most one positive element and each column of C
25 takes the unit norm. Many methods have been developed for solving (1.1) including
26 multiplicative update rules [7, 17], hierarchical alternating least squares (HALS) [13,
27 15, 21], and subspace exploration techniques [1]. Moreover, problem (1.1) can be
28 considered as a special case of optimization problems with nonnegative orthogonal
29 constraints [6, 12]. In this paper, we focus on problem (1.1) where the objective
30 function is a continuously differentiable bilinear composite function.

31 We say that a matrix $A \in \mathbb{R}_+^{m \times n}$ has an ONMF if the optimal value of (1.1) is 0,
32 that is, there exists a feasible point of (1.1) such that $A = BC^\top$. If A has an ONMF,
33 then $AC = B$ and thus $A = ACC^\top$. In [19], Pan and Ng proposed the following
34 minimization model

35 (1.2) $\min \|A - ACC^\top\|_F^2, \quad \text{subject to } C \in \mathbb{O}_+^{n \times r},$

*Department of Applied Mathematics, The Hong Kong Polytechnic University, Hung Hom, Kowloon, Hong Kong. Email: maxjchen@polyu.edu.hk. The research of this author is supported by the National Key R & D Program of China under grant 2023YFA1009303 and Hong Kong Research Grant Council project PolyU15300123.

†School of Mathematics and Statistics, Hanshan Normal University, Chaozhou, Guangdong, China. School of Mathematical Sciences, South China Normal University, Guangzhou, Guangdong, China. Email: liwen@scnu.edu.cn. The research is supported in part by the National Natural Science Foundation of China under grant 12471351.

‡School of Mathematical Sciences, South China Normal University, Guangzhou, Guangdong, China. Email: qilunluo92@gmail.com. The research is supported in part by the National Natural Science Foundation of China under grant 12301483, and the Natural Science Foundation of Guangdong Province of China under grants 2024A1515011527 and 2025A1515011611. Corresponding author.

36 where $\mathbb{O}_+^{n \times r} = \{C \in \mathbb{R}_+^{n \times r} \mid C^\top C = I\}$. Let $K = CC^\top$. Then $K = K^\top, K^2 =$
 37 $K, 0 \leq K_{i,j} \leq 1, \|K\|^2 = r$ [19, Theorem 3]. Using these properties of K , Pan and Ng
 38 proposed an algorithm to solve the following convex optimization model

$$39 \quad (1.3) \quad \min \|A - AK\|_F^2 + \theta \|K\|_1 + \beta \|K\|_*, \quad \text{subject to } K = K^\top, K \in \mathbb{R}_+^{n \times n}$$

40 and set eigenvectors corresponding to the r largest eigenvalues of the solution K^* of
 41 (1.3) to be C . In (1.3), $\|K\|_*$ is the nuclear norm of K , $\|K\|_1 = \sum_{i,j=1}^n |K_{i,j}|$, and $\theta >$
 42 $0, \beta > 0$. Under the assumption that A has an ONMF, some properties of the solution
 43 of (1.3) are established in [19]. Numerical results show Pan-Ng approach is better than
 44 some existing methods. Problem (1.3) is a nice convex approach for ONMF, but the
 45 solution of (1.3) may fail to satisfy $K^* = CC^\top$ and the eigenvectors corresponding to
 46 the r largest eigenvalues of K^* may have negative elements. Moreover, such approach
 47 arises an interesting question when A has an ONMF. To address this question and
 48 overcome difficulties in handling the orthogonal nonnegative constraints in ONMF,
 49 this paper presents new theoretical results in matrix analysis and an optimization
 50 algorithm for ONMF.

51 The main contributions of this paper are summarized as follows.

- 52 • We prove that a matrix $A \in \mathbb{R}_+^{m \times n}$ has an ONMF if and only if there is a
 53 matrix $Z \in \mathbb{R}^{n \times (n-r)}$ with $0 < r < n$ such that $AZ = 0, Z^\top Z = I, ZZ^\top \leq I$.
- 54 • Based on the necessary and sufficient condition for A to have an ONMF, we
 55 propose an optimization model

$$56 \quad (1.4) \quad \min \frac{1}{2} \|AZ\|_F^2,$$

subject to $Z \in \Omega := \{Z \in \mathbb{R}^{n \times (n-r)} \mid Z^\top Z = I, ZZ^\top \leq I\}$.

57 Problem (1.4) has a nonempty and bounded solution set, and its optimal
 58 value is 0 if and only if A has an ONMF.

- 59 • We propose an augmented Lagrangian algorithm for solving (1.4) and prove
 60 the convergence of the algorithm to a stationary point. Moreover, we de-
 61 scribe a simple method to obtain C and B from outputs of the algorithm
 62 by spectral decomposition. Preliminary numerical results are presented us-
 63 ing synthetic and real-world datasets, which show our approach substantially
 64 improves some existing approaches in clustering.

65 The rest of the paper is organized as follows. In Section 2, we derive the sufficient
 66 and necessary conditions for a nonnegative matrix A to admit an ONMF and establish
 67 the error bound. Section 3 proposes an augmented Lagrangian (AL) algorithm to solve
 68 the model (1.4) and prove its convergence. In Section 4, we present numerical results
 69 to demonstrate the effectiveness of our approach and compare with other approaches.

70 **Notation** Let I represent the identity matrix of appropriate dimensions and
 71 $\mathbb{O}^{n \times r} = \{C \in \mathbb{R}^{n \times r} \mid C^\top C = I\}$. The null space of A is denoted by $\text{Null}(A)$. We use
 72 $\|\cdot\|$ to denote the Frobenious norm $\|\cdot\|_F$ for matrices.

73 **2. Characterization of ONMF.** This section delves into the theoretical frame-
 74 work for matrices that admit ONMF. We first present Lemma 2.1, which serves as a
 75 key tool in the proof of Theorem 2.2, where we provide the necessary and sufficient
 76 conditions for a nonnegative matrix to admit an ONMF. At the end of this section,
 77 we give a perturbation error bound. These results provide a solid foundation for the
 78 development of algorithms, which are further explored in the next section.

79 LEMMA 2.1. Let $Q \in \mathbb{O}^{n \times r}$. If $QQ^\top \in \mathbb{R}_+^{n \times n}$, then there exists $C \in \mathbb{O}_+^{n \times r}$ such
80 that

$$81 \quad (2.1) \quad CC^\top = QQ^\top.$$

82 *Proof.* We prove this lemma by mathematical induction on n .

83 If $n = 1$, then $QQ^\top = 1$ and $C = 1$.

84 Assume that this lemma holds for each $n \in \{1, \dots, k\}$. We will prove it for
85 $n = k + 1$. Let $S = QQ^\top \in \mathbb{R}^{(k+1) \times (k+1)}$. The eigenvalues of S are either 1 or 0.

86 Case 1: If S is irreducible, the Perron-Frobenius theorem [2] implies that 1 is a
87 simple eigenvalue of S with a positive eigenvector v , while all other eigenvalues are
88 0. Note that $\text{rank}(S) = 1$, which implies $r = 1$. Let $C = \frac{v}{\|v\|_2} \in \mathbb{O}_+^{(k+1) \times 1}$ be the
89 normalized eigenvector. Then, S can be expressed as:

$$90 \quad S = \begin{bmatrix} C & U \end{bmatrix} \text{diag}(1, 0, \dots, 0) \begin{bmatrix} C & U \end{bmatrix}^\top = CC^\top,$$

91 where $U \in \mathbb{O}^{(k+1) \times k}$ is the orthogonal complement of C .

92 Case 2: Suppose that S is reducible. By the definition of the reducible matrix,
93 there is a permutation matrix P such that

$$94 \quad (2.2) \quad P^\top SP = \begin{bmatrix} K_1 & K_{12} \\ 0 & K_2 \end{bmatrix},$$

95 where K_i is a square matrix of order s_i and $s_1 + s_2 = k + 1$. Notice that S is
96 symmetric, so we have $K_{12} = 0$. Observe that S is idempotent and nonnegative, it
97 follows that K_1 and K_2 are also symmetric, idempotent and nonnegative, and thus the
98 eigenvalues of each K_1 and K_2 are 1 or 0. Denote $r_1 = \text{rank}(K_1)$ and $r_2 = \text{rank}(K_2)$
99 with $r = r_1 + r_2$. By the spectral decomposition, it is evident that there exist $Q_1 \in$
100 $\mathbb{O}^{s_1 \times r_1}, Q_2 \in \mathbb{O}^{s_2 \times r_2}$, such that $K_1 = Q_1 Q_1^\top$ and $K_2 = Q_2 Q_2^\top$.

If $r_1 r_2 > 0$, then the induction hypothesis ensures the existence of $C_1 \in \mathbb{O}_+^{s_1 \times r_1}$
and $C_2 \in \mathbb{O}_+^{s_2 \times r_2}$ such that

$$C_1 C_1^\top = K_1, \quad C_2 C_2^\top = K_2.$$

101 Hence this lemma holds with

$$102 \quad C = P \begin{bmatrix} C_1 & 0 \\ 0 & C_2 \end{bmatrix} \in \mathbb{O}_+^{(k+1) \times r}.$$

103 If $r_1 = r$ and $r_2 = 0$, then we can justify this lemma by taking

$$104 \quad C = P \begin{bmatrix} C_1 \\ 0 \end{bmatrix} \in \mathbb{O}_+^{(k+1) \times r}.$$

105 The similar result holds for $r_1 = 0$ and $r_2 = r$. The proof is complete. \square

106 THEOREM 2.2. A matrix $A \in \mathbb{R}_+^{m \times n}$ has an ONMF if and only if there exists
107 $Z \in \mathbb{R}^{n \times (n-r)}$ with $0 < r < n$ such that

$$108 \quad (2.3) \quad AZ = 0, \quad Z^\top Z = I_{n-r}, \quad ZZ^\top \leq I.$$

109 *Proof.* \Leftarrow : Assume there exists $Z \in \mathbb{R}^{n \times (n-r)}$ such that (2.3) holds. Let $Q \in \mathbb{O}^{n \times r}$
110 be the orthogonal complement of Z , then it yields $QQ^\top = I - ZZ^\top$. Since $ZZ^\top \leq I$,
111 it follows that $QQ^\top \in \mathbb{R}_+^{n \times n}$.

112 Following from Lemma 2.1, there exists $C \in \mathbb{O}_+^{n \times r}$ such that $CC^\top = QQ^\top$. Let
 113 $B = AC$. Then $B \in \mathbb{R}_+^{m \times r}$. Since $Z \in \text{Null}(A)$, we have:

$$114 \quad BC^\top = ACC^\top = AQQ^\top = A(I - ZZ^\top) = A,$$

115 which implies that A admits an ONMF.

116 \Rightarrow : If A has an ONMF, then there exist a matrix $B \in \mathbb{R}_+^{m \times r}$ and a matrix
 117 $C \in \mathbb{O}_+^{n \times r}$ such that $A = BC^\top$. Take $Z \in \mathbb{R}^{n \times (n-r)}$ such that $[C, Z]$ is orthogonal.
 118 Then $AZ = BC^\top Z = 0$, $Z^\top Z = I_{n-r}$ and $ZZ^\top = I - CC^\top \leq I$. \square

119 *Remark 2.3.* In Theorem 2.2, the matrix Z satisfying (2.3) may not be unique,
 120 but ZZ^\top is unique if $\text{rank}(A) = r$, as it is the orthogonal projection onto the null
 121 space of A .

122 If $AZ = 0$ holds only for the zero matrix Z , then $A \in \mathbb{R}_+^{m \times n}$ has full column rank.
 123 It is straightforward to verify that A admits an ONMF of the form $A = BC^\top$, where
 124 $B = A$ and $C = I$. For the remainder of this paper, we only consider $r < n$ unless
 125 explicitly stated otherwise.

126 **Example 2.1** Consider a matrix $A = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$, which admits an ONMF

127 with $B = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$ and $C^\top = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 1 \end{bmatrix}$. Both vectors $Z_1 = \left[\frac{1}{\sqrt{2}} \quad -\frac{1}{\sqrt{2}} \quad 0\right]^\top$

128 and $Z_2 = \left[-\frac{1}{\sqrt{2}} \quad \frac{1}{\sqrt{2}} \quad 0\right]^\top$ satisfy (2.3) in Theorem 2.2. Although $Z_1 \neq Z_2$, we have
 129 $Z_1 Z_1^\top = Z_2 Z_2^\top$.

