

Fast Projection-Based Methods for the Least Squares Nonnegative Matrix Approximation Problem

Dongmin Kim, Suvrit Sra, Inderjit S. Dhillon

Department of Computer Sciences, University of Texas Austin, TX 78712-1188, USA

Received 16 July 2007; accepted 7 August 2007

DOI:10.1002/sam.104

Published online 28 December 2007 in Wiley InterScience (www.interscience.wiley.com).

Abstract: Nonnegative matrix approximation (NNMA) is a popular matrix decomposition technique that has proven to be useful across a diverse variety of fields with applications ranging from document analysis and image processing to bioinformatics and signal processing. Over the years, several algorithms for NNMA have been proposed, e.g. Lee and Seung's multiplicative updates, alternating least squares (ALS), and gradient descent-based procedures. However, most of these procedures suffer from either slow convergence, numerical instability, or at worst, serious theoretical drawbacks. In this paper, we develop a new and improved algorithmic framework for the least-squares NNMA problem, which is not only theoretically well-founded, but also overcomes many deficiencies of other methods. Our framework readily admits powerful optimization techniques and as concrete realizations we present implementations based on the Newton, BFGS and conjugate gradient methods. Our algorithms provide numerical results superior to both Lee and Seung's method as well as to the alternating least squares heuristic, which was reported to work well in some situations but has no theoretical guarantees [1]. Our approach extends naturally to include regularization and box-constraints without sacrificing convergence guarantees. We present experimental results on both synthetic and real-world datasets that demonstrate the superiority of our methods, both in terms of better approximations as well as computational efficiency. © 2007 Wiley Periodicals, Inc. *Statistical Analy Data Mining* 1: 38–51, 2008

Keywords: nonnegative matrix approximation; factorization; projected Newton methods; active sets; least-squares

1. INTRODUCTION

Nonnegative matrix approximation (NNMA), also known as *nonnegative matrix factorization* [2] or *positive matrix factorization* [3], is a popular and effective matrix decomposition technique. It has become an established method for performing dimensionality reduction and related tasks such as clustering, image processing, and visualization—with applications across a diverse variety of fields. The NNMA problem setting is defined as follows. Let $A = [\mathbf{a}_1, \dots, \mathbf{a}_N]$ be the matrix of nonnegative inputs, where each column $\mathbf{a}_i \in \mathbb{R}_+^M$. NNMA seeks to approximate these input vectors by nonnegative linear (conic) combinations of a small number of *nonnegative representative vectors*, $\{\mathbf{b}_1, \dots, \mathbf{b}_K\}$, so that

$$\mathbf{a}_i \approx \sum_{k=1}^K \mathbf{b}_k c_{ki}, \quad (1)$$

where the coefficients c_{ki} are also nonnegative. We remark in passing that various alternative restrictions on \mathbf{b}_k or $\mathbf{c}_i = [c_{1i} \ c_{2i} \ \dots \ c_{Ki}]^T$ may be placed to obtain different

types of approximations. For the purpose of this paper, we focus only on the problem with nonnegativity constraints.

The quality of the approximation in Eq. (1) may be measured using an appropriate distortion function, for example, the Frobenius norm distortion or the Kullback-Leibler divergence. In this paper we focus on the former distortion, which leads to the following *least-squares NNMA* problem:

$$\underset{\mathbf{B}, \mathbf{C} \geq 0}{\text{minimize}} \mathcal{F}(\mathbf{B}; \mathbf{C}) = \frac{1}{2} \|\mathbf{A} - \mathbf{BC}\|_{\mathbb{F}}^2, \quad (2)$$

where A is the input matrix and B, C are the output (factor) matrices. The matrix B may be intuitively viewed as a set of basis vectors that can be conically combined using the coefficients in C to approximate the input A .

In this paper we develop two new Newton-type algorithms for solving Eq. (2) along with a theoretical analysis establishing their convergence. Both our algorithms improve upon the *de facto* procedure of Lee and Seung [4], hereafter referred to as *LS*, as well as upon the popular ALS heuristic, which has been reported to perform well in practice [1]. However, both LS and ALS have their drawbacks;

Correspondence to: Suvrit Sra (suvrit@cs.utexas.edu)

the former suffers from slow convergence, whereas the latter lacks theoretical guarantees on its performance—our new algorithms rectify both of these deficiencies.

