

# Ellipsoidal Rounding for Nonnegative Matrix Factorization Under Noisy Separability

**Tomohiko Mizutani**

MIZUTANI@KANAGAWA-U.AC.JP

*Department of Information Systems Creation  
Kanagawa University  
Yokohama, 221-8686, Japan*

**Editor:** Nathan Srebro

## Abstract

We present a numerical algorithm for nonnegative matrix factorization (NMF) problems under noisy separability. An NMF problem under separability can be stated as one of finding all vertices of the convex hull of data points. The research interest of this paper is to find the vectors as close to the vertices as possible in a situation in which noise is added to the data points. Our algorithm is designed to capture the shape of the convex hull of data points by using its enclosing ellipsoid. We show that the algorithm has correctness and robustness properties from theoretical and practical perspectives; correctness here means that if the data points do not contain any noise, the algorithm can find the vertices of their convex hull; robustness means that if the data points contain noise, the algorithm can find the near-vertices. Finally, we apply the algorithm to document clustering, and report the experimental results.

**Keywords:** nonnegative matrix factorization, separability, robustness to noise, enclosing ellipsoid, document clustering

## 1. Introduction

This paper presents a numerical algorithm for nonnegative matrix factorization (NMF) problems under noisy separability. The problem can be regarded as a special case of an NMF problem. Let  $\mathbb{R}_+^{d \times m}$  be the set of  $d$ -by- $m$  nonnegative matrices, and  $\mathbb{N}$  be the set of nonnegative integer numbers. A nonnegative matrix is a real matrix whose elements are all nonnegative. For a given  $\mathbf{A} \in \mathbb{R}_+^{d \times m}$  and  $r \in \mathbb{N}$ , the nonnegative matrix factorization (NMF) problem is to find  $\mathbf{F} \in \mathbb{R}_+^{d \times r}$  and  $\mathbf{W} \in \mathbb{R}_+^{r \times m}$  such that the product  $\mathbf{FW}$  is as close to  $\mathbf{A}$  as possible. The nonnegative matrices  $\mathbf{F}$  and  $\mathbf{W}$  give a factorization of  $\mathbf{A}$  of the form,

$$\mathbf{A} = \mathbf{FW} + \mathbf{N},$$

where  $\mathbf{N}$  is a  $d$ -by- $m$  matrix. This factorization is referred to as the NMF of  $\mathbf{A}$ .

Recent studies have shown that NMFs are useful for tackling various problems such as facial image analysis (Lee and Seung, 1999), topic modeling (Arora et al., 2012b, 2013; Ding et al., 2013), document clustering (Xu et al., 2003; Shahnaz et al., 2006), hyperspectral unmixing (Nascimento and Dias, 2005; Miao and Qi, 2007; Gillis and Vavasis, 2014), and blind source separation (Cichocki et al., 2009). Many algorithms have been developed in the context of solving such practical applications. However, there are some drawbacks in

the use of NMFs for such applications. One of them is in the hardness of solving an NMF problem. In fact, the problem has been shown to be NP-hard by Vavasis (2009).

As a remedy for the hardness of the problem, Arora et al. (2012a) proposed to exploit the notion of separability, which was originally introduced by Donoho and Stodden (2003) for the uniqueness of NMF. An NMF problem under separability becomes a tractable one. *Separability* assumes that  $\mathbf{A} \in \mathbb{R}_+^{d \times m}$  can be represented as

$$\mathbf{A} = \mathbf{F}\mathbf{W} \text{ for } \mathbf{F} \in \mathbb{R}_+^{d \times r} \text{ and } \mathbf{W} = (\mathbf{I}, \mathbf{K})\mathbf{\Pi} \in \mathbb{R}_+^{r \times m}, \quad (1)$$

where  $\mathbf{I}$  is an  $r$ -by- $r$  identity matrix,  $\mathbf{K}$  is an  $r$ -by- $(m - r)$  nonnegative matrix, and  $\mathbf{\Pi}$  is an  $m$ -by- $m$  permutation matrix. This means that each column of  $\mathbf{F}$  corresponds to that of  $\mathbf{A}$  up to a scaling factor. A matrix  $\mathbf{A}$  is said to be a *separable matrix* if it can be represented in the form (1). In this paper, we call  $\mathbf{F}$  the *basis matrix* of a separable matrix, and  $\mathbf{W}$ , as well as its submatrix  $\mathbf{K}$ , the *weight matrix*. *Noisy separability* assumes that a separable matrix  $\mathbf{A}$  contains a noise matrix  $\mathbf{N}$  such that  $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{N}$ , where  $\mathbf{N}$  is a  $d$ -by- $m$  matrix. Arora et al. showed that there exists an algorithm for finding the near-basis matrix of a noisy separable one if the noise is small in magnitude. Although a separability assumption restricts the fields of application for NMFs, it is known to be reasonable at least, in the contexts of document clustering (Kumar et al., 2013), topic modeling (Arora et al., 2012a,b, 2013), and hyperspectral unmixing (Gillis and Vavasis, 2014). In particular, this assumption is widely used as a pure-pixel assumption in hyperspectral unmixing (see, for instance, Nascimento and Dias, 2005; Miao and Qi, 2007; Gillis and Vavasis, 2014).

An NMF problem under noisy separability is to seek for the basis matrix of a noisy separable one. The problem is formally described as follows:

**Problem 1** *Let a data matrix  $\mathbf{M}$  be a noisy separable matrix of size  $d$ -by- $m$ . Find an index set  $\mathcal{I}$  with cardinality  $r$  on  $\{1, \dots, m\}$  such that  $\mathbf{M}(\mathcal{I})$  is as close to the basis matrix  $\mathbf{F}$  as possible.*

Here,  $\mathbf{M}(\mathcal{I})$  denotes a submatrix of  $\mathbf{M}$  that consists of every column vector with an index in  $\mathcal{I}$ . We call the column vector of  $\mathbf{M}$  a *data point* and that of the basis matrix  $\mathbf{F}$  a *basis vector*. An ideal algorithm for the problem should have correctness and robustness properties; correctness here means that, if the data matrix  $\mathbf{M}$  is just a separable one, the algorithm can find the basis matrix; robustness means that, if the data matrix  $\mathbf{M}$  is a noisy separable one, the algorithm can find the near-basis matrix. A formal description of the properties is given in Section 2.1

We present a novel algorithm for Problem 1. The main contribution of this paper is to show that the algorithm has correctness and robustness properties from theoretical and practical perspectives. It is designed on the basis of the geometry of a separable matrix. Under reasonable assumptions, the convex hull of the column vectors of a separable matrix forms a simplex, and in particular, the basis vectors correspond to the vertices. Therefore, if all vertices of a simplex can be found, we can obtain the basis matrix of the separable matrix. Our algorithm uses the fact that the vertices of simplex can be found by an ellipsoid. That is, if we draw the minimum-volume enclosing ellipsoid (MVEE) for a simplex, the ellipsoid only touches its vertices. More precisely, we give plus and minus signs to the vertices of a simplex, and take the convex hull; it becomes a crosspolytope having the simplex as one of

the facets. Then, the MVEE for the crosspolytope only touches the vertices of the simplex with plus and minus signs.

Consider Problem 1 without noise. In this case, the data matrix is just a separable one. Our algorithm computes the MVEE for the data points and outputs the points on the boundary of the ellipsoid. Then, the obtained points correspond to the basis vectors for a separable matrix. We show in Theorem 5 that the correctness property holds. Moreover, the algorithm works well even when the problem contains noise. We show in Theorem 9 that, if the noise is lower than a certain level, the algorithm correctly identifies the near-basis vectors for a noisy separable matrix, and hence, the robustness property holds. The existing algorithms (Arora et al., 2012a; Bittorf et al., 2012; Gillis, 2013; Gillis and Luce, 2013; Gillis and Vavasis, 2014; Kumar et al., 2013) are formally shown to have these correctness and robustness properties. In Section 2.4, our correctness and robustness properties are compared with those of the existing algorithms.

It is possible that noise will exceed the level that Theorem 9 guarantees. In such a situation, the MVEE for the data points may touch many points. Hence,  $r$  points need to be selected from the points on the boundary of the ellipsoid. We make the selection by using existing algorithms such as SPA (Gillis and Vavasis, 2014) and XRAY (Kumar et al., 2013). Our algorithm thus works as a preprocessor which filters out basis vector candidates from the data points and enhance the performance of existing algorithms.

We demonstrated the robustness of the algorithms to noise through experiments with synthetic data sets. In particular, we experimentally compared our algorithm with SPA and XRAY. We synthetically generated data sets with various noise levels, and measured the robustness of an algorithm by its recovery rate. The experimental results indicated that our algorithm can improve the recovery rates of SPA and XRAY.

Finally, we applied our algorithm to document clustering. Separability for a document-term matrix means that each topic has an anchor word. An anchor word is a word which is contained in one topic but not contained in the other topics. If an anchor word is found, it suggests the existence of its associated topic. We conducted experiments with document corpora and compared the clustering performances of our algorithm and SPA. The experimental results indicated that our algorithm would usually outperform SPA and can extract more recognizable topics.

The rest of this paper is organized as follows. Section 2 gives an outline of our algorithm and reviews related work. Then, the correctness and robustness properties of our algorithm are given, and a comparison with existing algorithms is described. Section 3 reviews the formulation and algorithm of computing the MVEE for a set of points. Sections 4 and 5 are the main part of this paper. We show the correctness and robustness properties of our algorithm in Section 4, and discuss its practical implementation in Section 5. Section 6 reports the numerical experiments for the robustness of algorithms and document clustering. Section 7 gives concluding remarks.

## 1.1 Notation and Symbols

We use  $\mathbb{R}^{d \times m}$  to denote a set of real matrices of size  $d$ -by- $m$ , and  $\mathbb{R}_+^{d \times m}$  to denote a set of nonnegative matrices of  $d$ -by- $m$ . Let  $\mathbf{A} \in \mathbb{R}^{d \times m}$ . The symbols  $\mathbf{A}^\top$  and  $\text{rank}(\mathbf{A})$  respectively denote the transposition and the rank. The symbols  $\|\mathbf{A}\|_p$  and  $\|\mathbf{A}\|_F$  are the matrix  $p$ -norm

and the Frobenius norm. The symbol  $\sigma_i(\mathbf{A})$  is the  $i$ th largest singular value. Let  $\mathbf{a}_i$  be the  $i$ th column vector of  $\mathbf{A}$ , and  $\mathcal{I}$  be a subset of  $\{1, \dots, m\}$ . The symbol  $\mathbf{A}(\mathcal{I})$  denotes a  $d$ -by- $|\mathcal{I}|$  submatrix of  $\mathbf{A}$  such that  $(\mathbf{a}_i : i \in \mathcal{I})$ . The convex hull of all the column vectors of  $\mathbf{A}$  is denoted by  $\text{conv}(\mathbf{A})$ , and referred to as the convex hull of  $\mathbf{A}$  for short. We denote an identity matrix and a vector of all ones by  $\mathbf{I}$  and  $\mathbf{e}$ , respectively.

We use  $\mathbb{S}^d$  to denote a set of real symmetric matrices of size  $d$ . Let  $\mathbf{A} \in \mathbb{S}^d$ . If the matrix is positive definite, we represent it as  $\mathbf{A} \succ \mathbf{0}$ . Let  $\mathbf{A}_1 \in \mathbb{S}^d$  and  $\mathbf{A}_2 \in \mathbb{S}^d$ . We denote by  $\langle \mathbf{A}_1, \mathbf{A}_2 \rangle$  the Frobenius inner product of the two matrices which is given as the trace of matrix  $\mathbf{A}_1 \mathbf{A}_2$ .

We use a MATLAB-like notation. Let  $\mathbf{A}_1 \in \mathbb{R}^{d \times m_1}$  and  $\mathbf{A}_2 \in \mathbb{R}^{d \times m_2}$ . We denote by  $(\mathbf{A}_1, \mathbf{A}_2)$  the horizontal concatenation of the two matrices, which is a  $d$ -by- $(m_1 + m_2)$  matrix. Let  $\mathbf{A}_1 \in \mathbb{R}^{d_1 \times m}$  and  $\mathbf{A}_2 \in \mathbb{R}^{d_2 \times m}$ . We denote by  $(\mathbf{A}_1; \mathbf{A}_2)$  the vertical concatenation of the two matrices, and it is a matrix of the form,

$$\begin{pmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{pmatrix} \in \mathbb{R}^{(d_1 + d_2) \times m}.$$

Let  $\mathbf{A}$  be a  $d$ -by- $m$  rectangular diagonal matrix having diagonal elements  $a_1, \dots, a_t$  where  $t = \min\{d, m\}$ . We use  $\text{diag}(a_1, \dots, a_t)$  to denote the matrix.

## 2. Outline of Proposed Algorithm and Comparison with Existing Algorithms

Here, we formally describe the properties mentioned in Section 1 that an algorithm is expected to have, and also describe the assumptions we place on Problem 1. Next, we give a geometric interpretation of a separable matrix under these assumptions, and then, outline the proposed algorithm. After reviewing the related work, we describe the correctness and robustness properties of our algorithm and compare with those of the existing algorithms.

### 2.1 Preliminaries

Consider Problem 1 whose data matrix  $\mathbf{M}$  is a noisy separable one of the form  $\mathbf{A} + \mathbf{N}$ . Here,  $\mathbf{A}$  is a separable matrix of (1) and  $\mathbf{N}$  is a noise matrix. We can rewrite it as

$$\begin{aligned} \mathbf{M} &= \mathbf{A} + \mathbf{N} \\ &= \mathbf{F}(\mathbf{I}, \mathbf{K})\mathbf{\Pi} + \mathbf{N} \\ &= (\mathbf{F} + \mathbf{N}^{(1)}, \mathbf{F}\mathbf{K} + \mathbf{N}^{(2)})\mathbf{\Pi} \end{aligned} \tag{2}$$

where  $\mathbf{N}^{(1)}$  and  $\mathbf{N}^{(2)}$  are  $d$ -by- $r$  and  $d$ -by- $\ell$  submatrices of  $\mathbf{N}$  such that  $\mathbf{N}\mathbf{\Pi}^{-1} = (\mathbf{N}^{(1)}, \mathbf{N}^{(2)})$ . Hereinafter, we use the notation  $\ell$  to denote  $m - r$ . The goal of Problem 1 is to identify an index set  $\mathcal{I}$  such that  $\mathbf{M}(\mathcal{I}) = \mathbf{F} + \mathbf{N}^{(1)}$ .