130

131 Suppose A is the ground truth matrix, and \tilde{A} is an observed matrix from the data
 132 with noise. Let Z and \tilde{Z} be solutions of problem (1.4) with A and \tilde{A} , respectively.
 133 From Theorem 2.2, $AZ = 0$ and $\tilde{A}\tilde{Z} = 0$ if and only if A and \tilde{A} have ONMF. The
 134 following theorem gives the error metric $\Delta M = \tilde{Z}\tilde{Z}^\top - ZZ^\top$. This error metric enables
 135 us to assess the accuracy of \tilde{Z} in approximating Z , thereby providing a rigorous means
 136 to evaluate the performance of algorithms for ONMF in preserving the structural
 137 properties of the original matrix A .

138 **THEOREM 2.4.** Let \tilde{Z} be a solution of problem (1.4) with $\tilde{A} \in \mathbb{R}_+^{m \times n}$. Assume
 139 that A has an ONMF with $A = BC^\top$, $B \in \mathbb{R}_+^{n \times r}$ and $C \in \mathbb{O}_+^{n \times r}$. Set $Z \in \mathbb{R}^{n \times (n-r)}$
 140 such that $[C, Z]$ is orthogonal, and let $M = ZZ^\top$, $\tilde{M} = \tilde{Z}\tilde{Z}^\top$. Define the difference
 141 $\Delta A = \tilde{A} - A$ and $\Delta M = \tilde{M} - M$. Then, the following bound holds:

$$142 \quad (2.4) \quad \|\Delta M\| \leq 2\sqrt{2}\|C^\top A^\dagger\|_2 \|\Delta A\|.$$

143 In addition, if \tilde{A} has an ONMF, then

$$144 \quad (2.5) \quad \|\Delta M\| \leq \sqrt{2}\|C^\top A^\dagger \Delta A \tilde{M}\|.$$

145 *Proof.* Since \tilde{Z} is in the feasible set Ω of problem (1.4), there is $\tilde{C} \in \mathbb{O}_+^{n \times r}$
 146 such that $\tilde{M} = I - \tilde{C}\tilde{C}^\top$ by Lemma 2.1. Let $U = [C, C']$ and $V = [\tilde{C}, \tilde{C}']$ be
 147 $n \times n$ orthogonal matrices. Noting $M = I - CC^\top$. It is easy to see that $\|\Delta M\|^2 =$
 148 $\|\Delta MU\|^2 = \|\Delta MC\|^2 + \|\Delta MC'\|^2$ and $\|\Delta MC'\| = \|\tilde{C}^\top C'\|$. By CS Decomposition
 149 Theorem [18], we have $\|C'^\top C\| = \|\tilde{C}^\top C'\|$. Note that

$$150 \quad \tilde{C}'^\top \Delta MC = \tilde{C}'^\top (CC^\top - \tilde{C}\tilde{C}^\top)C = \tilde{C}'^\top C,$$

151 we have

$$152 \quad \|\tilde{C}'^\top C\| = \|\tilde{C}'^\top \Delta MC\|.$$

Therefore,

$$\|\Delta M\|^2 = \|\Delta MC\|^2 + \|\tilde{C}'^\top \Delta MC\|^2 \leq 2\|\Delta MC\|^2,$$

which implies that

$$\|\Delta M\| \leq \sqrt{2}\|\Delta MC\| = \sqrt{2}\|C^\top \Delta M\|.$$

153 Since A has an ONMF, we have $AM = AZZ^\top = 0$. Using $\tilde{A}(I - \tilde{C}\tilde{C}^\top) = \tilde{A}\tilde{Z}\tilde{Z}^\top$, we
154 have

$$155 \quad A\Delta M = A(\tilde{M} - M) = A\tilde{M}.$$

156 Hence, we obtain

$$157 \quad C^\top \Delta M = B^\dagger A\Delta M = B^\dagger A\tilde{M},$$

158 which yields

$$159 \quad (2.6) \quad \|\Delta M\| \leq \sqrt{2}\|B^\dagger A\tilde{M}\|.$$

160 Note that C is column orthonormal, by the properties of Moore-Penrose inverse, we
161 have $C^\dagger = C^\top$, and thus $A^\dagger = (C^\dagger)^\top B^\dagger = CB^\dagger$, which leads to $B^\dagger = C^\top A^\dagger$. In view
162 of this and equation (2.6), we obtain

$$\begin{aligned} 163 \quad \|\Delta M\| &\leq \sqrt{2}\|C^\top A^\dagger A\tilde{M}\| \leq \sqrt{2}\|C^\top A^\dagger\|_2(\|\Delta A\tilde{M}\| + \|\tilde{A}\tilde{M}\|) \\ 164 &\leq \sqrt{2}\|C^\top A^\dagger\|_2(\|\Delta A\| + \|\tilde{A}\tilde{Z}\|) \\ 165 &\leq \sqrt{2}\|C^\top A^\dagger\|_2(\|\Delta A\| + \|\tilde{A}Z\|) \\ 166 &= \sqrt{2}\|C^\top A^\dagger\|_2(\|\Delta A\| + \|\Delta AZ\|) \\ 167 &= \sqrt{2}\|C^\top A^\dagger\|_2(\|\Delta A\| + \|\Delta A\|), \end{aligned}$$

168 where the fourth inequality uses that \tilde{Z} is a solution of problem (1.4) with $\tilde{A} \in \mathbb{R}_+^{m \times n}$.
169 Hence the bound (2.4) is obtained.

170 In addition, if \tilde{A} has an ONMF, then $\tilde{A}\tilde{M} = 0$, which implies $A\tilde{M} = (\tilde{A} - \Delta A)\tilde{M} =$
171 $-\Delta A\tilde{M}$, thus it obtains the bound (2.5). This completes the proof of the theorem. \square

172 **3. An augmented Lagrangian (AL) algorithm.** In contrast to previous al-
173 gorithms [8, 20] to solve (1.1) or solve a convex approximation problem (1.3) as in [19],
174 our approach is to solve (1.4) grounded in Theorem 2.2. The feasible region Ω of (1.4)
175 guarantees that $Z \in \mathbb{O}^{n \times (n-r)}$ and $ZZ^\top \leq I$, as required by the conditions in Theo-
176 rem 2.2. Moreover, it is straightforward to verify that Ω is non-empty, since $I_{n \times (n-r)}$
177 is a valid point, where $I_{n \times (n-r)}$ denotes the matrix formed by the first $(n-r)$ columns
178 of the identity matrix I .

Since the feasible region is compact and the objective function is continuous,
problem (1.4) admits a global minimum by the Weierstrass extreme value theorem.
To facilitate the application of the AL algorithm, we reformulate problem (1.4) by
using auxiliary variables. Let

$$\mathcal{W} = (X, Z) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times (n-r)}, \quad \mathbb{F} := \{\mathcal{W} \mid X \geq 0, Z^\top Z = I\}$$

179 and

$$180 \quad (3.1) \quad c(\mathcal{W}) = X - I + ZZ^\top.$$

181 Problem (1.4) can be reformulated as:

$$182 \quad (3.2) \quad \min_{\mathcal{W} \in \mathbb{F}} \frac{1}{2} \|AZ\|^2, \quad \text{subject to } c(\mathcal{W}) = 0.$$

183 For any given penalty parameter $\rho > 0$ and Lagrange multiplier $\Lambda \in \mathbb{R}^{n \times n}$, the
184 AL function for problem (3.2) is given by

$$185 \quad (3.3) \quad \mathcal{L}(\mathcal{W}, \Lambda, \rho) := \frac{1}{2} \|AZ\|^2 + \frac{1}{2\rho} (\|\Lambda + \rho c(\mathcal{W})\|^2 - \|\Lambda\|^2).$$

186 At each outer iteration, for fixed Λ, ρ , we solve the following AL subproblem:

$$187 \quad (3.4) \quad \min_{\mathcal{W} \in \mathbb{F}} \mathcal{L}(\mathcal{W}, \Lambda, \rho).$$

188 It can be shown that the following condition is a necessary condition [10, 23] for a
189 local minimizer of problem (3.4).

$$190 \quad (3.5) \quad H(\mathcal{W}, \Lambda, \rho) := \begin{pmatrix} (I - ZZ^\top) \left[A^\top A - 2\rho(I - X - \frac{1}{\rho}\Lambda) \right] Z \\ Z^\top Z - I \\ X - \left(I - ZZ^\top - \frac{1}{\rho}\Lambda \right)_+ \end{pmatrix} = 0,$$

191 where $(\cdot)_+ = \max\{\cdot, 0\}$ is the projection onto the nonnegative space. Let $\mathcal{W}^{\text{feas}}$
192 be an feasible point of problem (3.2). For instance, $\mathcal{W}^{\text{feas}} = (X^{\text{feas}}, Z^{\text{feas}})$, where
193 $Z^{\text{feas}} = I_{n \times (n-r)}$ is the first $n - r$ columns of the $n \times n$ identity matrix I , and
194 $X^{\text{feas}} = I - Z^{\text{feas}}(Z^{\text{feas}})^\top$.

195 Following the approach in [5], we present the AL algorithm for solving problem
196 (3.2), as detailed in Algorithm 3.1.

197 **3.1. Convergence Analysis.** To establish the convergence properties of Algo-
198 rithm 3.1, we first define a stationary point for problem (3.2). The convergence results
199 for Algorithm 3.1 are then presented in Theorem 3.2.

200 **DEFINITION 3.1 (Stationary point).** *Let $\Omega_w = \{\mathcal{W} \in \mathbb{F} \mid c(\mathcal{W}) = 0\}$ be the*
201 *feasible region of problem (3.2), we say $\mathcal{W}^* \in \Omega_w$ is a stationary point of (3.2) if*
202 *there exists $\Lambda^* \in \mathbb{R}^{n \times n}$ such that*

$$203 \quad (3.11) \quad (I - Z^*(Z^*)^\top) \left[\frac{1}{2} \nabla \|AZ\|^2 \Big|_{Z^*} + \nabla_Z \langle c(\mathcal{W}), \Lambda^* \rangle \Big|_{Z^*} \right] = 0.$$

204 If \mathcal{W}^* is a local minimizer of (3.2), then it is a stationary point of (3.2). See [3, 10].

205 **THEOREM 3.2.** *Let $\{\mathcal{W}^k\}$ be the sequence generated by Algorithm 3.1 and assume*
206 *that $\{\mathcal{W}^k\}$ is bounded¹. Then, the following statements hold:*

- 207 (i) $c(\mathcal{W}^k) \rightarrow 0$ as $k \rightarrow \infty$, where $c(\cdot)$ is defined in (3.1).
208 (ii) Any accumulation point \mathcal{W}^* of $\{\mathcal{W}^k\}$ is a stationary point of (3.2).

209 *Proof.* To prove statement (i), we consider the following two separate cases.

- 210 (i) Case (a): $\{\rho_k\}$ is bounded.

¹The boundness of $\{\mathcal{W}^k\}$ is proved in Theorem 3.4

Algorithm 3.1 Augmented Lagrangian Method for Model (3.2)

Input: Observed data $A \in \mathbb{R}_+^{m \times n}$, rank r , ρ_{\max} , tolerance ε , max iteration it_{\max} .

- 1: Choose $\Lambda^0 = \mathbf{0}_{n \times n} \in \mathbb{R}^{n \times n}$, $\rho_0 > 0$, $\gamma \in (1, \infty)$, $\tau \in (0, 1)$, $\eta \in (0, 1)$, a positive sequence $\{\epsilon_k\}$ with $\lim_{k \rightarrow \infty} \epsilon_k = 0$, and a constant

$$(3.6) \quad \Gamma \geq \max\{\|AZ^{\text{feas}}\|, \mathcal{L}(\mathcal{W}^{\text{feas}}, \Lambda^0, \rho_0)\}.$$

- 2: **for** $k = 0, 1, \dots, it_{\max}$ **do**

- 3: Solve problem (3.4) with $\Lambda = \Lambda^k$ and $\rho = \rho_k$ to find an approximate stationary point $\mathcal{W}^k \in \mathbb{F}$ of (3.4) such that

$$(3.7) \quad \|H(\mathcal{W}^k, \Lambda^k, \rho_k)\| \leq \epsilon_k, \quad \mathcal{L}(\mathcal{W}^k, \Lambda^k, \rho_k) \leq \Gamma.$$

- 4: Set

$$(3.8) \quad \Lambda^{k+1} = \Lambda^k + \rho_k c(\mathcal{W}^k).$$

- 5: **if** $k > 0$ and

$$(3.9) \quad \|c(\mathcal{W}^k)\| \leq \eta \|c(\mathcal{W}^{k-1})\|$$

then

- 6: Set $\rho_{k+1} = \rho_k$

- 7: **else**

- 8: Set

$$(3.10) \quad \rho_{k+1} = \max\{\gamma \rho_k, \|\Lambda^{k+1}\|^{1+\tau}\}.$$

- 9: **end if**

- 10: **if** $\|c(\mathcal{W}^k)\| < \varepsilon$ **or** $\rho_{k+1} > \rho_{\max}$ **then**

- 11: **break**

- 12: **end if**

- 13: **end for**

Output: \mathcal{W}^k

211 From the updating scheme in Algorithm 3.1 and $\gamma > 1$, we can see that $\{\rho_k\}$ is
 212 updated by (3.10) only for finite times. It thus implies that there is $k_0 > 0$ such that
 213 (3.9) holds for all $k \geq k_0$, that is,

$$214 \quad \|c(\mathcal{W}^k)\| \leq \eta \|c(\mathcal{W}^{k-1})\|, \quad \forall k \geq k_0.$$

215 Together with $\eta \in (0, 1)$, we obtain

$$216 \quad (3.12) \quad \lim_{k \rightarrow \infty} \|c(\mathcal{W}^k)\| = 0,$$

217 which clearly implies that statement (i) holds.