Researchers have also considered the following regularized NNMA problem

$$\underset{\mathbf{B}, \mathbf{C} \geq 0}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{A} - \mathbf{BC}\|_{\mathbb{F}}^2 + \lambda \|\mathbf{B}\|_{\mathbb{F}}^2 + \mu \|\mathbf{C}\|_{\mathbb{F}}^2, \quad (3)$$

where $\lambda > 0$, and $\mu > 0$ are regularization parameters. The motivation behind studying Eq. (3) can be ascribed to certain practical concerns. For example, the basic NNMA problem estimates the product \mathbf{BC} that has $(M + N)K$ parameters. Such a large number of parameters can lead to overfitting, which despite the apparent sparse representations yielded by NNMA, might be difficult to counter without regularization. Furthermore, the regularization terms also make the optimization problem more numerically stable.

Another interesting variation arises when one binds the solution values by imposing box-constraints on the variables. For NNMA this results in the problem

$$\begin{aligned} &\text{minimize} && \frac{1}{2} \|\mathbf{A} - \mathbf{BC}\|_{\mathbb{F}}^2, \\ &\text{subject to} && \mathbf{P} \leq \mathbf{B} \leq \mathbf{Q}, \\ &&& \mathbf{R} \leq \mathbf{C} \leq \mathbf{S}, \end{aligned} \quad (4)$$

where the inequalities are component-wise. Both Eqs. (3) and (4) can be handled by our methods without much additional difficulty.

2. BACKGROUND AND RELATED WORK

The NNMA objective function Eq. (2) is not simultaneously convex in both \mathbf{B} and \mathbf{C} due to the presence of the product term \mathbf{BC} . Hence, in general it is very difficult to find globally optimal solutions to Eq. (2). Fortunately, the objective function is at least individually convex in \mathbf{B} and in \mathbf{C} , which makes it possible to invoke an alternating minimization or descent procedure that takes the form:

1. Initialize \mathbf{B}^0 and/or \mathbf{C}^0 ; set $t \leftarrow 0$.
2. Fix \mathbf{B}^t and find \mathbf{C}^{t+1} such that

$$\mathcal{F}(\mathbf{B}^t, \mathbf{C}^{t+1}) \leq \mathcal{F}(\mathbf{B}^t, \mathbf{C}^t).$$

3. Fix \mathbf{C}^{t+1} and find \mathbf{B}^{t+1} such that

$$\mathcal{F}(\mathbf{B}^{t+1}, \mathbf{C}^{t+1}) \leq \mathcal{F}(\mathbf{B}^t, \mathbf{C}^{t+1}).$$

4. Let $t \leftarrow t + 1$, and repeat Steps 2 and 3 until convergence.

On the basis of the above procedure, we can categorize NNMA methods into two types, namely the *exact* and *inexact* methods. The former perform an exact minimization at each iterative step so that $\mathbf{C}^{t+1} = \text{argmin}_{\mathbf{C}} \mathcal{F}(\mathbf{B}^t, \mathbf{C})$ (similarly for \mathbf{B}^{t+1}), while the latter merely ensure descent, i.e. $\mathcal{F}(\mathbf{B}^t, \mathbf{C}^{t+1}) \leq \mathcal{F}(\mathbf{B}^t, \mathbf{C}^t)$ (similarly for \mathbf{B}^{t+1}).

Since the Frobenius norm of a matrix is just the sum of Euclidean norms over columns (or rows), minimization or descent over either \mathbf{B} or \mathbf{C} boils down to solving a sequence of nonnegative least squares (NNLS) problems of the form

$$\begin{aligned} &\underset{\mathbf{x}}{\text{minimize}} && f(\mathbf{x}) = \frac{1}{2} \|\mathbf{G}\mathbf{x} - \mathbf{h}\|_2^2, \\ &\text{subject to} && \mathbf{x} \geq \mathbf{0}. \end{aligned} \quad (5)$$

Exact methods find a global optimum of this subproblem, while inexact methods roughly approximate it. There do exist well-known methods for solving the NNLS problem, such as the Lawson-Hanson procedure [5], FNNLS [6], and other procedures mentioned in ref. [7]. However, as we show in ref. [8], our approach to solving NNLS outperforms the other methods, hence we favor it as the method of choice for solving Eq. (5). At this point, we alert the readers against a potential misinterpretation that could arise from our choice of nomenclature in terms of exact and inexact methods. It is not the case that the exact methods are superior to the inexact ones, or even that the exact methods could converge to a global optimum of Eq. (2). However, the exact methods do provide better theoretical properties and they tend to produce better quality solutions, even though there is still no guarantee on the global optimality due to the nonconvexity of Eq. (2). Inexact methods often provide great savings of computational effort by trading-off precision of the solutions for speed.