As mentioned in Section 1, it is ideal that an algorithm for Problem 1 has correctness and robustness properties. These properties are formally described as follows:

- **Correctness.** If the data matrix  $\mathbf{M}$  does not contain a noise matrix  $\mathbf{N}$  and is just a separable matrix, the algorithm returns an index set  $\mathcal{I}$  such that  $\mathbf{M}(\mathcal{I}) = \mathbf{F}$ .

- **Robustness.** If the data matrix  $\mathbf{M}$  contains a noise matrix  $\mathbf{N}$  and is a noisy separable matrix such that  $\|\mathbf{N}\|_p < \epsilon$ , the algorithm returns an index set  $\mathcal{I}$  such that  $\|\mathbf{M}(\mathcal{I}) - \mathbf{F}\|_p < \tau\epsilon$  for some constant real number  $\tau$ .

In particular, the robustness property has  $\tau = 1$ , if an algorithm can identify an index set  $\mathcal{I}$  such that  $\mathbf{M}(\mathcal{I}) = \mathbf{F} + \mathbf{N}^{(1)}$  where  $\mathbf{F}$  and  $\mathbf{N}^{(1)}$  are of (2) since  $\|\mathbf{M}(\mathcal{I}) - \mathbf{F}\|_p = \|\mathbf{N}^{(1)}\|_p < \epsilon$ .

In the design of the algorithm, some assumptions are usually placed on a separable matrix. Our algorithm uses Assumption 1.

**Assumption 1** *A separable matrix  $\mathbf{A}$  of (1) consists of an basis matrix  $\mathbf{F}$  and a weight matrix  $\mathbf{W}$  satisfying the following conditions.*

1-a) *Every column vector of weight matrix  $\mathbf{W}$  has unit 1-norm.*

1-b) *The basis matrix  $\mathbf{F}$  has full column rank.*

Assumption 1-a can be invoked without loss of generality. If the  $i$ th column of  $\mathbf{W}$  is zero, so is the  $i$ th column of  $\mathbf{A}$ . Therefore, we can construct a smaller separable matrix having  $\mathbf{W}$  with no zero column. Also, since we have

$$\mathbf{A} = \mathbf{F}\mathbf{W} \Leftrightarrow \mathbf{A}\mathbf{D} = \mathbf{F}\mathbf{W}\mathbf{D},$$

every column of  $\mathbf{W}$  can have unit 1-norm. Here,  $\mathbf{D}$  denotes a diagonal matrix having the  $(i, i)$ th diagonal element  $d_{ii} = 1/\|\mathbf{w}_i\|_1$ .

The same assumption is used by the algorithm in Gillis and Vavasis (2014). We may get the feeling that 1-b is strong. The algorithms (Arora et al., 2012a; Bittorf et al., 2012; Gillis, 2013; Gillis and Luce, 2013; Kumar et al., 2013) instead assume *simpliciality*, wherein no column vector of  $\mathbf{F}$  can be represented as a convex hull of the remaining vectors of  $\mathbf{F}$ . Although 1-b is a stronger assumption, it still seems reasonable for Problem 1 from the standpoint of practical application. This is because, in such cases, it is less common for the column vectors of the basis matrix  $\mathbf{F}$  to be linearly dependent.

## 2.2 Outline of Proposed Algorithm

Let us take a look at Problem 1 from a geometric point of view. For simplicity, consider the noiseless case first. Here, a data matrix is just a separable matrix  $\mathbf{A}$ . Separability implies that  $\mathbf{A}$  has a factorization of the form (1). Under Assumption 1,  $\text{conv}(\mathbf{A})$  becomes an  $(r - 1)$ -dimensional simplex in  $\mathbb{R}^d$ . The left part of Figure 1 visualizes a separable data matrix. The white points are data points, and the black ones are basis vectors. The key observation is that the basis vectors  $\mathbf{f}_1, \dots, \mathbf{f}_r$  of  $\mathbf{A}$  correspond to the vertices of  $\text{conv}(\mathbf{A})$ . This is due to separability. Therefore, if all vertices of  $\text{conv}(\mathbf{A})$  can be found, we can obtain the basis matrix  $\mathbf{F}$  of  $\mathbf{A}$ . This is not hard task, and we can design an efficient algorithm for doing it. But, if noise is added to a separable matrix, the task becomes hard. Let us suppose that the data matrix of Problem 1 is a noisy separable matrix  $\tilde{\mathbf{A}}$  of the form  $\mathbf{A} + \mathbf{N}$ . The vertices of  $\text{conv}(\tilde{\mathbf{A}})$  do not necessarily match the basis vectors  $\mathbf{f}_1, \dots, \mathbf{f}_r$  of  $\mathbf{A}$ . The right part of Figure 1 visualizes a noisy separable data matrix. This is the main reason why it is hard to identify the basis matrix from noisy separable one.

Our algorithm is designed on the basis of Proposition 3 in Section 4.1; it states that all vertices of a simplex can be found by using an ellipsoid. We here describe the proposition

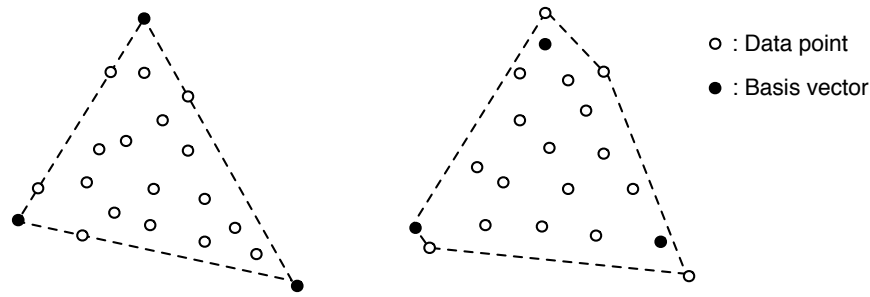


Figure 1: Convex hull of a separable data matrix with  $r = 3$  under Assumption 1. (Left) Noiseless case. (Right) Noisy case.

from a geometric point of view. Consider an  $(r - 1)$ -dimensional simplex  $\Delta$  in  $\mathbb{R}^r$ . Let  $\mathbf{g}_1, \dots, \mathbf{g}_r \in \mathbb{R}^r$  be the vertices of  $\Delta$ , and  $\mathbf{b}_1, \dots, \mathbf{b}_\ell \in \mathbb{R}^r$  be the points in  $\Delta$ . We draw the MVEE centered at the origin for a set  $\mathcal{S} = \{\pm\mathbf{g}_1, \dots, \pm\mathbf{g}_r, \pm\mathbf{b}_1, \dots, \pm\mathbf{b}_\ell\}$ . Then, the proposition says that the ellipsoid only touches the points  $\pm\mathbf{g}_1, \dots, \pm\mathbf{g}_r$  among all the points in  $\mathcal{S}$ . Therefore, the vertices of  $\Delta$  can be found by checking whether the points in  $\mathcal{S}$  lie on the boundary of ellipsoid. We should mention that the convex hull of the points in  $\mathcal{S}$  becomes a full-dimensional crosspolytope in  $\mathbb{R}^r$ . Figure 2 illustrates the MVEE for a crosspolytope in  $\mathbb{R}^3$ .

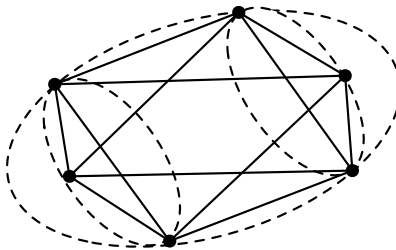


Figure 2: Minimum-volume enclosing ellipsoid for a full-dimensional crosspolytope in  $\mathbb{R}^3$ .

Under Assumption 1, the convex hull of a separable matrix  $\mathbf{A}$  becomes an  $(r - 1)$ -dimensional simplex in  $\mathbb{R}^d$ . Therefore, we rotate and embed the simplex in  $\mathbb{R}^r$  by using an orthogonal transformation. Such a transformation can be obtained by singular value decomposition (SVD) of  $\mathbf{A}$ .

Now let us outline our algorithm for Problem 1. In this description, we assume for simplicity that the data matrix is a separable one  $\mathbf{A} \in \mathbb{R}_+^{d \times m}$ . First, the algorithm constructs an orthogonal transformation through the SVD of  $\mathbf{A}$ . By applying the transformation, it transforms  $\mathbf{A}$  into a matrix  $\mathbf{P} \in \mathbb{R}^{r \times m}$  such that the  $\text{conv}(\mathbf{P})$  is an  $(r - 1)$ -dimensional sim-

plex in  $\mathbb{R}^r$ . Next, it draws the MVEE centered at the origin for a set  $\mathcal{S} = \{\pm \mathbf{p}_1, \dots, \pm \mathbf{p}_m\}$ , where  $\mathbf{p}_1, \dots, \mathbf{p}_m$  are the column vectors of  $\mathbf{P}$ , and outputs  $r$  points lying on the ellipsoid.

We call the algorithm *ellipsoidal rounding*, abbreviated as ER. The main computational costs of ER are in computing the SVD of  $\mathbf{A}$  and the MVEE for  $\mathcal{S}$ . The MVEE computation can be formulated as a tractable convex optimization problem with  $m$  variables. A polynomial-time algorithm exists, and it is also known that a hybrid of the interior-point algorithm and cutting plane algorithm works efficiently in practice.

In later sections, we will see that ER algorithm works well even if noise is added. In particular, we show that ER correctly identifies the near-basis vectors of a noisy separable matrix if the noise is smaller than some level. We consider a situation in which the noise exceeds that level. In such a situation, the shape of crosspolytope formed by the data points is considerably perturbed by the noise, and it is possible that the MVEE touches many points. We thus need to select  $r$  points from the points on the boundary of the ellipsoid. In this paper, we perform existing algorithms such as SPA (Gillis and Vavasis, 2014) and XRAY (Kumar et al., 2013) to make the selection. Hence, ER works as a preprocessor which filters out basis vector candidates from data points and enhances the performance of existing algorithms.

### 2.3 Related Work

First, we will review the algorithms for NMF of general nonnegative matrices. There are an enormous number of studies. A commonly used approach is to formulate it as a nonconvex optimization problem and compute the local solution. Let  $\mathbf{A}$  be a  $d$ -by- $m$  nonnegative matrix, and consider an optimization problem with matrix variables  $\mathbf{F} \in \mathbb{R}^{d \times r}$  and  $\mathbf{W} \in \mathbb{R}^{r \times m}$ ,

$$\text{minimize } \|\mathbf{F}\mathbf{W} - \mathbf{A}\|_{\mathbf{F}}^2 \text{ subject to } \mathbf{F} \geq \mathbf{0} \text{ and } \mathbf{W} \geq \mathbf{0}.$$

This is an intractable nonconvex optimization problem, and in fact, it was shown to be NP-hard by Vavasis (2009). Therefore, the research target is in how to compute the local solution efficiently. It is popular to use the block coordinate descent (BCD) algorithm for this purpose. The algorithm solves the problem by alternately fixing the variables  $\mathbf{F}$  and  $\mathbf{W}$ . The problem obtained by fixing either of  $\mathbf{F}$  and  $\mathbf{W}$  becomes a convex optimization problem. The existing studies propose to use, for instance, the projected gradient algorithm (Lin, 2007) and its variant (Lee and Seung, 2001), active set algorithm (Kim and Park, 2008, 2011), and projected quasi-Newton algorithm (Gong and Zhang, 2012). It is reported that the BCD algorithm shows good performance on average in computing NMFs. However, its performance depends on how we choose the initial point for starting the algorithm. We refer the reader to Kim et al. (2014) for a survey on the algorithms for NMF.

Next, we will survey the algorithms that work on noisy separable matrices. Four types of algorithm can be found:

- **AGKM (Arora et al., 2012a)**. The algorithm constructs  $r$  sets of data points such that all of the basis vectors are contained in the union of the sets and each set has one basis vector. The construction entails solving  $m$  linear programming (LP) problems with  $m - 1$  variables. Then, it chooses one element from each set, and outputs them.

- **Hottopixx** (Bittorf et al., 2012; Gillis, 2013; Gillis and Luce, 2013). Let  $\mathbf{A}$  be a separable matrix of the form  $\mathbf{F}(\mathbf{I}, \mathbf{K})\mathbf{\Pi}$ . Consider a matrix  $\mathbf{C}$  such that

$$\mathbf{C} = \mathbf{\Pi}^{-1} \begin{pmatrix} \mathbf{I} & \mathbf{K} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{\Pi} \in \mathbb{R}^{m \times m}.$$

It satisfies  $\mathbf{A} = \mathbf{AC}$ , and also, if the diagonal element is one, the position of its diagonal element indicates the index of basis vector in  $\mathbf{A}$ . The algorithm models  $\mathbf{C}$  as the variable of an LP problem. It entails solving a single LP problem with  $m^2$  variables.