218 Case (b): $\{\rho_k\}$ is unbounded.

219 By the updating scheme in Algorithm 3.1, we can observe that $\{\rho_k\}$ must be
 220 updated by (3.10) for infinite times. Let $\{\rho_{j_1}, \rho_{j_2}, \dots\}$ denote all elements in $\{\rho_k\}$
 221 that are updated by (3.10) and $\mathbb{J} = \{j_1, j_2, \dots\}$ is arranged in increasing order. It
 222 then follows that $\rho_{j_\ell} \rightarrow \infty$ as $\ell \rightarrow \infty$ and

$$223 \quad (3.13) \quad \rho_i = \rho_{j_\ell}, \quad j_\ell \leq i < j_{\ell+1}, \quad \ell \geq 1,$$

$$224 \quad (3.14) \quad \rho_{j_\ell} = \max\{\gamma \rho_{j_{\ell-1}}, \|\Lambda^{j_\ell}\|^{1+\tau}\}, \quad \ell \geq 1.$$

226 Let $\underline{j}(k) := \max\{j \in \mathbb{J} \mid k \geq j\}$ for every $k \geq j_1$.
 227 Next we claim that for every $k \geq j_1$,

$$228 \quad (3.15) \quad \frac{\|\Lambda^k\|}{\rho_k} \leq \frac{\|\Lambda^{\underline{j}(k)}\|}{\rho_{\underline{j}(k)}} + \sum_{i=1}^{k-\underline{j}(k)} \|c(\mathcal{W}^{k-i})\|.$$

229 In fact, $k \geq \underline{j}(k)$ by the definition of $\underline{j}(k)$. Clearly, (3.15) holds with equality when
 230 $k = \underline{j}(k)$. Now suppose $k > \underline{j}(k)$, we have $\rho_{k-i} = \rho_k$ for $1 \leq i \leq k - \underline{j}(k)$. In view of
 231 this and (3.8), it follows

$$232 \quad (3.16) \quad \begin{aligned} \frac{\|\Lambda^k\|}{\rho_k} &= \frac{\|\Lambda^k\|}{\rho_{k-1}} \leq \frac{\|\Lambda^{k-1}\|}{\rho_{k-1}} + \|c(\mathcal{W}^{k-1})\| \\ &\leq \cdots \leq \frac{\|\Lambda^{\underline{j}(k)}\|}{\rho_{\underline{j}(k)}} + \sum_{i=1}^{k-\underline{j}(k)} \|c(\mathcal{W}^{k-i})\|, \end{aligned}$$

233 which proves (3.15).

234 We then show that for every $k \geq j_1$,

$$235 \quad (3.17) \quad \frac{\|\Lambda^k\|}{\rho_k} \leq \frac{\|\Lambda^{\underline{j}(k)}\|}{\rho_{\underline{j}(k)}} + \frac{1}{1-\eta} \|c(\mathcal{W}^{\underline{j}(k)})\|.$$

Indeed, it follows from the definition of $\underline{j}(k)$ and the updating scheme of $\{\rho_\ell\}$ that

$$\|c(\mathcal{W}^{k-i})\| \leq \eta \|c(\mathcal{W}^{k-i-1})\|, \quad \text{for } 1 \leq i < k - \underline{j}(k),$$

236 which leads to

$$237 \quad (3.18) \quad \|c(\mathcal{W}^{k-i})\| \leq \eta^{k-\underline{j}(k)-i} \|c(\mathcal{W}^{\underline{j}(k)})\|, \quad \text{for } 1 \leq i \leq k - \underline{j}(k).$$

Then inequality (3.17) follows from (3.15) and the fact that

$$\sum_{\ell=0}^{k-\underline{j}(k)-1} \eta^\ell \leq \sum_{\ell=0}^{\infty} \eta^\ell = \frac{1}{1-\eta}.$$

238 In the following, we will show that

$$239 \quad (3.19) \quad \lim_{k \rightarrow \infty} \|\Lambda^{\underline{j}(k)}\| / \rho_{\underline{j}(k)} = 0,$$

$$240 \quad (3.20) \quad \lim_{k \rightarrow \infty} \|c(\mathcal{W}^{\underline{j}(k)})\| = 0.$$

241 Indeed, recall that $\rho_{j_\ell} \rightarrow \infty$ as $\ell \rightarrow \infty$. It thus follows that $\rho_{\underline{j}(k)} \rightarrow \infty$ as $k \rightarrow \infty$. By
 242 (3.10) and the definition of $\underline{j}(k)$, one has

$$243 \quad (3.21) \quad \rho_{\underline{j}(k)} = \max \left\{ \gamma \rho_{\underline{j}(k)-1}, \|\Lambda^{\underline{j}(k)}\|^{1+\tau} \right\}, \quad k \geq j_1.$$

244 This yields

$$245 \quad \|\Lambda^{\underline{j}(k)}\| \leq \rho_{\underline{j}(k)}^{\frac{1}{1+\tau}},$$

246 which implies that

$$247 \quad \|\Lambda^{\underline{j}(k)}\| / \rho_{\underline{j}(k)} \leq (\rho_{\underline{j}(k)})^{-\frac{\tau}{1+\tau}}, \quad k \geq j_1.$$

248 The relations (3.19) then follows from this and $\rho_{\underline{j}(k)} \rightarrow \infty$ as $k \rightarrow \infty$. Further, by the
 249 second inequality in (3.7) and the definition of the AL function (3.3), one has

$$250 \quad \frac{1}{2} \|AZ^{\underline{j}(k)}\|^2 + \frac{\|\Lambda^{\underline{j}(k)} + \rho_{\underline{j}(k)} c(\mathcal{W}^{\underline{j}(k)})\|^2 - \|\Lambda^{\underline{j}(k)}\|^2}{2\rho_{\underline{j}(k)}} \leq \Gamma,$$

251 which leads to

$$252 \quad (3.22) \quad \left\| c(\mathcal{W}^{\underline{j}(k)}) + \frac{\Lambda^{\underline{j}(k)}}{\rho_{\underline{j}(k)}} \right\|^2 \leq \frac{1}{\rho_{\underline{j}(k)}} [2\Gamma - \|AZ^{\underline{j}(k)}\|^2] + \frac{\|\Lambda^{\underline{j}(k)}\|^2}{\rho_{\underline{j}(k)}^2}.$$

253 Using this, (3.19), the boundedness of $\{\mathcal{W}^k\}$ and $\rho_{\underline{j}(k)} \rightarrow \infty$ as $k \rightarrow \infty$, one can see
 254 that

$$255 \quad (3.23) \quad \lim_{k \rightarrow \infty} \left\| c(\mathcal{W}^{\underline{j}(k)}) + \frac{\Lambda^{\underline{j}(k)}}{\rho_{\underline{j}(k)}} \right\| = 0,$$

256 which together with (3.19) implies that the relation (3.20) holds. In view of (3.17),
 257 (3.19) and (3.20), we conclude that

$$258 \quad (3.24) \quad \lim_{k \rightarrow \infty} \|\Lambda^k\|/\rho_k = 0.$$

259 By the same argument used to establish (3.22) for any k , we obtain

$$260 \quad (3.25) \quad \left\| c(\mathcal{W}^k) + \frac{\Lambda^k}{\rho_k} \right\|^2 \leq \frac{1}{\rho_k} [2\Gamma - \|AZ^k\|^2] + \frac{\|\Lambda^k\|^2}{\rho_k^2}.$$

261 Together with (3.24), the boundedness of $\{\mathcal{W}^k\}$, and $\rho_k \rightarrow \infty$ as $k \rightarrow \infty$, we conclude
 262 that $c(\mathcal{W}^k) \rightarrow 0$ as $k \rightarrow \infty$. This completes the proof of statement (i).

263 (ii) Since \mathcal{W}^* is an accumulation point and by the continuity of $c(\cdot)$, we have
 264 $c(\mathcal{W}^*) = 0$. Moreover, given that $\{\mathcal{W}^k\} \subseteq \mathbb{F}$ and is bounded, the closedness of \mathbb{F}
 265 ensures that $\mathcal{W}^* \in \mathbb{F}$, which implies that \mathcal{W}^* is a feasible point.

266 Let $\{\mathcal{W}^k\}_{\mathcal{K}} \rightarrow \mathcal{W}^*$ for some subsequence \mathcal{K} . To show that \mathcal{W}^* is a stationary
 267 point, we first establish that

$$268 \quad (3.26) \quad \|[I - Z^k(Z^k)^\top](A^\top A + 2\Lambda^{k+1})Z^k\| \rightarrow 0 \quad \text{as } k \rightarrow \infty, k \in \mathcal{K}.$$

269 Recall that the first relation in (3.7) holds for any positive sequence $\{\epsilon_k\}$. Letting
 270 $\epsilon_k \rightarrow 0$ as $k \rightarrow \infty$, it follows that $\{\epsilon_k\}_{\mathcal{K}} \rightarrow 0$ as well. Since $\{\mathcal{W}^k\} \subset \mathbb{F}$, one has
 271 $(I - Z^k(Z^k)^\top)Z^k = 0$. From (3.8), we have

$$\begin{aligned} & [I - Z^k(Z^k)^\top][A^\top A - 2\rho_k(I - X^k - \frac{\Lambda^k}{\rho_k})]Z^k \\ 272 \quad (3.27) \quad & = [I - Z^k(Z^k)^\top][A^\top A + 2\rho_k(\frac{\Lambda^k}{\rho_k} + X^k - I + Z^k(Z^k)^\top)]Z^k \\ & = [I - Z^k(Z^k)^\top]\{A^\top A + 2[\Lambda^k + \rho_k c(\mathcal{W}^k)]\}Z^k \\ & = [I - Z^k(Z^k)^\top](A^\top A + 2\Lambda^{k+1})Z^k. \end{aligned}$$

273 Then, we obtain (3.26) by the first inequality in (3.7).

274 Equation (3.26) implies that

$$275 \quad (3.28) \quad [I - Z^k(Z^k)^\top](A^\top A + 2\Lambda^{k+1})Z^k \rightarrow 0, \quad \text{as } k \in \mathcal{K} \rightarrow \infty.$$

276 Note that $\text{rank}(Z^k) = n - r$. Let $Q^k \in \mathbb{O}^{n \times r}$ be the orthogonal complement of Z ,
277 then $I - Z^k(Z^k)^\top = Q^k(Q^k)^\top$. Substituting this into equation (3.28), we obtain

$$278 \quad (3.29) \quad Q^k(Q^k)^\top(A^\top A + 2\Lambda^{k+1})Z^k = Q^k(Q^k)^\top A^\top A Z^k + 2Q^k \bar{\Lambda}^k \rightarrow 0,$$

279 as $k \rightarrow \infty, k \in \mathcal{K}$, where $\bar{\Lambda}^k = (Q^k)^\top \Lambda^{k+1} Z^k$.