In this paper we present a new exact method for NNMA, which we call FNMA^E. There have been other exact approaches in the literature. For example, Paatero and Tapper [3, 9, 10] introduced a set of algorithms for NNMA and provided a convergence proof for *one* of their methods that employs the preconditioned conjugate gradient method. However, their methods are described in a nebulous fashion, and they cite the need for considerable engineering effort [9] for an actual implementation. Bierlaire *et al.* [11] developed a projected gradient method for NNLS, which Lin [12] applied to solve problem Eq. (2). Recently, Merritt and Zhang [13] developed an interior-point gradient method for NNLS—a gradient descent-based method without projection that maintains feasibility of intermediate solutions throughout the iterations. They also provided a convergence proof for their method under the mild assumption that \mathbf{G} has full-rank. Though problem Eq. (5) can be solved by any constrained optimization technique, the above methods are all based on gradient descent since it

allows for efficient handling of simple nonnegativity constraints. However, gradient-based methods are known to have linear convergence rate at best, and often suffer from a phenomenon known as *zigzagging or jamming*. FNMA^E subsumes the projected gradient-based method as a special case and retains its algorithmic simplicity while overcoming its deficiencies by employing a nondiagonal gradient scaling matrix.

The group of *inexact* methods has witnessed greater popularity and it includes Lee and Seung's [4] multiplicative algorithms. Gonzalez and Zhang [14] proposed a variant of Lee and Seung's method that utilizes a different scaling scheme for negative gradients to get faster convergence. Berry *et al.* [1] report the alternating least squares (ALS) procedure to be a simple but effective method for performing NNMA. The ALS procedure is somewhat *ad-hoc*—it solves the unconstrained least squares problem at each step exactly, followed by a truncation of the negative entries to zero. However, ALS does not have any convergence guarantees, and we discuss this in more detail in Section 2.1. Another inexact approach is provided by Zdunek and Cichocki's, which we refer to as the ZC method, Zdunek and Cichocki [15] who proposed the combination of projection with a quasi-Newton procedure for NNMA.

Our FNMA^E procedure is a quasi-Newton method that remedies the theoretical deficiencies of both the ALS and ZC methods. As an alternative, we also present an inexact method called FNMA^I that shares the same algorithmic framework as its exact counterpart FNMA^E while providing a computationally more efficient procedure.

2.1. ALS and ZC Methods

As alluded to the above, both the ALS and ZC methods have theoretical deficiencies, which can lead to

nonmonotonic changes in the objective function value and to inferior approximations. We illustrate these deficiencies more clearly in this section providing further motivation for our algorithms.

Both ALS and ZC are closely related to FNMA^E and FNMA^I. A critical difference between FNMA^E and both these approaches (ZC and ALS) is that the former is an *exact* approach, whereas the latter two are *inexact* methods. To see why these methods are inexact, consider the NNLS subproblem Eq. (5) that they must solve. Let us denote a projection onto the nonnegative orthant by $\mathcal{P}_+[\cdot]$. Assuming \mathbf{G} to be of full rank, the ALS update for subproblem Eq. (5) may be written as

$$\mathbf{x} = \mathcal{P}_+[(\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T \mathbf{h}], \quad (6)$$

or equivalently,

$$\mathbf{x} = \mathcal{P}_+[\mathbf{x} - (\mathbf{G}^T \mathbf{G})^{-1} (\mathbf{G}^T \mathbf{G} \mathbf{x} - \mathbf{G}^T \mathbf{h})].$$

For the ZC approach, the update is

$$\mathbf{x}^{\text{new}} = \mathcal{P}_+[\mathbf{x}^{\text{old}} - \alpha \mathbf{D} (\mathbf{G}^T \mathbf{G} \mathbf{x}^{\text{old}} - \mathbf{G}^T \mathbf{h})], \quad (7)$$

where $\alpha > 0$ and \mathbf{D} is some positive definite matrix that approximates $(\mathbf{G}^T \mathbf{G})^{-1}$, i.e. the inverse of the Hessian. Note that the ALS update has $\mathbf{D} = (\mathbf{G}^T \mathbf{G})^{-1}$ and $\alpha = 1$ in this form. Figure 1 illustrates why the updates Eqs. (6) and (7) are inexact, moreover they fail to decrease the objective function for an arbitrary positive α . Observe that Eq. (6) performs an exact-Newton step followed by projection, while Eq. (7) does a quasi-Newton with projection. Hence, we see that both the ALS and the ZC approaches can lead to an increase in the objective function value (also see Fig. 6). Our exact method, FNMA^E, fixes this problem and is provably convergent unlike the ALS and ZC methods.