- **SPA** (Gillis and Vavasis, 2014). Let  $\mathbf{A}$  be a separable matrix of size  $d$ -by- $m$ , and  $\mathcal{S}$  be the set of the column vectors of  $\mathbf{A}$ . The algorithm is based on the following observation. Under Assumption 1, the maximum of a convex function over the elements in  $\mathcal{S}$  is attained at the vertex of  $\text{conv}(\mathbf{A})$ . The algorithm finds one element  $\mathbf{a}$  in  $\mathcal{S}$  that maximizes a convex function, and then, projects all elements in  $\mathcal{S}$  into the orthogonal space to  $\mathbf{a}$ . This procedure is repeated until  $r$  elements are found.
- **XRAY** (Kumar et al., 2013). The algorithm has a similar spirit as SPA, but it uses a linear function instead of a convex one. Let  $\mathbf{A}$  be a separable matrix of size  $d$ -by- $m$  and  $\mathcal{S}$  be the set of the column vectors of  $\mathbf{A}$ . Let  $\mathcal{I}_k$  be the index set obtained after the  $k$ th iteration. This is a subset of  $\{1, \dots, m\}$  with cardinality  $k$ . In the  $(k + 1)$ th iteration, it computes a residual matrix  $\mathbf{R} = \mathbf{A}(\mathcal{I}_k)\mathbf{X}^* - \mathbf{A}$ , where

$$\mathbf{X}^* = \arg \min_{\mathbf{X} \geq \mathbf{0}} \|\mathbf{A}(\mathcal{I}_k)\mathbf{X} - \mathbf{A}\|_2^2,$$

and picks up one of the column vectors  $\mathbf{r}_i$  of  $\mathbf{R}$ . Then, it finds one element from  $\mathcal{S}$  which maximizes a linear function having  $\mathbf{r}_i$  as the normal vector. Finally,  $\mathcal{I}_k$  is updated by adding the index of the obtained element. This procedure is repeated until  $r$  indices are found. The performance of XRAY depends on how we select the column vector of the residual matrix  $\mathbf{R}$  for making the linear function. Several ways of selection, called “rand”, “max”, “dist” and “greedy”, have been proposed by the authors.

The next section describes the properties of these algorithms.

## 2.4 Comparison with Existing Algorithm

We compare the correctness and robustness properties of ER with those of AGKM, Hottopixx, SPA, and XRAY. ER is shown to have the two properties in Theorems 5 and 9. In particular, our robustness property in Theorem 9 states that ER correctly identifies the near-basis matrix of a noisy separable one  $\tilde{\mathbf{A}}$ , and a robustness property with  $\tau = 1$  holds if we set

$$\epsilon = \frac{\sigma(1 - \mu)}{4} \tag{3}$$

and  $p = 2$  under Assumption 1. Here,  $\sigma$  is the minimum singular value of the basis matrix  $\mathbf{F}$  of a separable one  $\mathbf{A}$  in the  $\tilde{\mathbf{A}}$ , that is,  $\sigma = \sigma_r(\mathbf{F})$ , and  $\mu$  is  $\mu(\mathbf{K})$ :

$$\mu(\mathbf{K}) = \max_{i=1, \dots, \ell} \|\mathbf{k}_i\|_2$$



for a weight matrix  $\mathbf{K}$  of  $\mathbf{A}$ . Under Assumption 1-a, we have  $\mu \leq 1$ , and in particular, equality holds if and only if  $\mathbf{k}_i$  has only one nonzero element.

All four of the existing algorithms have been shown to have a correctness property, whereas every one except XRAY has a robustness property. Hottopixx is the most similar to ER. Bittorf et al. (2012) showed that it has the correctness and robustness with  $\tau = 1$  properties if one sets

$$\epsilon = \frac{\alpha \min\{d_0, \alpha\}}{9(r+1)} \quad (4)$$

and  $p = 1$  under simpliciality and other assumptions. Here,  $\alpha$  and  $d_0$  are as follows.  $\alpha$  is the minimum value of  $\delta_{\mathbf{F}}(j)$  for  $j = 1, \dots, r$ , where  $\delta_{\mathbf{F}}(j)$  denotes an  $\ell_1$ -distance between the  $j$ th column vector  $\mathbf{f}_j$  of  $\mathbf{F}$  and the convex hull of the remaining column vectors in  $\mathbf{F}$ .  $d_0$  is the minimum value of  $\|\mathbf{a}_i - \mathbf{f}_j\|_1$  for every  $i$  such that  $\mathbf{a}_i$  is not a basis vector, and every  $j = 1, \dots, r$ . The robustness of Hottopixx is further analyzed (Gillis, 2013; Gillis and Luce, 2013).

It can be interpreted that the  $\epsilon$  of ER (3) is given by the multiplication of two parameters representing flatness and closeness of a given data matrix since  $\sigma$  measures the flatness of the convex hull of data points, and  $1 - \mu$  measures the closeness between basis vectors and data points. Intuitively, we may say that an algorithm becomes sensitive to noise when a data matrix has the following features; one is that the convex hull of data points is close to a flat shape, and another is that there are data points close to basis vectors. The  $\epsilon$  of (3) well matches the intuition. We see a similar structure in the  $\epsilon$  of Hottopixx (4) since  $\alpha$  and  $d_0$  respectively measure the flatness and closeness of a given data.

Compared with Hottopixx, the  $\epsilon$  of ER (3) does not contain  $1/r$ , and hence, it does not decrease as  $r$  increases. However, Assumption 1-b of ER is stronger than the simpliciality of Hottopixx. In a practical implementation, ER can handle a large matrix, while Hottopixx may have limitations on the size of the matrix it can handle. Hottopixx entails solving an LP problem with  $m^2$  variables. In the NMFs arising in applications,  $m$  tends to be a large number. Although an LP is tractable, it becomes harder to solve as the size increases. Through experiments, we assessed the performance of Hottopixx with the CPLEX LP solver. The experiments showed that the algorithm had out of memory issues when  $m$  exceeded 2,000 with  $d = 100$ . Bittorf et al. (2012) proposed a parallel implementation to resolve these computational issues.

AGKM and SPA were shown to have a robustness property with  $\tau \geq 1$  for some  $\epsilon$  in Arora et al. (2012a) and Gillis and Vavasis (2014), respectively. In practical implementations, SPA and XRAY are scalable to the problem size and experimentally show good robustness. Section 6 reports a numerical comparison of ER with SPA and XRAY.

### 3. Review of Formulation and Algorithm for MVEE Computation

We review the formulation for computing the MVEE for a set of points, and survey the existing algorithms for the computation.

First of all, let us recall the terminology related to an ellipsoid. An ellipsoid in  $\mathbb{R}^d$  is defined as a set  $\mathcal{E}(\mathbf{L}, \mathbf{z}) = \{\mathbf{x} \in \mathbb{R}^d : (\mathbf{x} - \mathbf{z})^\top \mathbf{L}(\mathbf{x} - \mathbf{z}) \leq 1\}$  for a positive definite matrix  $\mathbf{L}$  of size  $d$  and a vector  $\mathbf{z} \in \mathbb{R}^d$ . Here,  $\mathbf{L}$  determines the shape of the ellipsoid and  $\mathbf{z}$  is the center. Let  $\mathbf{x}$  be a point in an ellipsoid  $\mathcal{E}(\mathbf{L}, \mathbf{z})$ . If the point  $\mathbf{x}$  satisfies the equality

$(\mathbf{x} - \mathbf{z})^\top \mathbf{L}(\mathbf{x} - \mathbf{z}) = 1$ , we call it an *active point* of the ellipsoid. In other words, an active point is one lying on the boundary of the ellipsoid.

The volume of the ellipsoid is given as  $c(d)/\sqrt{\det \mathbf{L}}$ , where  $c(d)$  represents the volume of a unit ball in  $\mathbb{R}^d$  and it is a real number depending on the dimension  $d$ . ER algorithm considers  $d$ -dimensional ellipsoids containing a set  $\mathcal{S}$  of points in  $\mathbb{R}^d$ , and in particular, finds the minimum volume ellipsoid centered at the origin. In this paper, such an ellipsoid is referred to as an origin-centered MVEE for short.

Now, we are ready to describe a formulation for computing the origin-centered MVEE for a set of points. For  $m$  points  $\mathbf{p}_1, \dots, \mathbf{p}_m \in \mathbb{R}^d$ , let  $\mathcal{S} = \{\pm \mathbf{p}_1, \dots, \pm \mathbf{p}_m\}$ . The computation of the origin-centered MVEE for  $\mathcal{S}$  is formulated as

$$\begin{aligned} \mathbb{Q}(\mathcal{S}) : \quad & \text{minimize} && -\log \det \mathbf{L}, \\ & \text{subject to} && \langle \mathbf{p}_i \mathbf{p}_i^\top, \mathbf{L} \rangle \leq 1, \quad i = 1, \dots, m, \\ & && \mathbf{L} \succ \mathbf{0}, \end{aligned}$$

where the matrix  $\mathbf{L}$  of size  $d$  is the decision variable. The optimal solution  $\mathbf{L}^*$  of  $\mathbb{Q}$  gives the origin-centered MVEE for  $\mathcal{S}$  as  $\mathcal{E}(\mathbf{L}^*) = \{\mathbf{x} : \mathbf{x}^\top \mathbf{L}^* \mathbf{x} \leq 1\}$ . We here introduce some terminology. An active point of  $\mathcal{E}(\mathbf{L}^*)$  is a vector  $\mathbf{p}_i \in \mathbb{R}^d$  satisfying  $\mathbf{p}_i^\top \mathbf{L}^* \mathbf{p}_i = 1$ . We call  $\mathbf{p}_i$  an *active point of*  $\mathbb{Q}(\mathcal{S})$ , and the index  $i$  of  $\mathbf{p}_i$  an *active index of*  $\mathbb{Q}(\mathcal{S})$ . The ellipsoid  $\mathcal{E}(\mathbf{L}^*)$  is centrally symmetric, and if a vector  $\mathbf{p}_i$  is an active point, so is  $-\mathbf{p}_i$ . The dual of  $\mathbb{Q}$  reads

$$\begin{aligned} \mathbb{Q}_*(\mathcal{S}) : \quad & \text{maximize} && \log \det \Omega(\mathbf{u}), \\ & \text{subject to} && \mathbf{e}^\top \mathbf{u} = 1, \\ & && \mathbf{u} \geq \mathbf{0}, \end{aligned}$$

where the vector  $\mathbf{u}$  is the decision variable. Here,  $\Omega : \mathbb{R}^m \rightarrow \mathbb{S}^d$  is a linear function given as  $\Omega(\mathbf{u}) = \sum_{i=1}^m \mathbf{p}_i \mathbf{p}_i^\top u_i$ ; equivalently,  $\Omega(\mathbf{u}) = \mathbf{P} \text{diag}(\mathbf{u}) \mathbf{P}^\top$  for  $\mathbf{P} = (\mathbf{p}_1, \dots, \mathbf{p}_m) \in \mathbb{R}^{d \times m}$ . It follows from the Karush-Kuhn-Tucker (KKT) conditions for these problems that the optimal solution  $\mathbf{L}^*$  of  $\mathbb{Q}$  is represented by  $\frac{1}{d} \Omega(\mathbf{u}^*)^{-1}$  for the optimal solution  $\mathbf{u}^*$  of  $\mathbb{Q}_*$ . We make the following assumption to ensure the existence of an optimal solution of  $\mathbb{Q}$ .

**Assumption 2**  $\text{rank}(\mathbf{P}) = d$  for  $\mathbf{P} = (\mathbf{p}_1, \dots, \mathbf{p}_m) \in \mathbb{R}^{d \times m}$ .

Later, the KKT conditions will play an important role in our discussion of the active points of  $\mathbb{Q}$ . Here though, we will describe the conditions:  $\mathbf{L}^* \in \mathbb{S}^d$  is an optimal solution for  $\mathbb{Q}$  and  $\mathbf{z}^* \in \mathbb{R}^m$  is the associated Lagrange multiplier vector if and only if there exist  $\mathbf{L}^* \in \mathbb{S}^d$  and  $\mathbf{z}^* \in \mathbb{R}^m$  such that

$$-(\mathbf{L}^*)^{-1} + \Omega(\mathbf{z}^*) = \mathbf{0}, \tag{5}$$

$$z_i^* (\langle \mathbf{p}_i \mathbf{p}_i^\top, \mathbf{L}^* \rangle - 1) = 0, \quad i = 1, \dots, m, \tag{6}$$

$$\langle \mathbf{p}_i \mathbf{p}_i^\top, \mathbf{L}^* \rangle \leq 1, \quad i = 1, \dots, m, \tag{7}$$

$$\mathbf{L}^* \succ \mathbf{0}, \tag{8}$$

$$z_i^* \geq 0, \quad i = 1, \dots, m. \tag{9}$$

Many algorithms have been proposed for solving problems  $\mathbb{Q}$  and  $\mathbb{Q}_*$ . These can be categorized into mainly two types: conditional gradient algorithms (also referred to as

Frank-Wolfe algorithms) and interior-point algorithms. Below, we survey the studies on these two algorithms.

Khachiyan (1996) proposed a barycentric coordinate descent algorithm, which can be interpreted as a conditional gradient algorithm. He showed that the algorithm has a polynomial-time iteration complexity. Several researchers investigated and revised Khachiyan’s algorithm. Kumar and Yildirim (2005) showed that the iteration complexity of Khachiyan’s algorithm can be slightly reduced if it starts from a well-selected initial point. Todd and Yildirim (2007) and Ahipasaoglu et al. (2008) incorporated a step called as a Wolfe’s away-step. The revised algorithm was shown to have a polynomial-time iteration complexity and a linear convergence rate.

A dual interior-point algorithm was given by Vandenberghe et al. (1998). A primal-dual interior-point algorithm was given by Toh (1999), and numerical experiments showed that this algorithm is efficient and can provide accurate solutions. A practical algorithm was designed by Sun and Freund (2004) for solving large-scale problems. In particular, a hybrid of the interior-point algorithm and cutting plane algorithm was shown to be efficient in numerical experiments. For instance, the paper reported that the hybrid algorithm can solve problems with  $d = 30$  and  $m = 30,000$  in under 30 seconds on a personal computer. Tsuchiya and Xia (2007) considered generalized forms of  $\mathbb{Q}$  and  $\mathbb{Q}_*$  and showed that a primal-dual interior-point algorithm for the generalized forms has a polynomial-time iteration complexity.