280 Given that Q^k and Z^k are column orthonormal, it follows that Q^k and Z^k are
281 bounded. Multiplying $(Q^k)^\top$ to the left of (3.29) shows the boundedness of $\bar{\Lambda}^k$. Thus
282 there exists $\mathcal{K}_1 \subset \mathcal{K}$ such that $\{Q^k\}_{\mathcal{K}_1} \rightarrow Q^*, \{Z^k\}_{\mathcal{K}_1} \rightarrow Z^*$ and $\{\bar{\Lambda}^k\}_{\mathcal{K}_1} \rightarrow \bar{\Lambda}^*$, where
283 Q^*, Z^* are column orthonormal. Observe that

$$284 \quad Q^*(Q^*)^\top = \lim_{k \in \mathcal{K}_1 \rightarrow \infty} Q^k(Q^k)^\top = \lim_{k \in \mathcal{K}_1 \rightarrow \infty} (I - Z^k(Z^k)^\top) = I - Z^*(Z^*)^\top.$$

285 It follows that

$$\begin{aligned} & \lim_{k \in \mathcal{K}_1 \rightarrow \infty} Q^k(Q^k)^\top A^\top A Z^k + 2Q^k \bar{\Lambda}^k \\ &= Q^*(Q^*)^\top A^\top A Z^* + 2Q^* \bar{\Lambda}^* \\ 286 \quad (3.30) \quad &= Q^*(Q^*)^\top A^\top A Z^* + 2Q^*(Q^*)^\top Q^* \bar{\Lambda}^* (Z^*)^\top Z^* \\ &= Q^*(Q^*)^\top (A^\top A + 2\Lambda^*) Z^* \\ &= (I - Z^*(Z^*)^\top) (A^\top A + 2\Lambda^*) Z^*, \end{aligned}$$

287 where $\Lambda^* = Q^* \bar{\Lambda}^* (Z^*)^\top$. Hence $[\nabla_Z \langle c(\mathcal{W}), \Lambda^* \rangle]_{Z^*} = 2\Lambda^* Z^*$. We prove that \mathcal{W}^* is a
288 stationary point. \square

289 **3.2. A block coordinate descent algorithm for subproblem (3.4).** In this
290 subsection, we propose a block coordinate descent (BCD) algorithm for solving the
291 subproblem (3.4) with fixed $\Lambda = \Lambda^k$ and $\rho = \rho_k$. The goal is to find an approximate
292 stationary point $\mathcal{W}^k \in \mathbb{F}$ that satisfies the conditions in (3.7).

293 For simplicity, we use \mathcal{W}_t for the iterations \mathcal{W}_t^k of the BCD algorithm. To guar-
294 antee the validity of the second relation in (3.7), one can select the initial point of the
295 BCD method as follows:

$$296 \quad (3.31) \quad \mathcal{W}_0 = \begin{cases} \mathcal{W}^{\text{feas}} & \text{if } k = 0 \text{ or } \mathcal{L}(\mathcal{W}^{k-1}, \Lambda, \rho) > \Gamma \\ \mathcal{W}^{k-1} & \text{otherwise,} \end{cases}$$

297 where $\mathcal{W}^{\text{feas}}$ denotes a feasible point for problem (3.2). For instance, let $Z = I_{n \times (n-r)}$
298 represent the first $n - r$ columns of the identity matrix I . Then, set $X = I - ZZ^\top$.

299 The subproblems for Z and X are solved alternately by fixing the other variable
300 at each step. Specifically, during the t -th iteration of the BCD algorithm, the updates
301 for Z and X are carried out as follows:

$$\begin{aligned} & Z_{t+1} \in \underset{Z^\top Z = I}{\text{argmin}} \mathcal{L}(X_t, Z, \Lambda, \rho) \\ 302 \quad (3.32) \quad &= \underset{Z^\top Z = I}{\text{argmin}} \text{tr} \left\{ Z^\top \left[A^\top A - 2\rho \left(I - X_t - \frac{\Lambda}{\rho} \right) \right] Z \right\}, \end{aligned}$$

$$\begin{aligned} & X_{t+1} = \underset{X \geq 0}{\text{argmin}} \mathcal{L}(X, Z_{t+1}, \Lambda, \rho) \\ 303 \quad (3.33) \quad &= \underset{X \geq 0}{\text{argmin}} \frac{1}{2} \|X - I + Z_{t+1} Z_{t+1}^\top + \frac{\Lambda}{\rho}\|^2. \end{aligned}$$

304 To find a solution of problem (3.32), we introduce the following lemma.

305 LEMMA 3.3 (Corollary 4.3.39 in [11]). *Let $S \in \mathbb{R}^{n \times n}$ be symmetric, and suppose*
 306 *$1 \leq \ell \leq n$. Then*

307 (3.34)
$$\lambda_1(S) + \cdots + \lambda_\ell(S) = \min_{Z \in \mathbb{R}^{n \times \ell}, Z^\top Z = I} \text{tr}(Z^\top S Z),$$

308 *where $\lambda_1(S) \leq \cdots \leq \lambda_n(S)$ are eigenvalues of S , and the minimum is*
 309 *achieved at a matrix Z whose columns are orthonormal eigenvectors associated with*
 310 *$\lambda_1, \dots, \lambda_\ell$.*

311 Lemma 3.3 gives a closed form solution of (3.32) by setting $S = A^\top A - 2\rho(I - X_t - \frac{\Lambda}{\rho})$.
 312 That is, Z_{t+1} can be constructed by the orthonormal eigenvectors associated with the
 313 $n - r$ smallest eigenvalues of S .

314 The subproblem (3.33) also has a closed-form solution, which are given by:

315 (3.35)
$$X_{t+1} = \max\{I - Z_{t+1} Z_{t+1}^\top - \frac{\Lambda}{\rho}, 0\}.$$

316 The BCD algorithm for subproblem (3.4) is summarized in Algorithm 3.2.

Algorithm 3.2 BCD algorithm to solve subproblem (3.4)

Input: Result from the $(k - 1)_{th}$ step of Algorithm 3.1 \mathcal{W}^{k-1} , $\Lambda = \Lambda_k$, $\rho = \rho_k$

- 1: Set tolerance ϵ_k , max iteration it_{max} .
- 2: Choose the initial point \mathcal{W}_0 by (3.31).
- 3: **for** $t = 0, 1, \dots, it_{max}$ **do**
- 4: Update Z_{t+1} by Lemma 3.3 with $S = A^\top A - 2\rho(I - X_t - \frac{\Lambda}{\rho})$.
- 5: Update X_{t+1} by using (3.35).
- 6: **if** $\|H(\mathcal{W}_{t+1}, \Lambda, \rho)\| < \epsilon_k$ **then**
- 7: **break**
- 8: **end if**
- 9: **end for**

Output: $\mathcal{W}^k = (X_{t+1}, Z_{t+1})$.

317 To ensure the boundedness of \mathcal{W}^k and that the conditions in (3.7) are satisfied
 318 within Algorithm 3.1, we demonstrate that these claims are achieved through the
 319 BCD algorithm outlined in Algorithm 3.2, as formalized in Theorem 3.4.

320 THEOREM 3.4. *Let $\{\mathcal{W}_t\}$ be a sequence generated by Algorithm 3.2. Then, the*
 321 *following statements hold:*

- 322 (i) *The sequence $\{\mathcal{W}_t\}$ is bounded.*
- 323 (ii) *For any $\epsilon_k > 0$, there is a positive integer \bar{t} , such that for any $t \geq \bar{t}$, we have*

324 (3.36)
$$\|H(\mathcal{W}_t, \Lambda, \rho)\| \leq \epsilon_k, \quad \mathcal{L}(\mathcal{W}_t, \Lambda, \rho) \leq \Gamma,$$

325 where Γ is defined in (3.6) within Algorithm 3.1.

326 *Proof.* (i) We first show $\Lambda = \Lambda^k \geq 0$, for any k . By using the update scheme of

327 Λ^k in (3.8) and X_{t+1} in (3.35), respectively, one has:

$$\begin{aligned}
\Lambda^{k+1} &= \rho_k \left(X^k - I + Z^k (Z^k)^\top + \frac{\Lambda^k}{\rho_k} \right) \\
(3.37) \quad &= \rho_k \left(\max\{I - Z^k (Z^k)^\top - \frac{\Lambda^k}{\rho_k}, 0\} - I + Z^k (Z^k)^\top + \frac{\Lambda^k}{\rho_k} \right) \\
&= \rho_k \max\{-I + Z^k (Z^k)^\top + \frac{\Lambda^k}{\rho_k}, 0\} \geq 0,
\end{aligned}$$

329 where the third equality follows $\max(a, 0) - a = \max(-a, 0)$. Thus $\Lambda^k \geq 0$, for any k .
330 Notice that Z_{t+1} is updated over the feasible set in (3.32), i.e., $Z_{t+1}^\top Z_{t+1} = I$, which
331 implies $\{Z_{t+1}\}$ is bounded.

332 Since $\frac{\Lambda}{\rho} \geq 0$, one can obtain

$$(3.38) \quad 0 \leq X_{t+1} = \max\{I - Z_{t+1} Z_{t+1}^\top - \frac{\Lambda}{\rho}, 0\} \leq \max\{I - Z_{t+1} Z_{t+1}^\top, 0\},$$

334 which implies $\{X_t\}$ is bounded. Consequently, we conclude that $\{\mathcal{W}_t\}$ is bounded.

335 Furthermore, since $\|\mathcal{W}_t\|_F^2 = \|Z_t\|_F^2 + \|X_t\|_F^2$, we obtain

$$(3.36) \quad \|\mathcal{W}_t\|_F^2 \leq (n - r) + 2(2n - r)^2,$$

337 where the bound is independent of both t and k . Noting that $\mathcal{W}^k = \lim_{t \rightarrow \infty} \mathcal{W}_t^k$, it
338 follows that the sequence $\{\mathcal{W}^k\}$ generated by Algorithm 3.1 is uniformly bounded.

339 (ii) From the scheme of Algorithm 3.2, it is easy to verify that for any $t \geq 0$,
(3.39)

$$(3.40) \quad -\frac{\|\Lambda\|^2}{2\rho} \leq \mathcal{L}(\mathcal{W}_{t+1}, \Lambda, \rho) \leq \mathcal{L}(Z_{t+1}, X_t, \Lambda, \rho) \leq \mathcal{L}(\mathcal{W}_t, \Lambda, \rho) \leq \mathcal{L}(\mathcal{W}_0, \Lambda, \rho) \leq \Gamma.$$

341 This shows that \mathcal{L} is bounded with respect to \mathcal{W}_t .

342 Now we prove the first inequality in (3.36). From the structure of Algorithm 3.2,
343 we have

$$(3.40) \quad \begin{pmatrix} (I - Z_{t+1} Z_{t+1}^\top) \left[A^\top A - 2\rho \left(I - X_t - \frac{\Lambda}{\rho} \right) \right] Z_{t+1} \\ Z_{t+1}^\top Z_{t+1} - I \\ X_{t+1} - \left(I - Z_{t+1} Z_{t+1}^\top - \frac{\Lambda}{\rho} \right)_+ \end{pmatrix} = 0.$$

Hence, we have

$$\|H(\mathcal{W}_{t+1}, \Lambda, \rho)\| \leq 2\rho \|X_t - X_{t+1}\|.$$

345 Since \mathcal{L} is strongly convex with respect to X , there is a constant α such that

$$\begin{aligned}
(3.46) \quad \|X_t - X_{t+1}\| &\leq \alpha [\mathcal{L}(Z_{t+1}, X_t, \Lambda, \rho) - \mathcal{L}(\mathcal{W}_{t+1}, \Lambda, \rho)] \\
(3.47) \quad &\leq \alpha [\mathcal{L}(\mathcal{W}_t, \Lambda, \rho) - \mathcal{L}(\mathcal{W}_{t+1}, \Lambda, \rho)].
\end{aligned}$$

348 Since $\mathcal{L}(\mathcal{W}_t, \Lambda, \rho)$ is decreasing and bounded below, one has $\mathcal{L}(\mathcal{W}_t, \Lambda, \rho)$ converges.
349 Thus $\mathcal{L}(\mathcal{W}_t, \Lambda, \rho) - \mathcal{L}(\mathcal{W}_{t+1}, \Lambda, \rho) \downarrow 0$, we can claim that there is $\bar{t} > 0$, such that for
350 $t > \bar{t}$, the first inequality in (3.36) holds. \square

351 Now, one can see $\mathcal{L}(\mathcal{W}^k, \Lambda, \rho) = \lim_{t \rightarrow \infty} \mathcal{L}(\mathcal{W}_t^{k-1}, \Lambda, \rho) \leq \Gamma$ which ensures that
352 the second condition in (3.7) holds at every iteration.