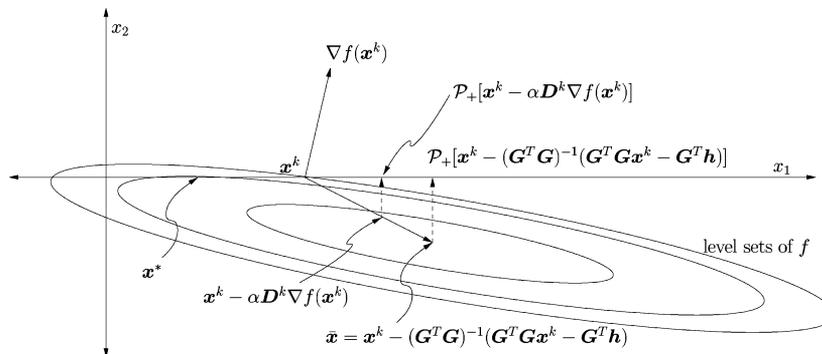


Fig. 1 Example where $\mathcal{P}_+[\mathbf{x}^k - \alpha \mathbf{D}^k \nabla f(\mathbf{x}^k)]$ fails to decrease the objective for an arbitrary $\alpha > 0$. In this figure, the ellipses represent level sets of f (the inner ellipses correspond to a smaller objective value), and \mathbf{D}^k is assumed to be exactly equal to the inverse of the Hessian. The current iterate is given as \mathbf{x}^k . Note that for Problem Eq. (5), the Newton method reaches the unconstrained optimum $\bar{\mathbf{x}}$ in a single iteration. However, the projected solution $\mathcal{P}_+[\mathbf{x}^k - \mathbf{D}^k \nabla f(\mathbf{x}^k)]$ for nonnegatively constrained problems leads to an increase in the objective since the current iterate (\mathbf{x}^k) moves from an inner ellipse to an outer one by the update rule.

3. ALGORITHMS AND THEORY

In this section, we develop an algorithm and associated supporting theory for solving Eq. (2). An efficient solution of the NNLS subproblem Eq. (5) forms the core of FNMA^E. Hence, first we focus our attention on efficiently solving the NNLS problem.

Broadly viewed, our method for solving NNLS may be viewed as combining the active set method with the projected gradient scheme. This approach is founded upon the observation that if the constraints active at the final solution are known in advance, the original problem can be solved by optimizing the objective in an equality-constrained manner over only the variables that correspond to the inactive constraints.

However, by itself, the projected gradient method, being a direct analog of steepest descent, suffers from deficiencies such as slow convergence and zigzagging. For *unconstrained* optimization problems, it is known that the use of nondiagonal positive definite gradient scaling matrices alleviates such problems. To overcome problems associated with gradient-based methods, Bertsekas [16] developed a projection framework for simply constrained cases based on the Newton-method. We build on that idea and employ nondiagonal gradient scaling based on the quasi-Newton method for Problem Eq. (5), which is a *constrained* minimization problem. However, since the constraints are particularly simple, this approach remains feasible and relatively simple.

3.1. Overview of our Method for NNLS

Our algorithm for solving Eq. (5) is iterative and at each iteration it partitions the variables into two groups, namely the *free* and *fixed* variables. The fixed variables are the components of \mathbf{x}^k with active constraints (equality satisfied) that have a corresponding positive derivative at iteration k . We index them by the *fixed set*, i.e.

$$I_+ = \{i | x_i^k = 0, [\nabla f(\mathbf{x}^k)]_i > 0\}. \quad (8)$$

For brevity, we will slightly manipulate the notation and say that $x_i^k \in I_+$ whenever $i \in I_+$.