Next, let us discuss the complexity of these two sorts of algorithms for  $\mathbb{Q}$  and  $\mathbb{Q}_*$ . In each iteration, the arithmetic operations of the conditional gradient algorithms are less than those of the interior-point algorithms. Each iteration of a conditional gradient algorithm (Khachiyan, 1996; Kumar and Yildirim, 2005; Todd and Yildirim, 2007; Ahipasaoglu et al., 2008) requires  $O(md)$  arithmetic operations. On the other hand, assuming that the number of data points  $m$  is sufficiently larger than the dimension of data points  $d$ , the main complexity of interior-point algorithms (Vandenberghe et al., 1998; Toh, 1999) comes from solving an  $m$ -by- $m$  system of linear equations in each iteration. The solution serves as the search direction for the next iteration. Solving these linear equations requires  $O(m^3)$  arithmetic operations. In practice, the number of iterations of conditional gradient algorithms is much larger than that of interior-point algorithms. Ahipasaoglu et al. (2008) reports that conditional gradient algorithms take several thousands iterations to solve problems such that  $d$  runs from 10 to 30 and  $m$  from 10,000 to 30,000. On the other hand, Sun and Freund (2004) reports that interior-point algorithms usually terminate after several dozen iterations and provide accurate solutions.

One of the concerns about interior-point algorithms is the computational cost of each iteration. It is possible to reduce the cost considerably by using a cutting plane strategy. A hybrid of interior-point algorithm and cutting plane algorithm has an advantage over conditional gradient algorithms. In fact, Ahipasaoglu et al. (2008) reports that the hybrid algorithm is faster than the conditional gradient algorithms and works well even on large problems. Therefore, we use the hybrid algorithm to solve  $\mathbb{Q}$  in our practical implementation of ER. The details are in Section 5.1.

Here, it should be mentioned that this paper uses a terminology “cutting plane strategy” for what other papers (Sun and Freund, 2004; Ahipasaoglu et al., 2008) have called the

“active set strategy”, since it might be confused with “active set algorithm” for solving a nonnegative least square problem.

#### 4. Description and Analysis of the Algorithm

The ER algorithm is presented below. Throughout of this paper, we use the notation  $\mathbb{N}$  to denote a set of nonnegative integer numbers.

---

**Algorithm 1** Ellipsoidal Rounding (ER) for Problem 1

---

**Input:**  $M \in \mathbb{R}_+^{d \times m}$  and  $r \in \mathbb{N}$ .

**Output:**  $\mathcal{I}$ .

- 1: Compute the SVD of  $M$ , and construct the reduced matrix  $P \in \mathbb{R}^{r \times m}$  associated with  $r$ .
  - 2: Let  $\mathcal{S} = \{\pm \mathbf{p}_1, \dots, \pm \mathbf{p}_m\}$  for the column vectors  $\mathbf{p}_1, \dots, \mathbf{p}_m$  of  $P$ . Solve  $\mathbb{Q}(\mathcal{S})$ , and construct the active index set  $\mathcal{I}$ .
- 

Step 1 needs to be explained in detail. Let  $M$  be a noisy separable matrix of size  $d$ -by- $m$ . In general, the  $M$  is a full-rank due to the existence of a noise matrix. However, the rank is close to  $r$  when the amount of noise is small, and in particular, it is  $r$  in the noiseless case. Accordingly, we construct a low-rank approximation matrix to  $M$  and reduce the redundancy in the space spanned by the column vectors of  $M$ .

We use an SVD for the construction of the low-rank approximation matrix. The SVD of  $M$  gives a decomposition of the form,

$$M = U \Sigma V^\top.$$

Here,  $U$  and  $V$  are  $d$ -by- $d$  and  $m$ -by- $m$  orthogonal matrices, respectively. In this paper, we call the  $U$  a *left orthogonal matrix* of the SVD of  $M$ . Let  $t = \min\{d, m\}$ .  $\Sigma$  is a rectangular diagonal matrix consisting of the singular values  $\sigma_1, \dots, \sigma_t$  of  $M$ , and it is of the form,

$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_t) \in \mathbb{R}^{d \times m}$$

with  $\sigma_1 \geq \dots \geq \sigma_t \geq 0$ . By choosing the top  $r$  singular values while setting the others to 0 in  $\Sigma$ , we construct

$$\Sigma^r = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0) \in \mathbb{R}^{d \times m}$$

and let

$$M^r = U \Sigma^r V^\top.$$

$M^r$  is the best rank- $r$  approximation to  $M$  as measured by the matrix 2-norm and satisfies  $\|M - M^r\|_2 = \sigma_{r+1}$  (see, for instance, Theorem 2.5.3 of Golub and Loan, 1996). By applying the left orthogonal matrix  $U^\top$  to  $M^r$ , we have

$$U^\top M^r = \begin{pmatrix} P \\ \mathbf{0} \end{pmatrix} \in \mathbb{R}^{d \times m},$$

where  $P$  is an  $r$ -by- $m$  matrix with  $\text{rank}(P) = r$ . We call such a matrix  $P$  a *reduced matrix of  $M$  associated with  $r$* . Since Assumption 2 holds for the  $P$ , it is possible to perform an MVEE computation for a set of the column vectors.

#### 4.1 Correctness for a Separable Matrix

We analyze the correctness property of Algorithm 1. Let  $\mathbf{A}$  be a separable matrix of size  $d$ -by- $m$ . Assume that Assumption 1 holds for  $\mathbf{A}$ . We run Algorithm 1 for  $(\mathbf{A}, \text{rank}(\mathbf{A}))$ . Step 1 computes the reduced matrix  $\mathbf{P}$  of  $\mathbf{A}$ . Since  $r = \text{rank}(\mathbf{A})$ , we have  $\mathbf{A} = \mathbf{A}^r$ , where  $\mathbf{A}^r$  is the best rank- $r$  approximation matrix to  $\mathbf{A}$ . Let  $\mathbf{U} \in \mathbb{R}^{d \times d}$  be the left orthogonal matrix of the SVD of  $\mathbf{A}$ . The reduced matrix  $\mathbf{P} \in \mathbb{R}^{r \times m}$  of  $\mathbf{A}$  is obtained as

$$\begin{aligned} \begin{pmatrix} \mathbf{P} \\ \mathbf{0} \end{pmatrix} &= \mathbf{U}^\top \mathbf{A} \\ &= \mathbf{U}^\top \mathbf{F}(\mathbf{I}, \mathbf{K})\mathbf{\Pi}. \end{aligned} \quad (10)$$

From the above, we see that

$$\mathbf{U}^\top \mathbf{F} = \begin{pmatrix} \mathbf{G} \\ \mathbf{0} \end{pmatrix} \in \mathbb{R}^{d \times m}, \text{ where } \mathbf{G} \in \mathbb{R}^{r \times r}. \quad (11)$$

Here, we have  $\text{rank}(\mathbf{G}) = r$  since  $\text{rank}(\mathbf{F}) = r$  by Assumption 1-b and  $\mathbf{U}$  is an orthogonal matrix. By using  $\mathbf{G}$ , we rewrite  $\mathbf{P}$  as

$$\mathbf{P} = (\mathbf{G}, \mathbf{G}\mathbf{K})\mathbf{\Pi}.$$

From Assumption 1-a, the column vectors  $\mathbf{k}_i$  of the weight matrix  $\mathbf{K} \in \mathbb{R}^{r \times \ell}$  satisfy the conditions

$$\|\mathbf{k}_i\|_1 = 1 \text{ and } \mathbf{k}_i \geq \mathbf{0}, \quad i = 1, \dots, \ell. \quad (12)$$

In Step 2, we collect the column vectors of  $\mathbf{P}$  and construct a set  $\mathcal{S}$  of them. Let  $\mathbf{B} = \mathbf{G}\mathbf{K}$ , and let  $\mathbf{g}_j$  and  $\mathbf{b}_i$  be the column vector of  $\mathbf{G}$  and  $\mathbf{B}$ , respectively.  $\mathcal{S}$  is a set of vectors  $\pm \mathbf{g}_1, \dots, \pm \mathbf{g}_r, \pm \mathbf{b}_1, \dots, \pm \mathbf{b}_\ell$ . The following proposition guarantees that the active points of  $\mathbb{Q}(\mathcal{S})$  are  $\mathbf{g}_1, \dots, \mathbf{g}_r$ . We can see from (10) and (11) that the index set of the column vectors of  $\mathbf{G}$  is identical to that of  $\mathbf{F}$ . Hence, the basis matrix  $\mathbf{F}$  of a separable one  $\mathbf{A}$  can be obtained by finding the active points of  $\mathbb{Q}(\mathcal{S})$ .

**Proposition 3** *Let  $\mathbf{G} \in \mathbb{R}^{r \times r}$  and  $\mathbf{B} = \mathbf{G}\mathbf{K} \in \mathbb{R}^{r \times \ell}$  for  $\mathbf{K} \in \mathbb{R}^{r \times \ell}$ . For the column vectors  $\mathbf{g}_j$  and  $\mathbf{b}_i$  of  $\mathbf{G}$  and  $\mathbf{B}$ , respectively, let  $\mathcal{S} = \{\pm \mathbf{g}_1, \dots, \pm \mathbf{g}_r, \pm \mathbf{b}_1, \dots, \pm \mathbf{b}_\ell\}$ . Suppose that  $\text{rank}(\mathbf{G}) = r$  and  $\mathbf{K}$  satisfies the condition (12). Then, the active point set of  $\mathbb{Q}(\mathcal{S})$  is  $\{\mathbf{g}_1, \dots, \mathbf{g}_r\}$ .*

**Proof** We show that an optimal solution  $\mathbf{L}^*$  of  $\mathbb{Q}(\mathcal{S})$  is  $(\mathbf{G}\mathbf{G}^\top)^{-1}$  and its associated Lagrange multiplier  $\mathbf{z}^*$  is  $(\mathbf{e}; \mathbf{0})$ , where  $\mathbf{e}$  is an  $r$ -dimensional all-ones vector and  $\mathbf{0}$  is an  $\ell$ -dimensional zero vector. Here, the Lagrange multipliers are one for the constraints  $\langle \mathbf{g}_j \mathbf{g}_j^\top, \mathbf{L} \rangle \leq 1$ , and these are zero for  $\langle \mathbf{b}_i \mathbf{b}_i^\top, \mathbf{L} \rangle \leq 1$ .

Since  $\mathbf{G}$  is nonsingular, the inverse of  $\mathbf{G}\mathbf{G}^\top$  exists and it is positive definite. Now we check that  $\mathbf{L}^* = (\mathbf{G}\mathbf{G}^\top)^{-1}$  and  $\mathbf{z}^* = (\mathbf{e}; \mathbf{0})$  satisfy the KKT conditions (5)-(9) for the problem. It was already seen that the conditions (5), (8), and (9) are satisfied. For the remaining conditions, we have

$$\langle \mathbf{g}_j \mathbf{g}_j^\top, (\mathbf{G}\mathbf{G}^\top)^{-1} \rangle = (\mathbf{G}^\top (\mathbf{G}\mathbf{G}^\top)^{-1} \mathbf{G})_{jj} = 1 \quad (13)$$

and

$$\begin{aligned}
 \langle \mathbf{b}_i \mathbf{b}_i^\top, (\mathbf{G}\mathbf{G}^\top)^{-1} \rangle &= (\mathbf{B}^\top (\mathbf{G}\mathbf{G}^\top)^{-1} \mathbf{B})_{ii} \\
 &= (\mathbf{K}^\top \mathbf{G}^\top (\mathbf{G}\mathbf{G}^\top)^{-1} \mathbf{G}\mathbf{K})_{ii} \\
 &= \mathbf{k}_i^\top \mathbf{k}_i \\
 &\leq \|\mathbf{k}_i\|_1^2 = 1.
 \end{aligned} \tag{14}$$

Here,  $(\cdot)_{ii}$  for a matrix denotes the  $(i, i)$ th element of the matrix. The inequality in (14) follows from condition (12). Also, the Lagrange multipliers are zero for the inequality constraints  $\langle \mathbf{b}_i \mathbf{b}_i^\top, (\mathbf{G}\mathbf{G}^\top)^{-1} \rangle \leq 1$ . Thus, conditions (6) and (7) are satisfied. Accordingly,  $(\mathbf{G}\mathbf{G}^\top)^{-1}$  is an optimal solution of  $\mathbb{Q}(\mathcal{S})$ .

We can see from (13) that  $\mathbf{g}_1, \dots, \mathbf{g}_r$  are the active points of the problem. Moreover, we may have equality in (14). In fact, equality holds if and only if  $\mathbf{k}_i$  has only one nonzero element. For such  $\mathbf{k}_i$ ,  $\mathbf{b}_i = \mathbf{G}\mathbf{k}_i$  coincides with some vector in  $\mathbf{g}_1, \dots, \mathbf{g}_r$ . ■

From the above discussion, we can immediately notice that this proposition holds if for a matrix  $\mathbf{K} \in \mathbb{R}^{r \times \ell}$ , the column vectors  $\mathbf{k}_i$  satisfy

$$\|\mathbf{k}_i\|_2 < 1, \quad i = 1, \dots, m. \tag{15}$$

Note that in contrast with condition (12), this condition does not require the matrix to be nonnegative.

**Corollary 4** *Proposition 3 holds even if we suppose that  $\mathbf{K} \in \mathbb{R}^{r \times \ell}$  satisfies condition (15), instead of condition (12).*

Note that this corollary is used to show the robustness of Algorithm 1 on a noisy separable matrix. The correctness of Algorithm 1 for a separable matrix follows from the above discussion and Proposition 3.

**Theorem 5** *Let  $\mathbf{A}$  be a separable matrix. Assume that Assumption 1 holds for  $\mathbf{A}$ . Then, Algorithm 1 for  $(\mathbf{A}, \text{rank}(\mathbf{A}))$  returns an index set  $\mathcal{I}$  such that  $\mathbf{A}(\mathcal{I}) = \mathbf{F}$ .*

## 4.2 Robustness for a Noisy Separable Matrix

Next, we analyze the robustness property of Algorithm 1. Let  $\mathbf{A}$  be a separable matrix of size  $d$ -by- $m$ . Assume that Assumption 1 holds for  $\mathbf{A}$ . Let  $\tilde{\mathbf{A}}$  be a noisy separable matrix of the form  $\tilde{\mathbf{A}} + \mathbf{N}$ . We run Algorithm 1 for  $(\tilde{\mathbf{A}}, \text{rank}(\mathbf{A}))$ . Step 1 computes the reduced matrix  $\tilde{\mathbf{P}}$  of  $\tilde{\mathbf{A}}$ . Let  $\tilde{\mathbf{U}} \in \mathbb{R}^{d \times d}$  be the left orthogonal matrix of the SVD of  $\tilde{\mathbf{A}}$ , and  $\tilde{\mathbf{A}}^r$  be the best rank- $r$  approximation matrix to  $\tilde{\mathbf{A}}$ . We denote the residual matrix  $\tilde{\mathbf{A}} - \tilde{\mathbf{A}}^r$  by  $\tilde{\mathbf{A}}_r^\perp$ .