3.3. Deriving B and C from output Z . From Theorem 3.2, an output Z of Algorithm 3.1 is in $\Omega = \{Z \in \mathbb{R}^{n \times (n-r)} \mid Z^\top Z = I, ZZ^\top \leq I\}$. We can use $Z \in \Omega$ to obtain $B \in \mathbb{R}_+^{m \times r}$ and $C \in \mathbb{O}_+^{n \times r}$ for ONMF of A . This can be achieved using the adjacency graph \mathcal{G} that constructed from $K = I - ZZ^\top$ with r strongly connected components. The complete procedure is detailed in Algorithm 3.3, and its main ideas are outlined below:

- (a) Since the adjacency graph \mathcal{G} is sensitive to small perturbations in $K = I - ZZ^\top$, even negligible entries can connect otherwise disjoint subgraphs. To eliminate these artifacts, we use the bisection method [22] to threshold K by setting entries with absolute values below δ to zero, ensuring that the resulting graph \mathcal{G} consists of exactly r strongly connected components, as described in lines 3–15 of Algorithm 3.3.
- (b) Identify the r strongly connected components $\mathcal{G}_1, \dots, \mathcal{G}_r$ of \mathcal{G} [9], with the corresponding adjacency matrices by K_1, \dots, K_r , where $K_i \in \mathbb{R}^{s_i \times s_i}$ and $\sum_{i=1}^r s_i = n$. These components can be computed using the `conncomp` function in MATLAB, see line 16 of Algorithm 3.3.
- (c) Since \mathcal{G}_i is strongly connected, its adjacency matrix K_i is irreducible [4]. Combined with the symmetry and idempotence of K_i , there exists a positive vector u_i such that $K_i = u_i u_i^\top$, $i = 1, \dots, r$.
- (d) Let π_i be the vertex set of \mathcal{G}_i , we have $K_i = K(\pi_i, \pi_i)$, where $\pi_i \subseteq \{1, \dots, n\}$, $\cup_{i=1}^r \pi_i = \{1, \dots, n\}$, and $\pi_i \cap \pi_j = \emptyset$ for $i \neq j$. The i -th column of C is obtained by normalizing the first column of K_i : $C(\pi_i, i) = \frac{u_i}{\|u_i\|_2}$, and setting $C(\pi_i^C, i) = 0$, where π_i^C is the complement of π_i [19]. This ensures that $C \in \mathbb{O}_+^{n \times r}$. Set $B = AC \in \mathbb{R}_+^{m \times r}$, we then obtain both factors B and C , as implemented in lines 17–20 of Algorithm 3.3.

4. Numerical Experiments. In this section, we evaluate the performance of our approach on both synthetic and real-world datasets. We denote the proposed model in Algorithm 3.1 as NS-ONMF. To ensure completeness and reproducibility, we specify the parameter initialization strategy for NS-ONMF as follows. In particular, the factor Z is set as $Z = I_{n \times (n-r)}$, i.e., the first $(n-r)$ columns of the identity matrix I_n . The penalty parameter ρ_0 is selected from $\{10^{-10}, 10^{-6}, 10^{-4}\}$ depending on the dataset type, where smaller values are used for near-ONMF synthetic data to allow gradual constraint enforcement, while larger values are used for real data to promote faster convergence under noise. The decay factor η , which controls the growth of the penalty, is chosen from $\{0.99, 0.995, 0.999\}$ to ensure stable updates. The update factor γ is chosen from the interval $(1, 2)$, with a typical value of 1.2; the parameter τ is taken from $(0, 1)$ with $\tau = 0.1$. In Algorithm 3.2, the maximum number of inner iterations is set to 10 for synthetic data to allow more accurate subproblem solutions, while for large-scale or noisy real datasets, a single inexact update is used to reduce computational cost without compromising performance. We further let $\Gamma = \max\{\|AZ^{\text{feas}}\|, \mathcal{L}(\mathcal{W}^{\text{feas}}, \Lambda^0, \rho_0)\}$ and $\epsilon_k = 1/k$. All experiments were conducted in MATLAB 2023b on an Intel(R) Xeon(R) Gold 5320 CPU @ 2.20GHz with 512GB of RAM, and the complete source code is publicly available on GitHub at <https://github.com/Qilun-Luo/ONMF>.

4.1. Synthetic Data. We will first test our approach with synthetic data.

4.1.1. A has an ONMF. Given a matrix A , our approach aims to minimize $\|AZ\|^2$ over the set Ω . We compare the performance with other ONMF methods, including BiOR-NM3F [8] (Bi-Orthogonal 3-Factor NMF), MU-ONMF [7] (Multi-

Algorithm 3.3 ONMF of A for deriving $B \in \mathbb{R}_+^{m \times r}$ and $C \in \mathbb{O}_+^{n \times r}$

Input: Null space $Z \in \mathbb{R}^{n \times (n-r)}$ obtained from Algorithm 3.1, factor r .

- 1: Initialize $\delta_{min} = 10^{-10}$, $\delta_{max} = 10^{-1}$, $it_{max} = 100$.
 - 2: Set $K = I - ZZ^\top$, $\delta = (\delta_{min} + \delta_{max})/2$.
 - 3: **for** $k = 1, \dots, it_{max}$ **do**
 - 4: Set $\tilde{K} = K$, and let $\tilde{K}(\text{abs}(\tilde{K}) < \delta) = 0$.
 - 5: Set $\tilde{G} = \tilde{K}$, and let $\tilde{G}(\tilde{G} \neq 0) = 1$.
 - 6: Construct the graph \mathcal{G} from \tilde{G} , and determine its number of strongly connected components n_c using the MATLAB function ‘conncomp’.
 - 7: **if** $n_c == r$ **then**
 - 8: **break**
 - 9: **else if** $n_c > r$ **then**
 - 10: $\delta_{max} = \delta$
 - 11: **else**
 - 12: $\delta_{min} = \delta$
 - 13: **end if**
 - 14: Compute $\delta = (\delta_{min} + \delta_{max})/2$.
 - 15: **end for**
 - 16: Find the irreducible components K_1, \dots, K_r and the corresponding vertex sets π_1, \dots, π_r of the graph \mathcal{G} using the ‘conncomp’ function.
 - 17: **for** $i = 1, \dots, r$ **do**
 - 18: Set $C(\pi_i, i) = \frac{u_i}{\|u_i\|_2}$ and $C(\pi_i^c, i) = 0$, where u_i is the first column of K_i .
 - 19: **end for**
 - 20: Compute $B = AC$.
- Output:** Factor matrices $B \in \mathbb{R}_+^{m \times r}$, $C \in \mathbb{O}_+^{n \times r}$.
-

401 plicative Updates on Stiefel Manifolds), EM-ONMF [20] (EM-like alternating algo-
 402 rithm), ONPMF [20] (Orthogonal Nonnegatively Penalized Matrix Factorization),
 403 SN-ONMF [19] (Sparsity and Nuclear norm minimization for ONMF). To ensure a
 404 fair comparison, the parameters for competing methods are set following their origi-
 405 nal recommendations. Particularly, the maximum number of iterations of the pro-
 406 posed method is 1000, and the stopping criteria are as follows: for BiOR-NM3F,
 407 MU-ONMF, HALS, and EM-ONMF, $\|C^k - C^{k-1}\| < \epsilon$, where $\epsilon = 10^{-8}$; for ONPMF,
 408 $\|\min(C^k, 0)\|/\|C^k\| < 10^{-3}$; for SN-ONMF, $\max(\|K^k - K^{k-1}\|, \|X^k - K^k\|, \|Z^k -$
 409 $K^k\|) < \epsilon$, where K, X, Z are defined in model (13) of [19]; and for NS-ONMF,
 410 $\|c(\mathcal{W}^k)\| < \epsilon$ or $\rho_{k+1} > \rho_{max}$.

411 We generate a matrix $A \in \mathbb{R}_+^{300 \times 100}$ for three distinct cases. In Case 1, A is
 412 constructed as $A = BC^\top$, where $C \in \mathbb{O}_+^{100 \times 10}$ has a block diagonal structure, and
 413 $B \in \mathbb{R}_+^{300 \times 10}$ is a random matrix. In Case 2, A is generated similarly to Case 1, but
 414 the orthogonal factor C does not possess the block diagonal structure. In Case 3,
 415 $A = BC^\top + R$, where B and C are generated in the same manner as in Case 1, and
 416 the noise $R = 0.01 \times \text{rand}(300, 100)$.

417 To assess the performance of the comparison methods, we use the following eval-
 418 uation metrics:

$$419 \quad (4.1) \quad \text{err}(\tilde{C}) = \|\tilde{C}\tilde{C}^\top - CC^\top\|, \text{orth}(\tilde{C}) = \|\tilde{C}^\top \tilde{C} - I\|^2, \text{res}(\tilde{B}, \tilde{C}) = \frac{\|A_{gt} - \tilde{B}\tilde{C}^\top\|}{\|A_{gt}\|},$$

420 where \tilde{B}, \tilde{C} are the factors computed by the algorithms and $A_{gt} = BC^\top$ denotes the

421 ground-truth orthogonal nonnegative decomposable matrix. All methods are run 720
 422 times, and the average results are reported in Table 1.

TABLE 1
 Results of different methods for finding C

Case #	Method	CPU	err(\tilde{C})	orth(\tilde{C})	res(\tilde{B}, \tilde{C})
Case 1	BiOR-NM3F [8]	3.3850	9.6405×10^{-1}	9.6405×10^{-1}	7.2555×10^{-2}
	MU-ONMF [7]	0.6624	4.6774×10^{-4}	4.6074×10^{-4}	1.2026×10^{-5}
	EM-ONMF [20]	0.2079	4.2426×10^{-1}	4.5142×10^{-16}	4.3893×10^{-2}
	ONPMF [20]	0.9867	1.5283×10^{-3}	5.9523×10^{-15}	1.8239×10^{-4}
	SN-ONMF [19]	0.3765	2.2225×10^{-4}	3.6828×10^{-15}	2.6307×10^{-5}
	NS-ONMF	0.0144	2.0297×10^{-12}	3.5108×10^{-16}	4.6018×10^{-13}
Case 2	BiOR-NM3F [8]	3.8633	9.6695×10^{-1}	9.6695×10^{-1}	7.3526×10^{-2}
	MU-ONMF [7]	0.6414	4.6894×10^{-4}	4.6123×10^{-4}	1.5922×10^{-5}
	EM-ONMF [20]	0.0947	1.4142×10^{-1}	3.7551×10^{-16}	1.5074×10^{-2}
	ONPMF [20]	0.9535	1.5661×10^{-3}	4.6853×10^{-15}	1.9442×10^{-4}
	SN-ONMF [19]	0.3541	2.5372×10^{-5}	5.5830×10^{-15}	2.9005×10^{-6}
	NS-ONMF	0.0049	4.5946×10^{-12}	4.0030×10^{-16}	1.0189×10^{-12}
Case 3	BiOR-NM3F [8]	3.4089	9.8347×10^{-1}	9.8256×10^{-1}	8.0480×10^{-2}
	MU-ONMF [7]	0.6570	2.6336×10^{-2}	6.3920×10^{-3}	2.7990×10^{-2}
	EM-ONMF [20]	0.0126	2.3252×10^0	4.1919×10^{-16}	2.5274×10^{-1}
	ONPMF [20]	0.9625	3.0425×10^{-2}	6.3940×10^{-15}	2.7678×10^{-2}
	SN-ONMF [19]	0.6528	2.3917×10^{-2}	4.3275×10^{-15}	2.8026×10^{-2}
	NS-ONMF	2.3704	1.9531×10^{-2}	2.2204×10^{-16}	2.5906×10^{-2}

423 Table 1 reports the results of different methods for recovering the factor C across
 424 three synthetic cases. Both NS-ONMF and EM-ONMF achieve orthogonality up to
 425 machine precision (on the order of 10^{-16}), confirming the correctness of their formula-
 426 tions. Beyond orthogonality, NS-ONMF consistently attains the best reconstruction
 427 error and relative residual across all cases. In the noise-free settings (Cases 1 and 2),
 428 NS-ONMF is also the fastest method, converging within milliseconds while simultane-
 429 ously achieving superior accuracy. This efficiency arises because the output produced
 430 by Algorithm 3.2 already yields factors close to the optimal solution of model (1.4),
 431 allowing Algorithm 3.1 to converge within only a few iterations. Moreover, Algo-
 432 rithm 3.2 employs closed-form updates for both subproblems and typically converges
 433 within ten iterations, further reducing runtime in clean data regimes. In the noisy
 434 setting (Case 3), NS-ONMF requires more runtime than some competitors, yet it still
 435 produces the lowest reconstruction error and residual, demonstrating greater robust-
 436 ness to noise. By contrast, methods such as EM-ONMF converge faster in noisy cases
 437 but at the expense of substantially larger reconstruction errors and residuals. These
 438 results highlight a clear trade-off: while several algorithms enforce orthogonality ef-
 439 fectively, NS-ONMF achieves the most accurate and stable overall decomposition,
 440 offering both efficiency in noise-free regimes and robustness under noise.