We denote the free variables and the fixed variables at iteration k by \mathbf{y}^k and \mathbf{z}^k respectively. Without loss of generality, we can assume that \mathbf{x}^k and $\nabla f(\mathbf{x}^k)$ are partitioned as

$$\mathbf{x}^k = \begin{bmatrix} \mathbf{y}^k \\ \mathbf{z}^k \end{bmatrix}, \quad \nabla f(\mathbf{x}^k) = \begin{bmatrix} \nabla f(\mathbf{y}^k) \\ \nabla f(\mathbf{z}^k) \end{bmatrix},$$

where $y_i^k \notin I_+$ and $z_i^k \in I_+$. Once the free variables at the current iteration are identified, we compute the projection

\mathbf{y} (onto the nonnegative orthant) as follows

$$\mathbf{y} = \mathcal{P}_+[\mathbf{y}^k - \alpha \bar{\mathbf{D}}^k \nabla f(\mathbf{y}^k)], \quad (9)$$

where $\alpha \geq 0$, and $\bar{\mathbf{D}}^k$ is an appropriate positive definite gradient scaling matrix. Note that $\nabla f(\mathbf{y}^k)$ is the gradient vector restricted to the free variables, and $\bar{\mathbf{D}}^k$ is a corresponding restricted scaling matrix.

Finally, given \mathbf{y} we update \mathbf{x}^k to obtain

$$\mathbf{x}^{k+1} \leftarrow \begin{bmatrix} \mathbf{y} \\ \mathbf{z}^k \end{bmatrix} = \begin{bmatrix} \mathcal{P}_+[\mathbf{y}^k - \alpha \bar{\mathbf{D}}^k \nabla f(\mathbf{y}^k)] \\ \mathbf{0} \end{bmatrix}, \quad (10)$$

where the last equality uses the fact that \mathbf{z}^k is fixed to zero. Now we can compute $\nabla f(\mathbf{x}^{k+1})$ and update the fixed set I_+ to obtain \mathbf{y}^{k+1} and \mathbf{z}^{k+1} .

Note that any algorithm that finds \mathbf{y} such that

$$g^k(\mathbf{y}) < g^k(\mathbf{y}^k), \quad \mathbf{y} \geq 0, \quad (11)$$

where

$$g^k(\mathbf{y}) = \frac{1}{2} \|\mathbf{G}[\mathbf{y}; \mathbf{z}^k] - \mathbf{h}\|_2^2, \quad (12)$$

can be used to update \mathbf{x}^k in Eq. (10), but since Eq. (11) is again a constrained problem, Eq. (9) remains a good choice for feasibility and efficiency of the overall algorithm. Furthermore, due to the resemblance of Eq. (9) to an iteration of the standard quasi-Newton update, it is possible to exploit the curvature information of g^k to obtain a faster convergence rate.

However, the computation of a proper $\bar{\mathbf{D}}^k$ at each iteration is not a trivial task as the size of \mathbf{y}^k may vary across iterations, and it may be necessary to vary the size of $\bar{\mathbf{D}}^k$ from one iteration to the next. To address this difficulty, we note that the curvature information from $\{\mathbf{y}^k\}$ is essentially captured by the sequence $\{\mathbf{x}^k\}$. Therefore, $\bar{\mathbf{D}}^k$ can be approximated by taking a proper sub-matrix of \mathbf{D}^k , which contains curvature information from the vectors $\{\mathbf{x}^k\}$. On the basis of the above rationale, we maintain a gradient scaling matrix \mathbf{D}^k that covers the entire vector \mathbf{x}^k at each iteration and build the restricted matrices $\bar{\mathbf{D}}^k$ from \mathbf{D}^k according to the free variables \mathbf{y}^k .

There are many possible choices for \mathbf{D}^k , ranging from the identity matrix to the inverse of the Hessian. We provide three well-established schemes for selecting the gradient scaling matrix \mathbf{D}^k : (i) the inverse of the Hessian, (ii) the BFGS update and (iii) the memoryless BFGS update [16]. Briefly, the use of inverse Hessian, which leads to the Newton method for unconstrained problems,

is suitable when the Hessian is available and sparse.¹ The BFGS update incrementally approximates the inverse of the Hessian using only gradient information at each iteration, whereas it is recommended for problems where the computation involving the Hessian is expensive. The memoryless BFGS is for larger problems where the storage for \mathbf{D}^k itself is too expensive. For the least squares objective function, the memoryless BFGS update becomes equivalent to the conjugate gradient method.