For the reduced matrix  $\mathbf{P}$  of  $\tilde{\mathbf{A}}$ , we have

$$\begin{aligned} \begin{pmatrix} \mathbf{P} \\ \mathbf{0} \end{pmatrix} &= \mathbf{U}^\top \tilde{\mathbf{A}}^r \\ &= \mathbf{U}^\top (\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\diamond^r) \end{aligned} \quad (16)$$

$$\begin{aligned} &= \mathbf{U}^\top (\mathbf{A} + \mathbf{N} - \tilde{\mathbf{A}}_\diamond^r) \\ &= \mathbf{U}^\top (\mathbf{A} + \bar{\mathbf{N}}) \end{aligned} \quad (17)$$

$$\begin{aligned} &= \mathbf{U}^\top ((\mathbf{F}, \mathbf{F}\mathbf{K})\mathbf{\Pi} + \bar{\mathbf{N}}) \\ &= \mathbf{U}^\top (\mathbf{F} + \bar{\mathbf{N}}^{(1)}, \mathbf{F}\mathbf{K} + \bar{\mathbf{N}}^{(2)})\mathbf{\Pi} \end{aligned} \quad (18)$$

$$= \mathbf{U}^\top (\hat{\mathbf{F}}, \hat{\mathbf{F}}\mathbf{K} + \hat{\mathbf{N}})\mathbf{\Pi}. \quad (19)$$

The following notation is used in the above:  $\bar{\mathbf{N}} = \mathbf{N} - \tilde{\mathbf{A}}_\diamond^r$  in (17);  $\bar{\mathbf{N}}^{(1)}$  and  $\bar{\mathbf{N}}^{(2)}$  in (18) are the  $d$ -by- $r$  and  $d$ -by- $\ell$  submatrices of  $\bar{\mathbf{N}}$  such that  $\bar{\mathbf{N}}\mathbf{\Pi}^{-1} = (\bar{\mathbf{N}}^{(1)}, \bar{\mathbf{N}}^{(2)})$ ;  $\hat{\mathbf{F}} = \mathbf{F} + \bar{\mathbf{N}}^{(1)}$  and  $\hat{\mathbf{N}} = -\bar{\mathbf{N}}^{(1)}\mathbf{K} + \bar{\mathbf{N}}^{(2)}$  in (19). This implies that

$$\mathbf{U}^\top \hat{\mathbf{F}} = \begin{pmatrix} \hat{\mathbf{G}} \\ \mathbf{0} \end{pmatrix}, \text{ where } \hat{\mathbf{G}} \in \mathbb{R}^{r \times r}, \quad (20)$$

and

$$\mathbf{U}^\top \hat{\mathbf{N}} = \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix}, \text{ where } \mathbf{R} \in \mathbb{R}^{r \times \ell}. \quad (21)$$

Hence, we can rewrite  $\mathbf{P}$  as

$$\mathbf{P} = (\hat{\mathbf{G}}, \hat{\mathbf{G}}\mathbf{K} + \mathbf{R})\mathbf{\Pi}.$$

$\tilde{\mathbf{A}}$  is represented by (2) as

$$\tilde{\mathbf{A}} = (\tilde{\mathbf{F}}, \tilde{\mathbf{F}}\mathbf{K} + \tilde{\mathbf{N}})\mathbf{\Pi},$$

where  $\tilde{\mathbf{F}}$  and  $\tilde{\mathbf{N}}$  denote  $\mathbf{F} + \mathbf{N}^{(1)}$  and  $-\mathbf{N}^{(1)}\mathbf{K} + \mathbf{N}^{(2)}$ , respectively. From (16), we have

$$\begin{pmatrix} (\hat{\mathbf{G}}, \hat{\mathbf{G}}\mathbf{K} + \mathbf{R})\mathbf{\Pi} \\ \mathbf{0} \end{pmatrix} = \mathbf{U}^\top ((\tilde{\mathbf{F}}, \tilde{\mathbf{F}}\mathbf{K} + \tilde{\mathbf{N}})\mathbf{\Pi} - \tilde{\mathbf{A}}_\diamond^r).$$

Therefore, the index set of the column vectors of  $\hat{\mathbf{G}}$  is identical to that of  $\tilde{\mathbf{F}}$ . If all the column vectors of  $\hat{\mathbf{G}}$  are found in  $\mathbf{P}$ , we can identify  $\tilde{\mathbf{F}}$  hidden in  $\tilde{\mathbf{A}}$ .

In Step 2, we collect the column vectors of  $\mathbf{P}$  and construct a set  $\mathcal{S}$  of them. Let

$$\hat{\mathbf{B}} = \hat{\mathbf{G}}\mathbf{K} + \mathbf{R}, \quad (22)$$

and let  $\hat{\mathbf{g}}_j$  and  $\hat{\mathbf{b}}_i$  respectively be the column vectors of  $\hat{\mathbf{G}}$  and  $\hat{\mathbf{B}}$ .  $\mathcal{S}$  is a set of vectors  $\pm\hat{\mathbf{g}}_1, \dots, \pm\hat{\mathbf{g}}_r, \pm\hat{\mathbf{b}}_1, \dots, \pm\hat{\mathbf{b}}_\ell$ . We can see from Corollary 4 that, if  $\text{rank}(\hat{\mathbf{G}}) = r$  and  $\hat{\mathbf{b}}_i$  is written as  $\hat{\mathbf{b}}_i = \hat{\mathbf{G}}\hat{\mathbf{k}}_i$  by using  $\hat{\mathbf{k}}_i \in \mathbb{R}^r$  with  $\|\hat{\mathbf{k}}_i\|_2 < 1$ , the active points of  $\mathbb{Q}(\mathcal{S})$  are given as the column vectors  $\hat{\mathbf{g}}_1, \dots, \hat{\mathbf{g}}_r$  of  $\hat{\mathbf{G}}$ . Below, we examine the amount of noise  $\mathbf{N}$  such that the conditions of Corollary 4 still hold.

**Lemma 6** *Let  $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{N} \in \mathbb{R}^{d \times m}$ . Then,  $|\sigma_i(\tilde{\mathbf{A}}) - \sigma_i(\mathbf{A})| \leq \|\mathbf{N}\|_2$  for each  $i = 1, \dots, t$  where  $t = \min\{d, m\}$ .*

**Proof** See Corollary 8.6.2 of Golub and Loan (1996). ■

**Lemma 7** Let  $n = \|\mathbf{N}\|_2$  and  $\mu = \mu(\mathbf{K})$ .

7-a) The matrix  $\bar{\mathbf{N}}$  of (17) satisfies  $\|\bar{\mathbf{N}}\|_2 \leq 2n$ .

7-b) The column vectors  $\mathbf{r}_i$  of matrix  $\mathbf{R}$  of (21) satisfy  $\|\mathbf{r}_i\|_2 \leq 2n(\mu + 1)$  for  $i = 1, \dots, m$ .

7-c) The singular values of matrix  $\hat{\mathbf{G}}$  of (20) satisfy  $|\sigma_i(\hat{\mathbf{G}}) - \sigma_i(\mathbf{F})| \leq 2n$  for  $i = 1, \dots, r$ .

**Proof** 7-a) Since  $\bar{\mathbf{N}} = \mathbf{N} - \tilde{\mathbf{A}}_\diamond^r$ ,

$$\|\bar{\mathbf{N}}\|_2 \leq \|\mathbf{N}\|_2 + \|\tilde{\mathbf{A}}_\diamond^r\|_2.$$

We have  $\|\tilde{\mathbf{A}}_\diamond^r\|_2 \leq n$  since  $\|\tilde{\mathbf{A}}_\diamond^r\|_2 = \sigma_{r+1}(\tilde{\mathbf{A}})$  and from Lemma 6,  $|\sigma_{r+1}(\tilde{\mathbf{A}}) - \sigma_{r+1}(\mathbf{A})| \leq n$ . Therefore,  $\|\bar{\mathbf{N}}\|_2 \leq 2n$ .

7-b) Let  $\hat{\mathbf{n}}_i$  be the column vector of the matrix  $\hat{\mathbf{N}}$  of (19). Since  $\mathbf{U}^\top \hat{\mathbf{n}}_i = (\mathbf{r}_i; \mathbf{0})$  for an orthogonal matrix  $\mathbf{U}$ , we have  $\|\hat{\mathbf{n}}_i\|_2 = \|\mathbf{r}_i\|_2$ . Therefore, we will evaluate  $\|\hat{\mathbf{n}}_i\|_2$ . Let  $\mathbf{k}_i$  and  $\bar{\mathbf{n}}_i^{(2)}$  be the column vectors of  $\mathbf{K}$  and  $\bar{\mathbf{N}}^{(2)}$ , respectively. Then,  $\hat{\mathbf{n}}_i$  can be represented as  $-\bar{\mathbf{N}}^{(1)}\mathbf{k}_i + \bar{\mathbf{n}}_i^{(2)}$ . Thus, by Lemma 7-a, we have

$$\|\mathbf{r}_i\|_2 = \|\hat{\mathbf{n}}_i\|_2 \leq \|\bar{\mathbf{N}}^{(1)}\|_2 \|\mathbf{k}_i\|_2 + \|\bar{\mathbf{n}}_i^{(2)}\|_2 \leq 2n(\mu + 1).$$

7-c) Since  $\mathbf{U}^\top \hat{\mathbf{F}} = (\hat{\mathbf{G}}; \mathbf{0})$  for an orthogonal matrix  $\mathbf{U}$ , the singular values of  $\hat{\mathbf{F}}$  and  $\hat{\mathbf{G}}$  are identical. Also, since  $\hat{\mathbf{F}} = \mathbf{F} + \bar{\mathbf{N}}^{(1)}$  and Lemma 6, we have

$$|\sigma_i(\hat{\mathbf{G}}) - \sigma_i(\mathbf{F})| = |\sigma_i(\hat{\mathbf{F}}) - \sigma_i(\mathbf{F})| \leq \|\bar{\mathbf{N}}^{(1)}\|_2 \leq 2n. \quad \blacksquare$$

The following lemma ensures that the conditions of Corollary 4 hold if the amount of noise is smaller than a certain level.

**Lemma 8** Let  $\hat{\mathbf{G}}$  be the matrix of (20), and let  $\hat{\mathbf{b}}_i$  be the column vector of  $\hat{\mathbf{B}}$  of (22). Suppose that  $\|\mathbf{N}\|_2 < \epsilon$  for  $\epsilon = \frac{1}{4}\sigma(1 - \mu)$  where  $\sigma = \sigma_r(\mathbf{F})$  and  $\mu = \mu(\mathbf{K})$ . Then,

8-a)  $\text{rank}(\hat{\mathbf{G}}) = r$ .

8-b)  $\hat{\mathbf{b}}_i$  is represented as  $\hat{\mathbf{G}}\hat{\mathbf{k}}_i = \hat{\mathbf{b}}_i$  by using  $\hat{\mathbf{k}}_i$  such that  $\|\hat{\mathbf{k}}_i\|_2 < 1$ .

In the proof below,  $n$  denotes  $\|\mathbf{N}\|_2$ .

**Proof** 8-a) From Lemma 7-c, the minimum singular value of  $\hat{\mathbf{G}}$  satisfies

$$\begin{aligned} \sigma_r(\hat{\mathbf{G}}) &\geq \sigma - 2n \\ &> \sigma - 2\epsilon = \frac{1}{2}\sigma(1 + \mu) > 0. \end{aligned}$$



The final inequality follows from  $\sigma > 0$  due to Assumption 1-b. Hence, we have  $\text{rank}(\widehat{\mathbf{G}}) = r$ .

8-b) Let  $\mathbf{k}_i$  and  $\mathbf{r}_i$  be the column vectors of  $\mathbf{K}$  and  $\mathbf{R}$ , respectively. Then, we have  $\widehat{\mathbf{b}}_i = \widehat{\mathbf{G}}\mathbf{k}_i + \mathbf{r}_i$ . Since Lemma 8-a guarantees that  $\widehat{\mathbf{G}}$  has an inverse, it can be represented as  $\widehat{\mathbf{b}}_i = \widehat{\mathbf{G}}\widehat{\mathbf{k}}_i$  by  $\widehat{\mathbf{k}}_i = \mathbf{k}_i + \widehat{\mathbf{G}}^{-1}\mathbf{r}_i$ . It follows from Lemmas 7-b and 7-c that

$$\begin{aligned} \|\widehat{\mathbf{k}}_i\|_2 &\leq \|\mathbf{k}_i\|_2 + \|\widehat{\mathbf{G}}^{-1}\|_2\|\mathbf{r}_i\|_2 \\ &\leq \mu + \frac{2n(\mu+1)}{\sigma-2n}. \end{aligned}$$

Since  $n < \frac{1}{4}\sigma(1-\mu)$ , we have  $\|\widehat{\mathbf{k}}_i\|_2 < 1$ . ■

The robustness of Algorithm 1 for a noisy separable matrix follows from the above discussion, Corollary 4, and Lemma 8.