441 In the following, we analyze the relationship between the noise level ε and the
 442 relative residual $\text{res}(\tilde{B}, \tilde{C})$. Similar to Case 3, where $A = BC^T + R$, we set $R = \varepsilon \times$
 443 $\text{rand}(300, 100)$, with $\varepsilon = \text{logspace}(-8, 0, 20)$, meaning ε is generated between 10^{-8} and
 444 1 on a logarithmic scale with 20 points. Figure 1(a) shows that the proposed method
 445 outperforms the competitors, especially for small noise levels. This also explains why,
 446 in Case 1 and Case 2 (noise-free cases), the method NS-ONMF requires less CPU
 447 time, as lower noise leads to higher accuracy and fewer iterations.

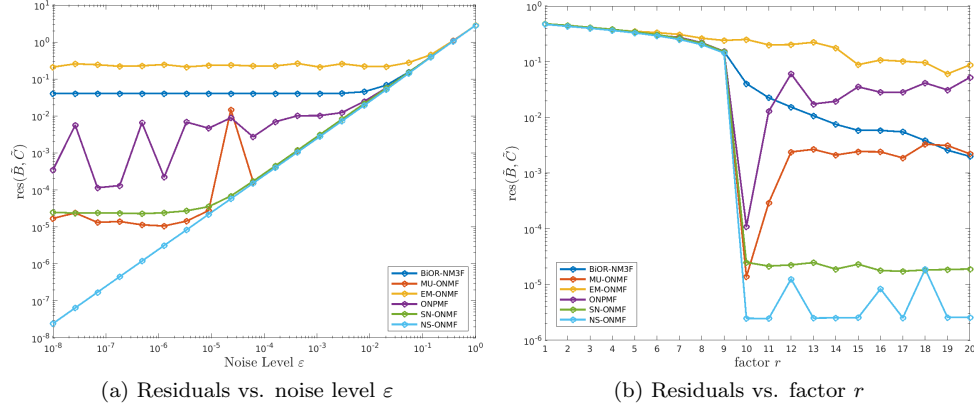


FIG. 1. Residual $\text{res}(\tilde{B}, \tilde{C})$ comparisons across noise levels (a) and factor values (b).

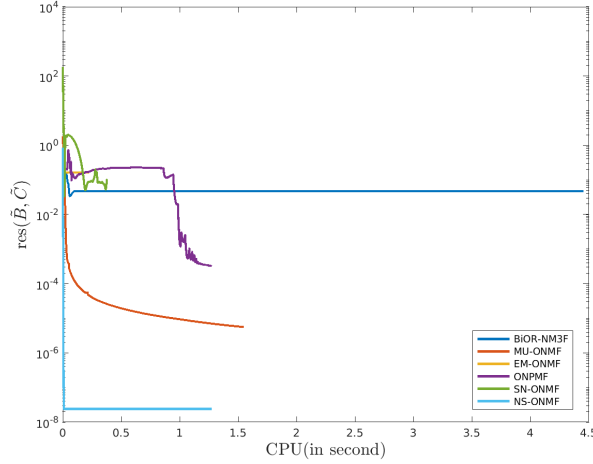


FIG. 2. Relative residual error vs. CPU time

448 Next, we examine the impact of the factor r on the residual $\text{res}(\tilde{B}, \tilde{C})$. Let
 449 $A \in \mathbb{R}_+^{300 \times 100}$ be generated as in Case 3, where the true rank of A is 10. However, we
 450 do not assume the rank of A is known and instead vary the factor r from 1 to 20 to
 451 evaluate how the methods handle different values of r . From Figure 1(b), we observe
 452 that for $r < \text{rank}(A)$, all methods perform similarly, with the residual near 1. When
 453 $r = \text{rank}(A) = 10$, most methods improve their performance, with the proposed NS-
 454 ONMF achieving the best outcome, reaching a residual of 10^{-5} . For $r > \text{rank}(A)$,
 455 the performance of the methods slightly deteriorates, but NS-ONMF still maintains
 456 a residual of 10^{-5} . This result further demonstrates that our model does not require
 457 the rank of A to be known, in accordance with Theorem 2.2.

458 In this part, we examine the efficiency and convergence behavior of the compared
 459 methods, we construct the synthetic matrix $A = BC^T + R$ using Case 3 with noise
 460 level 10^{-8} . Figure 2 plots the relative residual error $\text{res}(\tilde{B}, \tilde{C})$ against CPU time for
 461 all algorithms. This runtime-based comparison provides a more direct illustration of
 462 convergence efficiency. The results show that NS-ONMF rapidly decreases the residual

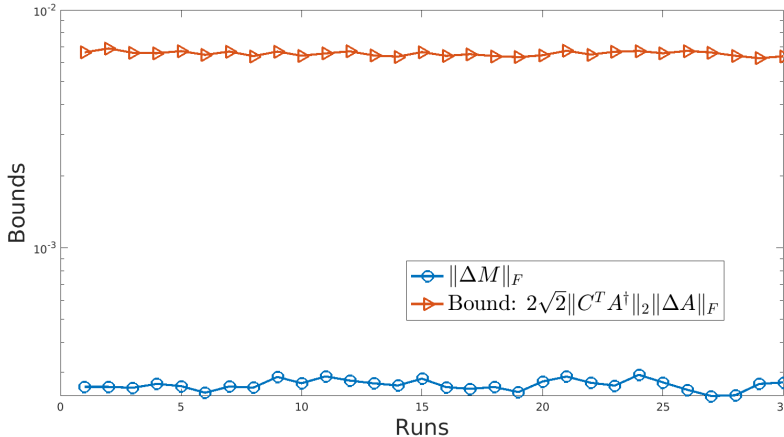


FIG. 3. Error bounds computed by Theorem 2.4

463 error, reaching nearly 10^{-8} within a fraction of a second. Although MU-ONMF also
 464 decreases steadily, it requires more runtime and does not achieve the same accuracy.
 465 ONPMF and SN-ONMF stagnate at higher error levels, while BiOR-NM3F maintains
 466 a relatively large error throughout. EM-ONMF is the fastest method though it attains
 467 a higher residual compared with NS-ONMF. These observations confirm that NS-
 468 ONMF is both effective and efficient: it converges quickly in terms of CPU time and
 469 consistently attains the smallest residual error among all tested methods.

470 To further validate Theorem 2.4, we construct a ground-truth orthogonal non-
 471 negative matrix decomposition $A = BC^\top$, where $C \in \mathbb{O}_+^{100 \times 10}$ and $B \in \mathbb{R}_+^{300 \times 10}$. We
 472 then perturb A by defining $\tilde{A} = A + R$, with $R = 10^{-4} \times \text{rand}(300, 100)$. Applying
 473 the proposed method to \tilde{A} yields the factors $\tilde{B}, \tilde{C}, \tilde{Z}$. We evaluate the deviation

474
$$\Delta M = ZZ^\top - \tilde{Z}\tilde{Z}^\top = CC^\top - \tilde{C}\tilde{C}^\top,$$

475 and compare the theoretical bound in (2.4) with the empirical value $\|\Delta M\|_F$. This
 476 experiment is repeated 30 times to ensure robustness. The results, shown in Figure 3,
 477 confirm that the empirical deviations are consistent with the theoretical bounds es-
 478 tablished in Theorem 2.4.

479 **4.1.2. Synthetic clustering data.** In this part, we generate the synthetic clus-
 480 tering data similar to [20]. Each dataset consists of six clusters, denoted as $\{\pi_i\}_{i=1}^6$,
 481 with each cluster containing $100 - (i - 1) \times 10$ data points, for a total of 450 data
 482 points. Each cluster centroid $u_i \in \mathbb{R}^{1000}$ for $1 \leq i \leq 6$ is generated uniformly at
 483 random within the unit cube $[0, 1]^{1000}$. Each data point m_j for $1 \leq j \leq 450$ is a scalar
 484 multiple of its corresponding cluster centroid, given by $m_j = \alpha u_k$, where α is chosen
 485 uniformly at random from the interval $[0.2, 1]$. The data points are then perturbed
 486 by adding noise drawn from a normal distribution with mean zero and standard devi-
 487 ation ε (negative entries are set to zero). Finally, A is constructed by collecting these
 488 data points into a matrix of size 1000×450 .

489 To obtain the clusters using ONMF methods, we utilize the nonnegative orthonor-
 490 mal factor C to determine the cluster assignments. Each row of C serves as a feature
 491 vector for a sample, and the column index of the positive entry in each row indicates
 492 the cluster to which the sample belongs. However, since some methods may yield
 493 multiple positive entries in a row, we assign the cluster based on the index of the

494 maximum entry. To evaluate the clustering results of the compared algorithms, we
 495 use accuracy as the performance metric. We examine the clustering performance of the
 496 proposed NS-ONMF method against six existing ONMF methods: BiOR-NM3F [8],
 497 MU-ONMF [7], HALS [21], EM-ONMF [20], ONPMF [20], and SN-ONMF [19].

498 The noise level ε is selected from the interval $[10^{-2}, 1]$ with a total of 50 points.
 499 Figure 4 presents the average accuracy of each method across different noise levels
 500 over 10 runs. The results indicate that the proposed NS-ONMF method and ONPMF
 501 consistently achieve the best performance. Both methods successfully identify the
 502 true clusters across a wide range of noise levels, demonstrating their robustness and
 503 effectiveness in clustering tasks.

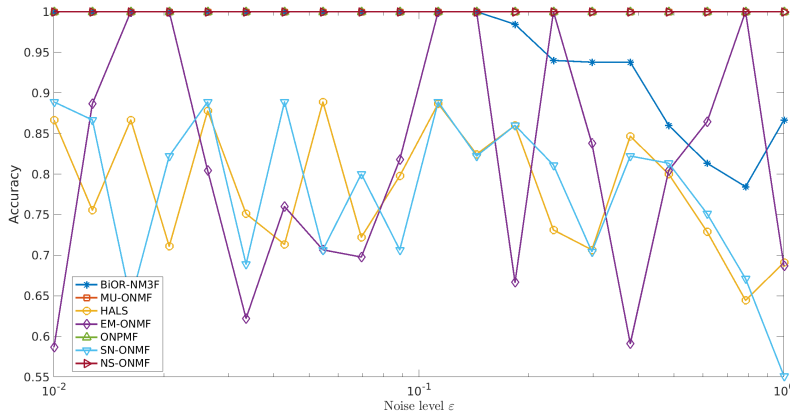


FIG. 4. Comparison of ONMF methods on synthetic clustering data with varying noise levels.

504 To examine the convergence behavior of Algorithm 3.1, we set the noise level
 505 $\varepsilon = 0.1$ and generate synthetic clustering data as previously described. We plot
 506 $c(\mathcal{W}^k)$ and $\mathcal{L}(\mathcal{W}^k, \Lambda^k, \rho_k)$ in Figure 5. The first plot, Figure 5(a), shows that $c(\mathcal{W}^k)$
 507 steadily decreases and converges to zero, consistent with the results presented in
 508 Theorem 3.2. This convergence indicates that the algorithm is progressively satisfying
 509 the constraints. The second plot, Figure 5(b), shows that $\mathcal{L}(\mathcal{W}^k, \Lambda^k, \rho_k)$ remains
 510 bounded throughout the process. We note that in this experiment the noise level
 511 is relatively large ($\varepsilon = 0.1$), which explains the higher iteration counts compared
 512 with the small-noise cases (e.g., $\varepsilon = 10^{-8}$). For small-noise settings, the Figure 2
 513 demonstrates that Algorithm 3.1 converges within only a few iterations and rapidly
 514 reduces the residual error to the order of 10^{-7} . By contrast, Figure 5 highlights that
 515 even under larger noise Algorithm 3.1 still converges in the sense that $c(\mathcal{W}^k) \rightarrow 0$,
 516 confirming robustness. We further emphasize that the augmented Lagrangian does not
 517 exhibit a monotone convergence trend because the penalty parameter ρ^k is adaptively
 518 increased (see Algorithm 3.1, line 8 and Eq. (3.10)), which may cause $\mathcal{L}(\mathcal{W}^k, \Lambda^k, \rho^k)$
 519 to increase even when the iterates approach a stationary point. Nevertheless, the
 520 stopping criterion $\|c(\mathcal{W}^k)\| < \epsilon$ is theoretically justified, since $\|c(\mathcal{W}^k)\| \rightarrow 0$ implies
 521 that \mathcal{W}^k converges to a stationary point of problem (3.2) by Theorem 3.2. This
 522 bounded behavior confirms the algorithm's stability under the given conditions, and
 523 the overall results support the effectiveness of Algorithm 3.1 in converging to a solution
 524 that satisfies the imposed constraints while maintaining stability in the objective
 525 function.