**Theorem 9** *Let  $\widetilde{\mathbf{A}}$  be a noisy separable matrix of the form  $\mathbf{A} + \mathbf{N}$ . Assume that Assumption 1 holds for the separable matrix  $\mathbf{A}$  in  $\widetilde{\mathbf{A}}$ . Set  $\epsilon = \frac{1}{4}\sigma(1-\mu)$  where  $\sigma = \sigma_r(\mathbf{F})$  and  $\mu = \mu(\mathbf{K})$  for the basis and weight matrices  $\mathbf{F}$  and  $\mathbf{K}$  of  $\mathbf{A}$ . If  $\|\mathbf{N}\|_2 < \epsilon$ , Algorithm 1 for  $(\widetilde{\mathbf{A}}, \text{rank}(\mathbf{A}))$  returns an index set  $\mathcal{I}$  such that  $\|\widetilde{\mathbf{A}}(\mathcal{I}) - \mathbf{F}\|_2 < \epsilon$ .*

In Theorem 9, let  $\mathbf{F}^* = \widetilde{\mathbf{A}}(\mathcal{I})$ , and  $\mathbf{W}^*$  be an optimal solution of the convex optimization problem,

$$\text{minimize } \|\widetilde{\mathbf{A}}(\mathcal{I})\mathbf{X} - \widetilde{\mathbf{A}}\|_{\mathbf{F}}^2 \text{ subject to } \mathbf{X} \geq \mathbf{0},$$

where the matrix  $\mathbf{X}$  of size  $r$ -by- $m$  is the decision variable. Then,  $(\mathbf{F}^*, \mathbf{W}^*)$  serves as the NMF factor of  $\widetilde{\mathbf{A}}$ . It is possible to evaluate the residual error of this factorization in a similar way to the proof of Theorem 4 by Gillis and Vavasis (2014).

**Corollary 10** *Let  $\mathbf{w}_i^*$  and  $\widetilde{\mathbf{a}}_i$  be the column vectors of  $\mathbf{W}^*$  and  $\widetilde{\mathbf{A}}$ , respectively. Then,  $\|\mathbf{F}^*\mathbf{w}_i^* - \widetilde{\mathbf{a}}_i\|_2 < 2\epsilon$  for  $i = 1, \dots, m$ .*

**Proof** From Assumption 1-a, the column vectors  $\mathbf{w}_i$  of  $\mathbf{W}$  satisfy  $\|\mathbf{w}_i\|_2 \leq 1$  for  $i = 1, \dots, m$ . Therefore, for  $i = 1, \dots, m$ ,

$$\begin{aligned} \|\mathbf{F}^*\mathbf{w}_i^* - \widetilde{\mathbf{a}}_i\|_2 &\leq \|\mathbf{F}^*\mathbf{w}_i - \widetilde{\mathbf{a}}_i\|_2 \\ &= \|\mathbf{F}^*\mathbf{w}_i - \mathbf{F}\mathbf{w}_i + \mathbf{F}\mathbf{w}_i - \mathbf{a}_i - \mathbf{n}_i\|_2 \\ &= \|(\mathbf{F}^* - \mathbf{F})\mathbf{w}_i - \mathbf{n}_i\|_2 \\ &\leq \|\mathbf{F}^* - \mathbf{F}\|_2\|\mathbf{w}_i\|_2 + \|\mathbf{n}_i\|_2 < 2\epsilon, \end{aligned}$$

where  $\mathbf{a}_i$  and  $\mathbf{n}_i$  denote the  $i$ th column vector of  $\mathbf{A}$  and  $\mathbf{N}$ , respectively. ■

## 5. Implementation in Practice

Theorem 9 guarantees that Algorithm 1 correctly identifies the near-basis matrix of a noisy separable matrix if the noise is smaller than some level. But in the NMFs of matrices arising from practical applications, it seems that the noise level would likely exceed the level for

which the theorem is valid. In such a situation, the algorithm might generate more active points than hoped. Therefore, we need to add a selection step in which  $r$  points are selected from the active points. Also, the number of active points depends on which dimension we choose in the computation of the reduced matrix  $\mathbf{P}$ . Algorithm 1 computes the reduced matrix  $\mathbf{P}$  of the data matrix and draws an origin-centered MVEE for the column vectors  $\mathbf{p}_1, \dots, \mathbf{p}_m$  of  $\mathbf{P}$ . As we will see in Lemma 11, the number of active points depends on the dimension of  $\mathbf{p}_1, \dots, \mathbf{p}_m$ . Therefore, we introduce an input parameter  $\rho$  to control the dimension. By taking account of these considerations, we design a practical implementation of Algorithm 1.

---

**Algorithm 2** Practical Implementation of Algorithm 1

---

**Input:**  $\mathbf{M} \in \mathbb{R}_+^{d \times m}$ ,  $r \in \mathbb{N}$ , and  $\rho \in \mathbb{N}$ .

**Output:**  $\mathcal{I}$ .

- 1: Run Algorithm 1 for  $(\mathbf{M}, \rho)$ . Let  $\mathcal{J}$  be the index set returned by the algorithm.
  - 2: If  $|\mathcal{J}| < r$ , increase  $\rho$  by 1 and go back to Step 1. Otherwise, select  $r$  elements from  $\mathcal{J}$  and construct the set  $\mathcal{I}$  of these elements.
- 

One may wonder whether Algorithm 2 infinitely loops or not. In fact, we can show that under some conditions, infinite loops do not occur.

**Lemma 11** For  $\mathbf{p}_1, \dots, \mathbf{p}_m \in \mathbb{R}^\rho$ , let  $\mathcal{S} = \{\pm \mathbf{p}_1, \dots, \pm \mathbf{p}_m\}$ . Suppose that Assumption 2 holds. Then,  $\mathbb{Q}(\mathcal{S})$  has at least  $\rho$  active points.

**Proof** Consider the KKT conditions (5)-(9) for  $\mathbb{Q}(\mathcal{S})$ . Condition (5) requires  $\Omega(\mathbf{z}^*)$  to be nonsingular. Since  $\text{rank}(\mathbf{P}) = \rho$  from the assumption, at least  $\rho$  nonzero  $z_i^*$  exist. Therefore, we see from condition (6) that  $\mathbb{Q}(\mathcal{S})$  has at least  $\rho$  active points. ■

**Proposition 12** Suppose that we choose  $r$  such that  $r \leq \text{rank}(\mathbf{M})$ . Then, Algorithm 2 terminates after a finite number of iterations.

**Proof** For the active index set  $\mathcal{J}$  constructed in Step 1, Lemma 11 guarantees that  $|\mathcal{J}| \geq \rho$ . The parameter  $\rho$  increases by 1 if  $|\mathcal{J}| < r$  in Step 2 and can continue to increase up to  $\rho = \text{rank}(\mathbf{M})$ . Since  $r \leq \text{rank}(\mathbf{M})$ , it is necessarily to satisfy  $|\mathcal{J}| \geq \rho \geq r$  after a finite number of iterations. ■

Proposition 12 implies that  $\rho$  may not be an essential input parameter since Algorithm 2 always terminates under  $r \leq \text{rank}(\mathbf{M})$  even if starting with  $\rho = 1$ .

There are some concerns about Algorithm 2. One is in how to select  $r$  elements from an active index set  $\mathcal{J}$  in Step 2. It is possible to have various ways to make the selection. We rely on existing algorithms, such as XRAY and SPA, and perform these existing algorithms for  $(\mathbf{M}(\mathcal{J}), \rho)$ . Thus, Algorithm 1 can be regarded as a preprocessor which filters out basis vector candidates from the data points and enhance the performance of existing algorithms. Another concern is in the computational cost of solving  $\mathbb{Q}$ . In the next section, we describe a cutting plane strategy for efficiently performing an interior-point algorithm.

### 5.1 Cutting Plane Strategy for Solving $\mathbb{Q}$

Let  $\mathcal{S}$  be a set of  $m$  points in  $\mathbb{R}^d$ . As mentioned in Section 3,  $O(m^3)$  arithmetic operations are required in each iteration of an interior-point algorithm for  $\mathbb{Q}(\mathcal{S})$ . A cutting plane strategy is a way to reduce the number of points which we need to deal with in solving  $\mathbb{Q}(\mathcal{S})$ . The strategy was originally used by Sun and Freund (2004). In this section, we describe the details of our implementation.

The cutting plane strategy for solving  $\mathbb{Q}$  has a geometric interpretation. It is thought of that active points contribute a lot to the drawing the MVEE for a set of points but inactive points make less of a contribution. This geometric intuition can be justified by the following proposition. Let  $\mathbf{L}$  be a  $d$ -by- $d$  matrix. We use the notation  $\delta_{\mathbf{L}}(\mathbf{p})$  to denote  $\langle \mathbf{p}\mathbf{p}^\top, \mathbf{L} \rangle$  for an element  $\mathbf{p} \in \mathbb{R}^d$  of  $\mathcal{S}$ .

**Proposition 13** *Let  $\bar{\mathcal{S}}$  be a subset of  $\mathcal{S}$ . If an optimal solution  $\bar{\mathbf{L}}^*$  of  $\mathbb{Q}(\bar{\mathcal{S}})$  satisfies  $\delta_{\bar{\mathbf{L}}^*}(\mathbf{p}) \leq 1$  for all  $\mathbf{p} \in \mathcal{S} \setminus \bar{\mathcal{S}}$ , then  $\bar{\mathbf{L}}^*$  is an optimal solution of  $\mathbb{Q}(\mathcal{S})$ .*

The proof is omitted since it is obvious. The proposition implies that  $\mathbb{Q}(\mathcal{S})$  can be solved by using its subset  $\bar{\mathcal{S}}$  instead of  $\mathcal{S}$ . The cutting plane strategy offers a way of finding such a  $\bar{\mathcal{S}}$ , in which a smaller problem  $\mathbb{Q}(\bar{\mathcal{S}})$  has the same optimal solution as  $\mathbb{Q}(\mathcal{S})$ . In this strategy, we first choose some points from  $\mathcal{S}$  and construct a set  $\mathcal{S}^1$  containing these points. Let  $\mathcal{S}^k$  be the set constructed in the  $k$ th iteration. In the  $(k+1)$ th iteration, we choose some points from  $\mathcal{S} \setminus \mathcal{S}^k$  and expand  $\mathcal{S}^k$  to  $\mathcal{S}^{k+1}$  by adding these points to  $\mathcal{S}^k$ . Besides expanding, we also shrink  $\mathcal{S}^k$  by discarding some points which can be regarded as useless for drawing the origin-centered MVEE. These expanding and shrinking phases play an important role in constructing a small set. Algorithm 3 describes a cutting plane strategy for solving  $\mathbb{Q}(\mathcal{S})$ .

---

#### Algorithm 3 Cutting Plane Strategy for Solving $\mathbb{Q}(\mathcal{S})$

---

**Input:**  $\mathcal{S} = \{\mathbf{p}_1, \dots, \mathbf{p}_m\}$ .

**Output:**  $\mathbf{L}^*$ .

- 1: Choose an initial set  $\mathcal{S}^1$  from  $\mathcal{S}$  and let  $k = 1$ .
  - 2: Solve  $\mathbb{Q}(\mathcal{S}^k)$  and find the optimal solution  $\mathbf{L}^k$ . If  $\delta_{\mathbf{L}^k}(\mathbf{p}) \leq 1$  holds for all  $\mathbf{p} \in \mathcal{S} \setminus \mathcal{S}^k$ , let  $\mathbf{L}^* = \mathbf{L}^k$ , and stop.
  - 3: Choose a subset  $\mathcal{F}$  of  $\mathcal{S}^k$  and a subset  $\mathcal{G}$  of  $\{\mathbf{p} \in \mathcal{S} \setminus \mathcal{S}^k : \delta_{\mathbf{L}^k}(\mathbf{p}) > 1\}$ . Update  $\mathcal{S}^k$  as  $\mathcal{S}^{k+1} = (\mathcal{S}^k \setminus \mathcal{F}) \cup \mathcal{G}$  and increase  $k$  by 1. Then, go back to Step 2.
- 

Now, we give a description of our implementation of Algorithm 3. To construct the initial set  $\mathcal{S}^1$  in Step 1, our implementation employs the algorithm used in the papers (Kumar and Yildirim, 2005; Todd and Yildirim, 2007; Ahipasaoglu et al., 2008). The algorithm constructs a set  $\mathcal{S}^1$  by greedily choosing  $2d$  points in a step-by-step manner such that the convex hull is a  $d$ -dimensional crosspolytope containing as many points in  $\mathcal{S}$  as possible. We refer the reader to Algorithm 3.1 of Kumar and Yildirim (2005) for the precise description.

To shrink and expand  $\mathcal{S}^k$  in Step 3, we use a shrinking threshold parameter  $\theta$  such that  $\theta < 1$ , and an expanding size parameter  $\eta$  such that  $\eta \geq 1$ . These parameters are set before running Algorithm 3. For shrinking, we construct  $\mathcal{F} = \{\mathbf{p} \in \mathcal{S}^k : \delta_{\mathbf{L}^k}(\mathbf{p}) \leq \theta\}$  by using  $\theta$ .

For expanding, we arrange the points of  $\{\mathbf{p} \in \mathcal{S} \setminus \mathcal{S}^k : \delta_{\mathcal{L}^k}(\mathbf{p}) > 1\}$  in descending order, as measured by  $\delta_{\mathcal{L}^k}(\cdot)$ , and construct  $\mathcal{G}$  by choosing the top  $(m - 2d)/\eta$  points. If the set  $\{\mathbf{p} \in \mathcal{S} \setminus \mathcal{S}^k : \delta_{\mathcal{L}^k}(\mathbf{p}) > 1\}$  has less than  $(m - 2d)/\eta$  points, we choose all the points and construct  $\mathcal{G}$ .

## 6. Experiments

We experimentally compared Algorithm 2 with SPA and the variants of XRAY. These two existing algorithms were chosen because their studies (Bittorf et al., 2012; Gillis and Luce, 2013; Kumar et al., 2013) report that they outperform AGKM and Hottopixx, and scale to the problem size. Two types of experiments were conducted: one is the evaluation for the robustness of the algorithms to noise on synthetic data sets, and the other is the application of the algorithms to clustering of real-world document corpora.

We implemented Algorithm 2, and three variants of XRAY, “max”, “dist” and “greedy”, in MATLAB. We put Algorithm 3 in Algorithm 2 so it would solve  $\mathbb{Q}$  efficiently. The software package SDPT3 (Toh et al., 1999) was used for solving  $\mathbb{Q}(\mathcal{S}^k)$  in Step 2 of Algorithm 3. The shrinking parameter  $\theta$  and expanding size parameter  $\eta$  were set as 0.9999 and 5, respectively. The implementation of XRAY formulated the computation of the residual matrix  $\mathbf{R} = \mathbf{A}(\mathcal{I}_k)\mathbf{X}^* - \mathbf{A}$  as a convex optimization problem,

$$\mathbf{X}^* = \arg \min_{\mathbf{X} \geq \mathbf{0}} \|\mathbf{A}(\mathcal{I}_k)\mathbf{X} - \mathbf{A}\|_F^2.$$

For the implementation of SPA (Gillis and Vavasis, 2014), we used code from the first author’s website. Note that SPA and XRAY are sensitive to the normalization of the column vectors of the data matrix (see, for instance, Kumar et al., 2013), and for this reason, we used a data matrix whose column vectors were not normalized. All experiments were done in MATLAB on a 3.2 GHz CPU processor and 12 GB memory.

We will use the following abbreviations to represent the variants of algorithms. For instance, Algorithm 2 with SPA for an index selection of Step 2 is referred to as ER-SPA. Also, the variant of XRAY with “max” selection policy is referred to as XRAY(max).

### 6.1 Synthetic Data

Experiments were conducted for the purpose of seeing how well Algorithm 2 could improve the robustness of SPA and XRAY to noise. Specifically, we compared it with SPA, XRAY(max), XRAY(dist), and XRAY(greedy). The robustness of algorithm was measured by a recovery rate. Let  $\mathcal{I}$  be an index set of basis vectors in a noisy separable matrix, and  $\mathcal{I}^*$  be an index set returned by an algorithm. The recovery rate is the ratio given by  $|\mathcal{I} \cap \mathcal{I}^*| / |\mathcal{I}|$ .

We used synthetic data sets of the form  $\mathbf{F}(\mathbf{I}, \mathbf{K})\mathbf{\Pi} + \mathbf{N}$  with  $d = 250$ ,  $m = 5,000$ , and  $r = 10$ . The matrices  $\mathbf{F}$ ,  $\mathbf{K}$ ,  $\mathbf{\Pi}$  and  $\mathbf{N}$  were synthetically generated as follows. The entries of  $\mathbf{W} \in \mathbb{R}_+^{d \times r}$  were drawn from a uniform distribution on the interval  $[0, 1]$ . The column vectors of  $\mathbf{K} \in \mathbb{R}_+^{r \times \ell}$  were from a Dirichlet distribution whose  $r$  parameters were uniformly from the interval  $[0, 1]$ . The permutation matrix  $\mathbf{\Pi}$  was randomly generated. The entries of the noise matrix  $\mathbf{N} \in \mathbb{R}^{d \times m}$  were from a normal distribution with mean 0 and standard deviation  $\delta$ . The parameter  $\delta$  determined the intensity of the noise, and it was chosen from

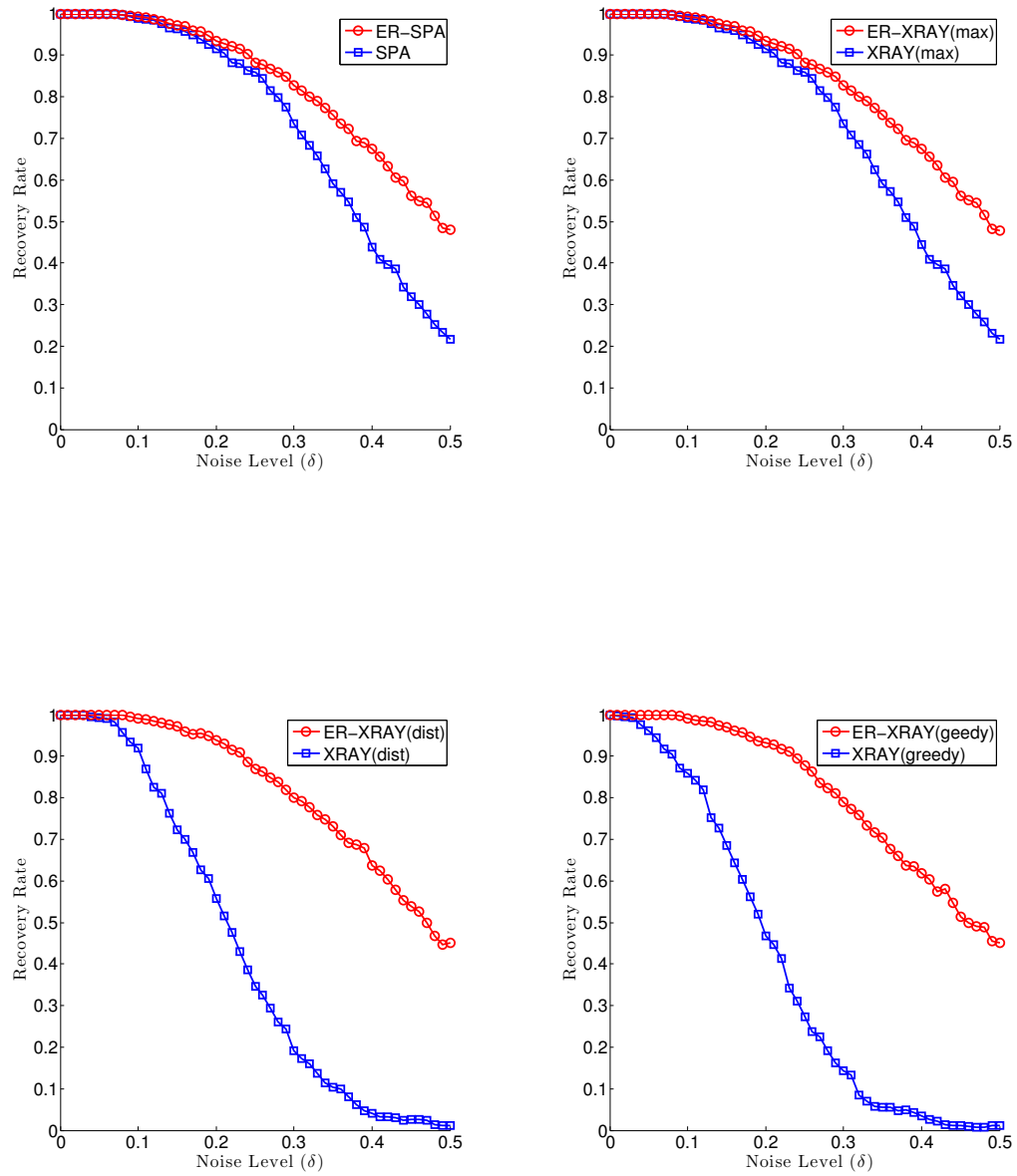


Figure 3: Comparison of the recovery rates of Algorithm 2 with SPA and XRAY.

Recovery rate	100%	90%	80%	70%
ER-SPA	0.06	0.24	0.32	0.37
SPA	0.05	0.21	0.27	0.31
ER-XRAY(max)	0.06	0.24	0.32	0.37
XRAY(max)	0.05	0.21	0.27	0.31
ER-XRAY(dist)	0.07	0.23	0.29	0.36
XRAY(dist)	0.03	0.10	0.13	0.16
ER-XRAY(greedy)	0.07	0.23	0.29	0.35
XRAY(greedy)	0.00	0.08	0.12	0.14

Table 1: Maximum values of noise level  $\delta$  for different recovery rates in percentage.

0 to 0.5 in 0.01 increments. A single data set consisted of 51 matrices with various amounts of noise, and we made 50 different data sets. Algorithm 2 was performed in the setting that  $\mathbf{M}$  is a matrix in the data set and  $r$  and  $\rho$  are each 10.

Figure 3 depicts the average recovery rate on the 50 data sets for Algorithm 2, SPA and XRAY. Table 1 summarizes the maximum values of noise level  $\delta$  for different recovery rates in percentage. The noise level was measured by 0.01, and hence, for instance, the entry “0.00” at XRAY(greedy) for 100% recovery rate means that the maximum value is in the interval  $[0.00, 0.01)$ . We see from the figure that Algorithm 2 improved the recovery rates of the existing algorithms. In particular, the recovery rates of XRAY(dist) and XRAY(greedy) rapidly decrease as the noise level increases, but Algorithm 2 significantly improved them. Also, the figure shows that Algorithm 2 tended to slow the decrease in the recovery rate. We see from the table that Algorithm 2 is more robust to noise than SPA and XRAY.

Table 2 summarizes the average number of active points and elapsed time for 50 data sets taken by Algorithm 2 with  $\delta = 0, 0.25$  and  $0.5$ . We read from the table that the elapsed time increases with the number of active points. The average elapsed times of SPA, XRAY(max), XRAY(dist), and XRAY(greedy) was respectively 0.03, 1.18, 16.80 and 15.85 in seconds. Therefore, we see that the elapsed time of Algorithm 2 was within a reasonable range.

$\delta$	Active points	Elapsed time (second)			
		ER-SPA	ER-XRAY(max)	ER-XRAY(dist)	ER-XRAY(greedy)
0	10	1.05	1.07	1.07	1.07
0.25	12	3.08	3.10	3.10	3.10
0.5	23	4.70	4.71	4.71	4.71

Table 2: Average number of active points and elapsed time of Algorithm 2.

## 6.2 Application to Document Clustering

Consider a set of  $d$  documents. Let  $m$  be the total number of words appearing in the document set. We represent the documents by a bag-of-words. That is, the  $i$ th document is represented as an  $m$ -dimensional vector  $\mathbf{a}_i$ , whose elements are the appearance frequencies of words in the document. A document vector  $\mathbf{a}_i$  can be assumed to be generated by a

convex combination of several topic vectors  $\mathbf{w}_1, \dots, \mathbf{w}_r$ . This type of generative model has been used in many papers (for instance, Xu et al., 2003; Shahnaz et al., 2006; Arora et al., 2012b, 2013; Ding et al., 2013; Kumar et al., 2013).

Let  $\mathbf{W}$  be an  $r$ -by- $m$  topic matrix such that  $\mathbf{w}_1^\top, \dots, \mathbf{w}_r^\top$  are stacked from top to bottom and are of the form  $(\mathbf{w}_1; \dots; \mathbf{w}_r)$ . The model allows us to write a document vector in the form  $\mathbf{a}_i^\top = \mathbf{f}_i^\top \mathbf{W}$  by using a coefficient vector  $\mathbf{f}_i \in \mathbb{R}^r$  such that  $\mathbf{e}^\top \mathbf{f}_i = 1$  and  $\mathbf{f}_i \geq \mathbf{0}$ . This means that we have  $\mathbf{A} = \mathbf{F}\mathbf{W}$  for a document-by-word matrix  $\mathbf{A} = (\mathbf{a}_1; \dots; \mathbf{a}_d) \in \mathbb{R}_+^{d \times m}$ , a coefficient matrix  $\mathbf{F} = (\mathbf{f}_1; \dots; \mathbf{f}_d) \in \mathbb{R}_+^{d \times r}$ , and a topic matrix  $\mathbf{W} = (\mathbf{w}_1; \dots; \mathbf{w}_r) \in \mathbb{R}_+^{r \times m}$ . In the same way as the papers (Arora et al., 2012b, 2013; Ding et al., 2013; Kumar et al., 2013), we assume that a document-by-word matrix  $\mathbf{A}$  is separable. This requires that  $\mathbf{W}$  is of  $(\mathbf{I}, \mathbf{K})\mathbf{\Pi}$ , and it means that each topic has an *anchor word*. An anchor word is a word that is contained in one topic but not contained in the other topics. If an anchor word is found, it suggests that the associated topic exists.

Algorithms for Problem 1 can be used for clustering documents and finding topics for the above generative model. The algorithms for a document-word matrix  $\mathbf{A}$  return an index set  $\mathcal{I}$ . Let  $\mathbf{F} = \mathbf{A}(\mathcal{I})$ . The row vector elements  $f_{i1}, \dots, f_{ir}$  of  $\mathbf{F}$  can be thought of as the contribution rate of topics  $\mathbf{w}_1, \dots, \mathbf{w}_r$  for generating a document  $\mathbf{a}_i$ . The highest value  $f_{ij^*}$  among the elements implies that the topic  $\mathbf{w}_{j^*}$  contributes the most to the generation of document  $\mathbf{a}_i$ . Hence, we assign document  $\mathbf{a}_i$  to a cluster having the topic  $\mathbf{w}_{j^*}$ . There is an alternative to using  $\mathbf{F}$  for measuring the contribution rates of the topics. Step 1 of Algorithm 1 produces a rank- $r$  approximation matrix  $\mathbf{A}^r$  to  $\mathbf{A}$  as a by-product. Let  $\mathbf{F}' = \mathbf{A}^r(\mathcal{I})$ , and use it as an alternative to  $\mathbf{F}$ . We say that clustering with  $\mathbf{F}$  is clustering with the original data matrix, and that clustering with  $\mathbf{F}'$  is clustering with a low-rank approximation data matrix.

Experiments were conducted in the purpose of investigating clustering performance of algorithms and also checking whether meaningful topics could be extracted. To investigate the clustering performance, we used only SPA since our experimental results implied that XRAY would underperform. We assigned the values of the document-word matrix on the basis of the tf-idf weighting scheme, for which we refer the reader to Manning et al. (2008), and normalized the row vectors to the unit 1-norm.

To evaluate the clustering performance, we measured the accuracy (AC) and normalized mutual information (NMI). These measures are often used for this purpose (see, for instance, Xu et al., 2003; Manning et al., 2008). Let  $\Omega_1, \dots, \Omega_r$  be the manually classified classes and  $\mathcal{C}_1, \dots, \mathcal{C}_r$  be the clusters constructed by an algorithm. Both  $\Omega_i$  and  $\mathcal{C}_j$  are the subsets of the document set  $\{\mathbf{a}_1, \dots, \mathbf{a}_m\}$  such that each subset does not share any documents and the union of all subsets coincides with the document set. AC is computed as follows. First, compute the correspondence between classes  $\Omega_1, \dots, \Omega_r$  and clusters  $\mathcal{C}_1, \dots, \mathcal{C}_r$  such that the total number of common documents  $|\Omega_i \cap \mathcal{C}_j|$  is maximized. This computation can be done by solving an assignment problem. After that, rearrange the classes and clusters in the obtained order and compute

$$\frac{1}{d} \sum_{k=1}^r |\Omega_k \cap \mathcal{C}_k|.$$

This value is the AC for the clusters constructed by an algorithm. NMI is computed as

$$\frac{I(\Omega, \mathcal{C})}{\frac{1}{2}(E(\Omega) + E(\mathcal{C}))}.$$

$I$  and  $E$  denote the mutual information and entropy for the class family  $\Omega$  and cluster family  $\mathcal{C}$  where  $\Omega = \{\Omega_1, \dots, \Omega_r\}$  and  $\mathcal{C} = \{C_1, \dots, C_r\}$ . We refer the reader to Section 16.3 of Manning et al. (2008) for the precise forms of  $I$  and  $E$ .