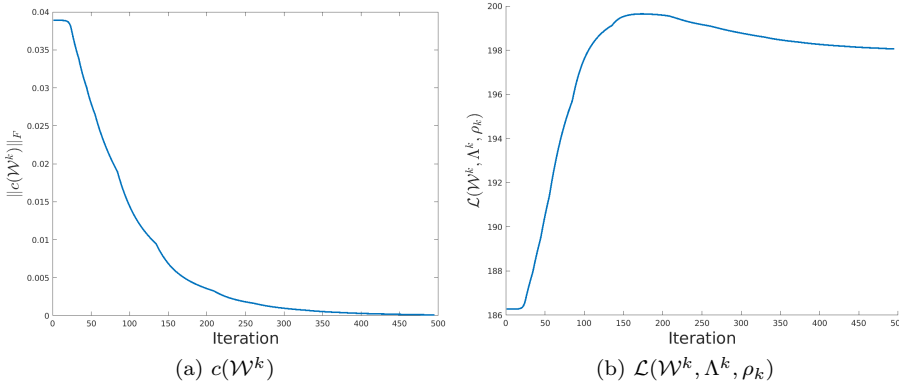


FIG. 5. Convergence analysis of Algorithm 3.1

526 **4.2. Real-world Data.** In this section, we will apply the proposed NS-ONMF
 527 method to real-world datasets, focusing on clustering. We compare the performance
 528 of NS-ONMF with six other ONMF methods: BiOR-NM3F [8], MU-ONMF [7], HALS
 529 [21], EM-ONMF [20], ONPMF [20], and SN-ONMF [19]. To evaluate the effectiveness
 530 of each method, we also use four metrics [16]: Purity, NMI, MIhat, and Accuracy.
 531 These metrics provide a comprehensive assessment of the clustering quality. To ensure
 532 comparisons focus on problem (1.1), clustering results are derived from the orthogonal
 533 nonnegative factor \tilde{C} , the computed solution of (1.1). Hard clustering is applied across
 534 all methods, assigning each sample to the cluster with the highest probability given
 535 by the factor matrix \tilde{C} .

536 We will evaluate the ONMF methods on real-world data clustering tasks, which
 537 include face clustering, document clustering, and object clustering. Each category
 538 comprises three well-known datasets, resulting in a total of nine datasets. The sum-
 539 mary of these datasets is provided in Table 2.

TABLE 2
 Summary of clustering datasets

Dataset	Type	Features (m)	Samples (n)	Classes (K)
<i>PIE_pose27</i>	Face	1024	2856	68
<i>Yale</i>	Face	1024	165	15
<i>ORL</i>	Face	1024	400	40
<i>bbcsport</i>	Document	4613	737	5
<i>tr11</i>	Document	6429	414	9
<i>tr12</i>	Document	5804	313	8
<i>seeds</i>	Object	7	210	3
<i>wine</i>	Object	13	173	3
<i>glass</i>	Object	9	214	6

540 Table 3 reports the performance of several ONMF methods on the *PIE_pose27*,
 541 *Yale*, and *ORL* face-clustering datasets. Across all three datasets, NS-ONMF con-
 542 sistentlly achieves the best clustering quality, attaining the highest Purity, NMI, MI-
 543 hat, and Accuracy. In particular, on the challenging *PIE_pose27* dataset, NS-ONMF
 544 reaches a Purity of 0.8207 and an Accuracy of 0.7990, outperforming all competing

545 methods by a clear margin. Similar improvements are observed on *Yale* and *ORL*,
 546 where NS-ONMF delivers better NMI and MIhat values, indicating more coherent
 547 and informative cluster structures. Figure 6, which displays the basis images for the
 548 *PIE_pose27* dataset, further supports these findings. The basis images generated by
 549 NS-ONMF are significantly clearer and more distinguishable than those produced by
 550 other methods, indicating NS-ONMF’s superior ability to capture essential clustering
 551 features. Although NS-ONMF incurs a higher computational cost, this is offset by its
 552 substantial improvements in clustering accuracy and the quality of the basis images.

TABLE 3
Results of different methods for face clustering

Dataset	Method	CPU	Purity	NMI	MIhat	Accuracy
<i>PIE_pose27</i>	BiOR-NM3F [8]	58.6085	0.3442	0.5843	0.5580	0.2773
	MU-ONMF [7]	15.9989	0.2868	0.4616	0.3919	0.2202
	HALS [21]	15.8625	0.1299	0.3239	0.2617	0.1278
	EM-ONMF [20]	2.6473	0.2721	0.5401	0.5322	0.2465
	ONPMF [20]	57.2348	0.1635	0.3228	0.2979	0.1488
	SN-ONMF [19]	156.2487	0.7122	0.8911	0.8672	0.6835
	NS-ONMF	422.0823	0.8207	0.8785	0.8704	0.7990
<i>Yale</i>	BiOR-NM3F [8]	0.6950	0.3818	0.4291	0.4120	0.3697
	MU-ONMF [7]	0.4034	0.3576	0.4154	0.3844	0.3091
	HALS [21]	0.7649	0.2606	0.3069	0.3110	0.2364
	EM-ONMF [20]	0.1460	0.3455	0.4375	0.4195	0.3394
	ONPMF [20]	2.7489	0.3091	0.3386	0.2962	0.3030
	SN-ONMF [19]	0.6933	0.3455	0.4028	0.3950	0.3333
	NS-ONMF	2.4502	0.4909	0.5305	0.5253	0.4848
<i>ORL</i>	BiOR-NM3F [8]	1.8026	0.4775	0.6495	0.6058	0.4225
	MU-ONMF [7]	1.0387	0.5350	0.7188	0.6773	0.4725
	HALS [21]	3.8143	0.1475	0.3815	0.3030	0.1425
	EM-ONMF [20]	0.4160	0.5275	0.7040	0.6770	0.4675
	ONPMF [20]	8.3403	0.1375	0.3167	0.2240	0.1325
	SN-ONMF [19]	2.9851	0.5200	0.6841	0.6801	0.4800
	NS-ONMF	8.2375	0.7000	0.8048	0.7878	0.6500

553 Table 4 compares the performance of ONMF methods on document clustering
 554 across the *bbsport*, *tr11*, and *tr12* datasets. In the *bbsport* dataset, NS-ONMF
 555 achieves the highest Purity (0.9457), NMI (0.8485), MIhat (0.8473), and Accuracy
 556 (0.9457), slightly outperforming ONPMF, which also attains competitive clustering
 557 quality. Figure 7 illustrating the *bbsport* clustering results reinforces these met-
 558 rics, showing that NS-ONMF produces well-separated clusters that closely reflect the
 559 ground truth. NS-ONMF consistently delivers superior separation and clarity, par-
 560 ticularly in more challenging datasets like *tr11* and *tr12*, making it the most robust
 561 method for document clustering among those evaluated.

562 Table 5 presents a comparison of various ONMF methods applied to object clus-
 563 tering tasks, specifically for the *seeds*, *wine*, and *glass* datasets. NS-ONMF emerges as
 564 the most effective method, delivering the highest values in Purity, NMI, MIhat, and
 565 Accuracy across all datasets. In the *seeds* dataset, for example, NS-ONMF achieves
 566 a standout Purity of 0.8810 and an NMI of 0.6518, demonstrating its robustness in
 567 differentiating between object classes. Figure 8 depicting the *seeds* dataset clustering
 568 further illustrates NS-ONMF’s capability to delineate distinct and coherent clusters,

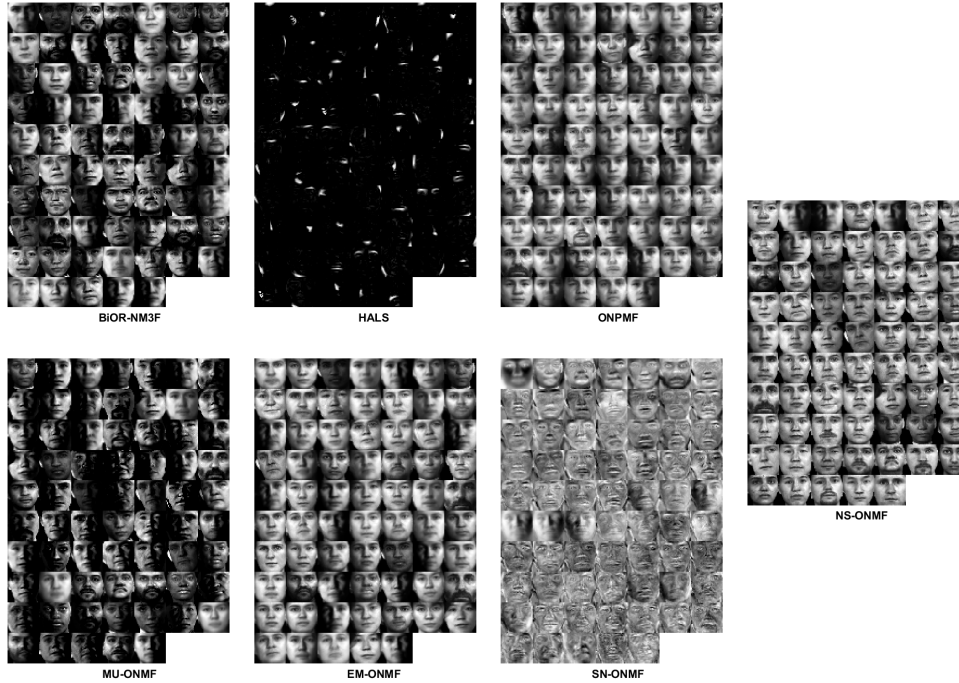


FIG. 6. Basis images generated by different ONMF methods for the *PIE_pose27* Dataset.

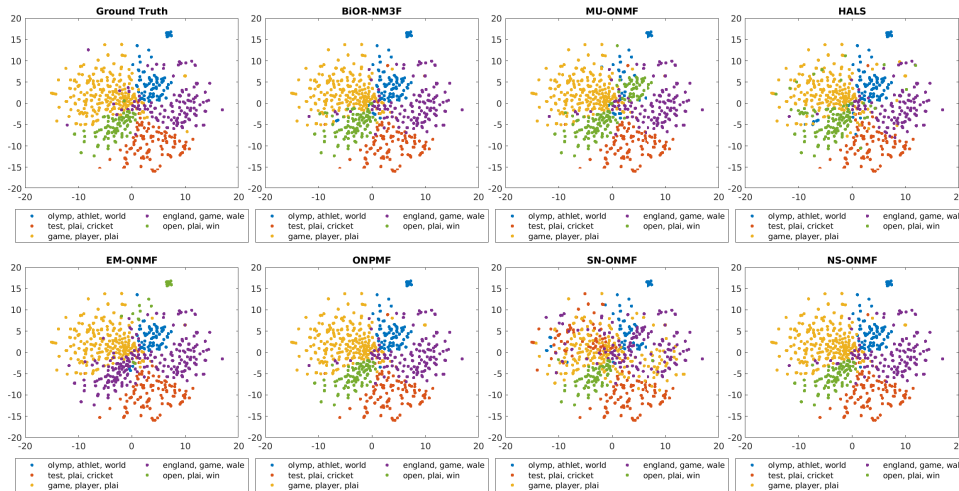


FIG. 7. Document clustering by ONMF methods for the *bbc sport* Dataset.

569 closely mirroring the ground truth distribution.