Two document corpora were used for the clustering-performance evaluation: Reuters-21578 and 20 Newsgroups. These corpora are publicly available from the UCI Knowledge Discovery in Databases Archive (<http://kdd.ics.uci.edu>). In particular, we used the data preprocessing of Deng Cai, in which multiple classes are discarded. The data sets are available from the website (<http://www.cad.zju.edu.cn/home/dengcai>). The Reuters-21578 corpus consists of 21,578 documents appearing in the Reuters newswire in 1987, and these documents are manually classified into 135 classes. The text corpus is reduced by the preprocessing to 8,293 documents in 65 classes. Furthermore, we cut off classes with less than 5 documents. The resulting corpus contains 8,258 documents with 18,931 words in 48 classes, and the sizes of the classes range from 5 to 3,713. The 20 Newsgroups corpus consists of 18,846 documents with 26,213 words appearing in 20 different newsgroups. The size of each class is about 1,000.

We randomly picked some classes from the corpora and evaluated the clustering performance 50 times. Algorithm 2 was performed in the setting that  $\mathbf{M}$  is a document-word matrix and  $r$  and  $\rho$  each are the number of classes. In clustering with a low-rank approximation data matrix, we used the rank- $r$  approximation matrix to a document-word matrix.

# Classes	AC				NMI			
	Original		Low-rank approx.		Original		Low-rank approx.	
	ER-SPA	SPA	ER-SPA	SPA	ER-SPA	SPA	ER-SPA	SPA
6	0.605	0.586	0.658	0.636	0.407	0.397	0.532	0.466
8	0.534	0.539	0.583	0.581	0.388	0.387	0.491	0.456
10	0.515	0.508	0.572	0.560	0.406	0.393	0.511	0.475
12	0.482	0.467	0.532	0.522	0.399	0.388	0.492	0.469

Table 3: (Reuters-21578) Average AC and NMI of ER-SPA and SPA with the original data matrix and low-rank approximation data matrix.

Tables 3 and 4 show the results for Reuters-21578 and 20 Newsgroups, respectively. They summarize the average ACs and NMIs of ER-SPA and SPA. The column with “# Classes” lists the number of classes we chose. The columns labeled “Original” and “Low-rank approx.” are respectively the averages of the corresponding clustering measurements with the original data matrix and low-rank approximation data matrix. The tables suggest that clustering with a low-rank approximation data matrix performed better than clustering with the original data matrix. We see from Table 3 that ER-SPA could achieve improvements in the AC and NMI of SPA on Reuters-21578 when the clustering was done with a



# Classes	AC				NMI			
	Original		Low-rank approx.		Original		Low-rank approx.	
	ER-SPA	SPA	ER-SPA	SPA	ER-SPA	SPA	ER-SPA	SPA
6	0.441	0.350	0.652	0.508	0.314	0.237	0.573	0.411
8	0.391	0.313	0.612	0.474	0.306	0.242	0.555	0.415
10	0.356	0.278	0.559	0.439	0.291	0.228	0.515	0.397
12	0.319	0.240	0.517	0.395	0.268	0.205	0.486	0.372

Table 4: (20 Newsgroups) Average AC and NMI of ER-SPA and SPA with the original data matrix and low-rank approximation data matrix.

low-rank approximation data matrix. Table 4 indicates that ER-SPA outperformed SPA in AC and NMI on 20 Newsgroups.

Finally, we compared the topics obtained by ER-SPA and SPA. We used the BBC corpus of Greene and Cunningham (2006), which is available from the website (<http://mlg.ucd.ie/datasets/bbc.html>). The documents in the corpus have been subjected by preprocessed such as stemming, stop-word removal, and low word frequency filtering. It consists of 2,225 documents with 9,636 words that appeared on the BBC news website in 2004-2005. The documents were news on 5 topics: “business”, “entertainment”, “politics”, “sport” and “tech”.

AC		NMI	
ER-SPA	SPA	ER-SPA	SPA
0.939	0.675	0.831	0.472

Table 5: AC and NMI of ER-SPA and SPA with low-rank approximation data matrix for BBC.

	Anchor word	1	2	3	4	5
ER-SPA	film	award	best	oscar	nomin	actor
SPA	film	award	best	oscar	nomin	star
ER-SPA	mobil	phone	user	softwar	microsoft	technolog
SPA	mobil	phone	user	microsoft	music	download
ER-SPA	bank	growth	economi	price	rate	oil
SPA	bank	growth	economi	price	rate	oil
ER-SPA	game	plai	player	win	england	club
SPA	fiat	sale	profit	euro	japan	firm
ER-SPA	elect	labour	parti	blair	tori	tax
SPA	blog	servic	peopl	site	firm	game

Table 6: Anchor words and top-5 frequent words in topics grouped by ER-SPA and SPA for BBC.

Table 5 shows the ACs and NMIs of ER-SPA and SPA on the low-rank approximation data matrix for the BBC corpus. The table indicates that the AC and NMI of ER-SPA are higher than those of SPA. Table 6 summarizes the words in the topics obtained by ER-SPA and SPA. The topics were computed by using a low-rank approximation data matrix. The table lists the anchor word and the 5 most frequent words in each topic from left to right. We computed the correspondence between topics obtained by ER-SPA and SPA and grouped the topics for each algorithm. Concretely, we measured the 2-norm of each topic vector and computed the correspondence by solving an assignment problem. We can see from the table that the topics obtained by these two algorithms are almost the same from the first to the third panel, and they seem to correspond to “entertainment”, “tech” and “business”. The topics in the fourth and fifth panels, however, are different. The topic in the fifth panel by ER-SPA seems to correspond to “politics”. In contrast, it is difficult to find the topic corresponding to “politics” in the panels by SPA. These show that ER-SPA could extract more recognizable topics than SPA.

**Remark 14** *Sparsity plays an important role in computing the SVD for a large document corpus. In general, a document-word matrix arising from a text corpus is quite sparse. Our implementation of Algorithm 2 used the MATLAB command `svds` that exploits the sparsity of a matrix in the SVD computation. The implementation could work on all data of 20 Newsgroups corpus, which formed a document-word matrix of size 18,846-by-26,213.*

## 7. Concluding Remarks

We presented Algorithm 1 for Problem 1 and formally showed that it has correctness and robustness properties. Numerical experiments on synthetic data sets demonstrated that Algorithm 2, which is the practical implementation of Algorithm 1, is robustness to noise. The robustness of the algorithm was measured in terms of the recovery rate. The results indicated that Algorithm 2 can improve the recovery rates of SPA and XRAY. The algorithm was then applied to document clustering. The experimental results implied that it outperformed SPA and extracted more recognizable topics.

We will conclude by suggesting a direction for future research. Algorithm 2 needs to do two computations: one is the SVD of the data matrix and the other is the MVEE for a set of reduced-dimensional data points. It would be ideal to have a single computation that could be parallelized. The MVEE computation requires that the convex hull of data points is full-dimensional. Hence, the SVD computation should be carried out on data points. However, if we could devise an alternative convex set for MVEE, it would be possible to avoid SVD computation. It would be interesting to investigate the possibility of algorithms that find near-basis vectors by using the other type of convex set for data points.

## Acknowledgments

The author would like to thank Akiko Takeda of the University of Tokyo for her insightful and enthusiastic discussions, and thank the referees for careful reading and helpful suggestions that considerably improved the quality of this paper.

## References

- S. D. Ahipasaoglu, P. Sun, and M. J. Todd. Linear convergence of a modified Frank-Wolfe algorithm for computing minimum-volume enclosing ellipsoids. *Optimization Methods and Software*, 23(1):5–19, 2008.
- S. Arora, R. Ge, R. Kannan, and A. Moitra. Computing a nonnegative matrix factorization – Provably. In *Proceedings of the 44th Symposium on Theory of Computing (STOC)*, pages 145–162, 2012a.
- S. Arora, R. Ge, and A. Moitra. Learning topic models – Going beyond SVD. In *Proceedings of the 2012 IEEE 53rd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 1–10, 2012b.
- S. Arora, R. Ge, Y. Halpern, D. Mimno, and A. Moitra. A practical algorithm for topic modeling with provable guarantees. In *Proceedings of the 30th International Conference on Machine Learning (ICML)*, 2013.
- V. Bittorf, B. Recht, C. Re, and J. A. Tropp. Factoring nonnegative matrices with linear programs. In *Advances in Neural Information Processing Systems 25 (NIPS)*, pages 1223–1231, 2012.
- A. Cichocki, R. Zdunek, A. H. Phan, and S. Amari. *Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation*. Wiley, 2009.
- W. Ding, M. H. Rohban, P. Ishwar, and V. Saligrama. Topic discovery through data dependent and random projections. In *Proceedings of the 30th International Conference on Machine Learning (ICML)*, 2013.
- D. Donoho and V. Stodden. When does non-negative matrix factorization give a correct decomposition into parts? In *Advances in Neural Information Processing Systems 16 (NIPS)*, pages 1141–1148, 2003.
- N. Gillis. Robustness analysis of Hottopixx, a linear programming model for factoring nonnegative matrices. *SIAM Journal on Matrix Analysis and Applications*, 34(3):1189–1212, 2013.
- N. Gillis and R. Luce. Robust near-separable nonnegative matrix factorization using linear optimization. arXiv:1302.4385v1, 2013.
- N. Gillis and S. A. Vavasis. Fast and robust recursive algorithms for separable nonnegative matrix factorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(4):698–714, 2014.
- G. H. Golub and C. F. Van Loan. *Matrix Computation*. The Johns Hopkins University Press, 3rd edition, 1996.
- P. Gong and C. Zhang. Efficient nonnegative matrix factorization via projected newton method. *Pattern Recognition*, 45(9):3557–3565, 2012.

- D. Greene and P. Cunningham. Practical solutions to the problem of diagonal dominance in kernel document clustering. In *Proceedings of the 23th International Conference on Machine Learning (ICML)*, 2006.
- L. G. Khachiyan. Rounding of polytopes in the real number model of computation. *Mathematics of Operations Research*, 21(2):307–320, 1996.
- H. Kim and H. Park. Non-negative matrix factorization based on alternating non-negativity constrained least squares and active set method. *SIAM Journal on Matrix Analysis and Applications*, 30(2):713–730, 2008.
- H. Kim and H. Park. Fast nonnegative matrix factorization: An active-set-like method and comparisons. *SIAM Journal on Scientific Computing*, 33(6):3261–3281, 2011.
- J. Kim, Y. He, and H. Park. Algorithms for nonnegative matrix and tensor factorizations: a unified view based on block coordinate descent framework. *Journal of Global Optimization*, 58(2):285–319, 2014.
- A. Kumar, V. Sindhwani, and P. Kambadur. Fast conical hull algorithms for near-separable non-negative matrix factorization. In *Proceedings of the 30th International Conference on Machine Learning (ICML)*, 2013.
- P. Kumar and E. A. Yildirim. Minimum-volume enclosing ellipsoids and core sets. *Journal of Optimization Theory and Applications*, 126(1), 2005.
- D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401:788–791, 1999.
- D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *Advances in Neural Information Processing Systems 13 (NIPS)*, pages 556–562, 2001.
- C.-J. Lin. Projected gradient methods for non-negative matrix factorization. *Neural Computation*, 19(10):2756–2779, 2007.
- C. D. Manning, P. Raghavan, and H. Schuetze. *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- L. Miao and H. Qi. Endmember extraction from highly mixed data using minimum volume constrained nonnegative matrix factorization. *IEEE Transactions on Geoscience and Remote Sensing*, 45(2):765–777, 2007.
- J. M. P. Nascimento and J. M. B. Dias. Vertex component analysis: A fast algorithm to unmix hyperspectral data. *IEEE Transactions on Geoscience and Remote Sensing*, 43(4):898–910, 2005.
- F. Shahnaz, M. W. Berry, V. P. Pauca, and R. J. Plemmons. Document clustering using nonnegative matrix factorization. *Information Processing and Management*, 42(2):373–386, 2006.
- P. Sun and R. M. Freund. Computation of minimum-volume covering ellipsoids. *Operations Research*, 52(5):690–706, 2004.

- M. J. Todd and E. A. Yildirim. On Khachiyan’s algorithm for the computation of minimum-volume enclosing ellipsoids. *Discrete Applied Mathematics*, 155(13):1731–1744, 2007.
- K.-C. Toh. Primal-dual path-following algorithms for determinant maximization problems with linear matrix inequalities. *Computational Optimization and Applications*, 14(3):309–330, 1999.
- K.-C. Toh, M. J. Todd, and R. H. Tütüncü. SDPT3 – a MATLAB software package for semidefinite programming. *Optimization Methods and Software*, 11:545–581, 1999.
- T. Tsuchiya and Y. Xia. An extension of the standard polynomial-time primal-dual path-following algorithm to the weighted determinant maximization problem with semidefinite constraints. *Pacific Journal of Optimization*, 3(1):165–182, 2007.
- L. Vandenberghe, S. Boyd, and S. P. Wu. Determinant maximization with linear matrix inequality constraints. *SIAM Journal on Matrix Analysis and Applications*, 19(2):499–533, 1998.
- S. A. Vavasis. On the complexity of nonnegative matrix factorization. *SIAM Journal of Optimization*, 20(3):1364–1377, 2009.
- W. Xu, X. Liu, and Y. Gong. Document clustering based on non-negative matrix factorization. In *Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 267–273, 2003.