570 For these real-world datasets, it is unlikely that the data exhibit a near-ONMF
 571 structure; instead, they are more accurately modeled as $A = BC^T + R$, where R
 572 represents a substantial noise component. Under such conditions, NS-ONMF generally
 573 requires more iterations to converge, and the increased runtime should not be inter-
 574 preted as the sole reason for its improved clustering performance. As illustrated in
 575 Figure 2, several competing methods (e.g., BiOR-NM3F, MU-ONMF, and ONPMF)

TABLE 4
Results of different methods for document clustering

Dataset	Method	CPU	Purity	NMI	MIhat	Accuracy
<i>bbcspot</i>	BiOR-NM3F [8]	6.7560	0.9199	0.7901	0.7847	0.9199
	MU-ONMF [7]	6.1448	0.8412	0.7034	0.6966	0.8412
	HALS [21]	3.8879	0.8602	0.6834	0.6771	0.8602
	EM-ONMF [20]	0.9328	0.8100	0.7226	0.6896	0.7788
	ONPMF [20]	12.6639	0.9444	0.8424	0.8421	0.9444
	SN-ONMF [19]	10.3098	0.6418	0.4667	0.4653	0.6418
	NS-ONMF	29.1887	0.9457	0.8485	0.8473	0.9457
<i>tr11</i>	BiOR-NM3F [8]	7.9396	0.7029	0.5655	0.5273	0.5097
	MU-ONMF [7]	9.6800	0.6957	0.5267	0.4946	0.4638
	HALS [21]	7.6255	0.6184	0.4841	0.4780	0.5072
	EM-ONMF [20]	0.6722	0.7343	0.5848	0.5430	0.4783
	ONPMF [20]	23.6050	0.7367	0.5762	0.5410	0.4734
	SN-ONMF [19]	3.8672	0.4324	0.2928	0.2783	0.3478
	NS-ONMF	9.1732	0.7488	0.6173	0.5812	0.6256
<i>tr12</i>	BiOR-NM3F [8]	5.2975	0.6869	0.5247	0.5106	0.5495
	MU-ONMF [7]	5.7893	0.7125	0.5576	0.5463	0.6070
	HALS [21]	5.1425	0.4984	0.3510	0.3268	0.4153
	EM-ONMF [20]	1.3501	0.5623	0.4556	0.4246	0.4505
	ONPMF [20]	13.7598	0.6454	0.5304	0.5219	0.5495
	SN-ONMF [19]	2.2722	0.3834	0.2040	0.2029	0.3259
	NS-ONMF	4.8724	0.7157	0.6106	0.6006	0.6134

576 consume comparable or even greater CPU times, yet their residual errors stagnate
577 at much higher levels, resulting in weaker clustering outcomes. This confirms that
578 clustering quality is determined by the quality of the feasible solution rather than by
579 runtime alone. Moreover, subproblem (3.32) requires the computation of $n - r$ eigen-
580 pairs, which becomes computationally expensive for large n (as in document and face
581 datasets). This further contributes to the higher CPU times observed in Tables 3-5.
582 From another perspective, the computed matrix $Z \in \mathbb{R}^{n \times (n-r)}$ captures additional
583 informative structure by including a larger set of basis directions. Furthermore, Al-
584 gorithm 3.3 guarantees that C is both nonnegative and orthogonal, a property not
585 strictly enforced in some competing methods (e.g., SN-ONMF). These structural ad-
586 vantages enhance the quality of both Z and C , which in turn explain the consistently
587 superior clustering performance of NS-ONMF. The additional runtime observed in
588 noisy or large-scale datasets should therefore be regarded as a reasonable trade-off for
589 obtaining more accurate and reliable clustering results.

590 **5. Conclusions.** This paper advances the theoretical framework of ONMF by
591 characterizing properties of the null space of a nonnegative matrix $A \in \mathbb{R}^{m \times n}$
592 for the existence of ONMF. We prove that A has an ONMF if and only if there is a
593 matrix $Z \in \mathbb{R}^{n \times (n-r)}$ such that $AZ = 0$, $Z^\top Z = I$ and $ZZ^\top \leq I$. The new necessary
594 and sufficient condition provides solid mathematical foundation to study optimization
595 model (1.4) for ONMF. We propose an augmented Lagrangian algorithm for solving
596 (1.4) and prove the convergence of the algorithm. Numerical evaluations on both
597 synthetic and real-world datasets demonstrate the efficacy of the proposed methods,

TABLE 5
Results of different methods for object clustering

Dataset	Method	CPU	Purity	NMI	MIhat	Accuracy
<i>seeds</i>	BiOR-NM3F [8]	0.1996	0.7429	0.4272	0.4252	0.7429
	MU-ONMF [7]	0.0445	0.6190	0.3396	0.3112	0.6190
	HALS [21]	0.0463	0.6667	0.5543	0.4420	0.6667
	EM-ONMF [20]	0.0043	0.8000	0.5095	0.5082	0.8000
	ONPMF [20]	0.1725	0.4571	0.0437	0.0430	0.4571
	SN-ONMF [19]	0.7741	0.5571	0.1499	0.1495	0.5571
	NS-ONMF	0.7995	0.8810	0.6518	0.6508	0.8810
<i>wine</i>	BiOR-NM3F [8]	0.1719	0.6348	0.3791	0.3764	0.5843
	MU-ONMF [7]	0.0518	0.6348	0.3711	0.3488	0.6348
	HALS [21]	0.0581	0.4045	0.0081	0.0142	0.4045
	EM-ONMF [20]	0.0103	0.6966	0.3823	0.3762	0.6966
	ONPMF [20]	0.2531	0.4607	0.0352	0.0340	0.4045
	SN-ONMF [19]	0.6889	0.4775	0.0771	0.0771	0.4775
	NS-ONMF	0.8507	0.7022	0.3857	0.3803	0.7022
<i>glass</i>	BiOR-NM3F [8]	0.2234	0.4907	0.1895	0.1805	0.3692
	MU-ONMF [7]	0.0683	0.5374	0.2804	0.2736	0.4252
	HALS [21]	0.0859	0.3598	0.0111	0.0293	0.3131
	EM-ONMF [20]	0.0109	0.5467	0.3096	0.3055	0.4673
	ONPMF [20]	0.2499	0.3598	0.0308	0.0245	0.3084
	SN-ONMF [19]	0.7417	0.4579	0.2652	0.2538	0.4252
	NS-ONMF	2.2538	0.5701	0.3976	0.3615	0.5374

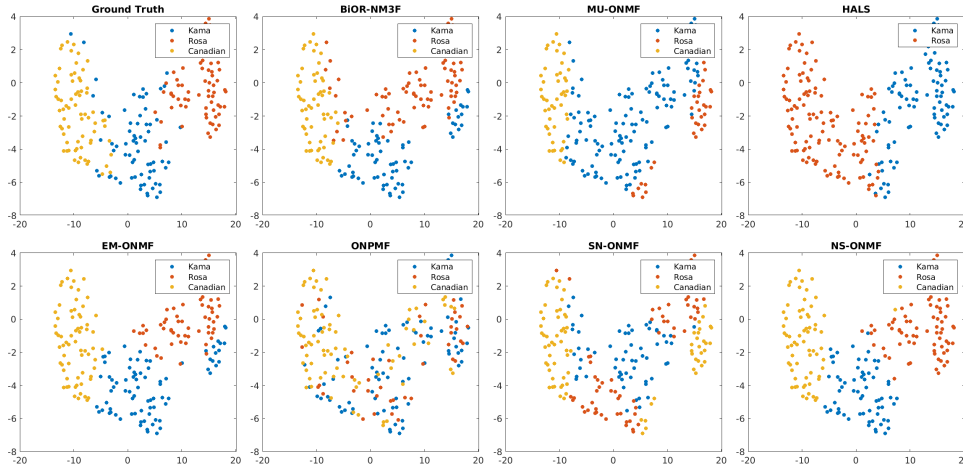


FIG. 8. Object clustering by ONMF methods for the seeds Dataset.

598 showcasing substantial improvements in clustering. These findings not only validate
 599 the robustness and accuracy of our approach but also underscore its utility in enhanc-
 600 ing interpretability in practical applications.

601 **Acknowledgement.** The authors are grateful to the two anonymous reviewers
 602 for their valuable comments and for verifying the implementation of the code, and to

603 Prof. Delin Chu and Prof. Zaikun Zhang for their insightful suggestions.

604

REFERENCES

- 605 [1] M. ASTERIS, D. PAPAIOPOULOS, AND A. G. DIMAKIS, *Orthogonal NMF through subspace*
606 *exploration*, Adv. Neural Inf. Process. Syst., 28 (2015).
- 607 [2] A. BERMAN AND R. J. PLEMMONS, *Nonnegative Matrices in the Mathematical Sciences*, SIAM,
608 1994.
- 609 [3] E. G. BIRGIN AND J. M. MARTÍNEZ, *Practical Augmented Lagrangian Methods for Constrained*
610 *Optimization*, SIAM, 2014.
- 611 [4] R. A. BRUALDI, H. J. RYSER, ET AL., *Combinatorial matrix theory*, vol. 39, Springer, 1991.
- 612 [5] X. CHEN, L. GUO, Z. LU, AND J. J. YE, *An augmented Lagrangian method for non-Lipschitz*
613 *nonconvex programming*, SIAM J. Numer. Anal., 55 (2017), pp. 168–193.
- 614 [6] X. CHEN, Y. HE, AND Z. ZHANG, *Tight error bounds for the sign-constrained Stiefel manifold*,
615 SIAM J. Optim., 35 (2025), pp. 302–329.
- 616 [7] S. CHOI, *Algorithms for orthogonal nonnegative matrix factorization*, in 2008 IEEE Interna-
617 tional Joint Conference on Neural Networks (IEEE World Congress on Computational
618 Intelligence), 2008, pp. 1828–1832, <https://doi.org/10.1109/IJCNN.2008.4634046>.
- 619 [8] C. DING, T. LI, W. PENG, AND H. PARK, *Orthogonal nonnegative matrix t-factorizations for*
620 *clustering*, in Proceedings of the 12th ACM SIGKDD international conference on Knowl-
621 edge discovery and data mining, 2006, pp. 126–135.
- 622 [9] L. K. FLEISCHER, B. HENDRICKSON, AND A. PINAR, *On identifying strongly connected com-*
623 *ponents in parallel*, in Parallel and Distributed Processing: 15 IPDPS 2000 Workshops
624 Cancun, Mexico, May 1–5, 2000 Proceedings 14, Springer, 2000, pp. 505–511.
- 625 [10] B. GAO, X. LIU, X. CHEN, AND Y.-X. YUAN, *A new first-order algorithmic framework for*
626 *optimization problems with orthogonality constraints*, SIAM J. Optim., 28 (2018), pp. 302–
627 332.
- 628 [11] R. A. HORN AND C. R. JOHNSON, *Matrix Analysis*, Cambridge University Press, 2012.
- 629 [12] B. JIANG, X. MENG, Z. WEN, AND X. CHEN, *An exact penalty approach for optimization with*
630 *nonnegative orthogonality constraints*, Math. Program., 198 (2023), pp. 855–897.
- 631 [13] K. KIMURA, Y. TANAKA, AND M. KUDO, *A fast hierarchical alternating least squares algorithm*
632 *for orthogonal nonnegative matrix factorization*, in Asian Conference on Machine Learning,
633 PMLR, 2015, pp. 129–141.
- 634 [14] B. LI, G. ZHOU, AND A. CICHOCKI, *Two efficient algorithms for approximately orthogonal*
635 *nonnegative matrix factorization*, IEEE Signal Process. Lett., 22 (2014), pp. 843–846.
- 636 [15] W. LI, J. LI, X. LIU, AND L. DONG, *Two fast vector-wise update algorithms for orthogonal*
637 *nonnegative matrix factorization with sparsity constraint*, J. Comput. Appl. Math., 375
638 (2020), p. 112785.
- 639 [16] C. D. MANNING, P. RAGHAVAN, AND H. SCHÜTZE, *Introduction to Information Retrieval*,
640 vol. 39, Cambridge University Press, 2008.
- 641 [17] A. MIRZAL, *A convergent algorithm for orthogonal nonnegative matrix factorization*, J. Com-
642 put. Appl. Math., 260 (2014), pp. 149–166.
- 643 [18] C. C. PAIGE AND M. A. SAUNDERS, *Towards a generalized singular value decomposition*, SIAM
644 J. Numer. Anal., 18 (1981), pp. 398–405.
- 645 [19] J. PAN AND M. NG, *Orthogonal nonnegative matrix factorization by sparsity and nuclear norm*
646 *optimization*, SIAM J. Matrix Anal. Appl., 39 (2018), pp. 856–875.
- 647 [20] F. POMPILI, N. GILLIS, P.-A. ABSIL, AND F. GLINEUR, *Two algorithms for orthogonal non-*
648 *negative matrix factorization with application to clustering*, Neurocomputing, 141 (2014),
649 pp. 15–25.
- 650 [21] M. SHIGA, K. TATSUMI, S. MUTO, K. TSUDA, Y. YAMAMOTO, T. MORI, AND T. TANJI, *Sparse*
651 *modeling of EELS and EDX spectral imaging data by nonnegative matrix factorization*,
652 Ultramicroscopy, 170 (2016), pp. 43–59.
- 653 [22] K. SIKORSKI, *Bisection is optimal*, Numer. Math., 40 (1982), pp. 111–117.
- 654 [23] N. XIAO, X. LIU, AND Y. YUAN, *Exact penalty function for $2,1$ norm minimization over the*
655 *stiefel manifold*, SIAM J. Optim., 31 (2021), pp. 3097–3126.