3.1.1. The BFGS update

Suppose \mathbf{H}^k is the current approximation to the Hessian. The BFGS update adds a rank-two correction to \mathbf{H}^k to obtain

$$\mathbf{H}^{k+1} = \mathbf{H}^k - \frac{\mathbf{H}^k \mathbf{u} \mathbf{u}^T \mathbf{H}^k}{\mathbf{u}^T \mathbf{H}^k \mathbf{u}} + \frac{\mathbf{w} \mathbf{w}^T}{\mathbf{u}^T \mathbf{w}}, \quad (13)$$

where $\mathbf{w} = \nabla f(\mathbf{x}^{k+1}) - \nabla f(\mathbf{x}^k)$, and $\mathbf{u} = \mathbf{x}^{k+1} - \mathbf{x}^k$. Let \mathbf{D}^k denote the inverse of \mathbf{H}^k , then applying the Sherman-Morrison-Woodbury formula to Eq. (13) yields

$$\begin{aligned} \mathbf{D}^{k+1} &= \mathbf{D}^k - \frac{(\mathbf{D}^k \mathbf{w} \mathbf{u}^T + \mathbf{u} \mathbf{w}^T \mathbf{D}^k)}{\mathbf{u}^T \mathbf{w}} \\ &\quad + \left(1 + \frac{\mathbf{w}^T \mathbf{D}^k \mathbf{w}}{\mathbf{u}^T \mathbf{w}}\right) \frac{\mathbf{u} \mathbf{u}^T}{\mathbf{u}^T \mathbf{w}}. \end{aligned} \quad (14)$$

Since, $\nabla f(\mathbf{x}^k) = \mathbf{G}^T \mathbf{G} \mathbf{x}^k - \mathbf{G}^T \mathbf{h}$ for the NNLS problem, Eq. (14) can be rewritten as

$$\begin{aligned} \mathbf{D}^{k+1} &= \mathbf{D}^k - \frac{(\mathbf{D}^k \mathbf{G}^T \mathbf{G} \mathbf{u} \mathbf{u}^T + \mathbf{u} \mathbf{u}^T \mathbf{G}^T \mathbf{G} \mathbf{D}^k)}{\mathbf{u}^T \mathbf{G}^T \mathbf{G} \mathbf{u}} \\ &\quad + \left(1 + \frac{\mathbf{u}^T \mathbf{G}^T \mathbf{G} \mathbf{D}^k \mathbf{G}^T \mathbf{G} \mathbf{u}}{\mathbf{u}^T \mathbf{G}^T \mathbf{G} \mathbf{u}}\right) \frac{\mathbf{u} \mathbf{u}^T}{\mathbf{u}^T \mathbf{G}^T \mathbf{G} \mathbf{u}}. \end{aligned} \quad (15)$$

REMARK 1: Note that $\|\mathbf{G} \mathbf{u}\|_2^2$ appears as the denominator in the last two terms of Eq. (15). When \mathbf{G} is of full-rank, Eq. (15) is always well-defined, since

$$\|\mathbf{G} \mathbf{u}\|_2 = 0, \quad \text{iff } \mathbf{u} = \mathbf{0},$$

which in turn implies that the method has converged and the update Eq. (15) is not needed anymore. In general, even if \mathbf{G} is rank-deficient, we can avoid trouble by simply bypassing the update. The only requirement that we need to satisfy is that \mathbf{D}^{k+1} remains positive and definite, which can be easily satisfied by simply setting $\mathbf{D}^{k+1} = \mathbf{D}^k$. However, in practice Eq. (15) remains well-defined, and we do not usually encounter $\|\mathbf{G} \mathbf{u}\| = 0$.

¹ In the implementation, we use the Q-less QR-decomposition of the Hessian instead of computing the inverse Hessian.

3.1.2. Memoryless BFGS update

By simply resetting the matrix \mathbf{D}^k to an identity matrix at each iteration k , we obtain the memoryless BFGS update. Assume that at each iteration, we compute an optimal parameter α in Eq. (9). It can be shown that the equality

$$\nabla f(\mathbf{x}^{k+1})^T \mathbf{u} = 0$$

holds under this assumption. Let $\mathbf{d}^{k+1} = -\mathbf{D}^{k+1} \nabla f(\mathbf{x}^{k+1})$, then from Eq. (14),

$$\begin{aligned} \mathbf{d}^{k+1} &= -\mathbf{D}^{k+1} \nabla f(\mathbf{x}^{k+1}) \\ &= -\nabla f(\mathbf{x}^{k+1}) - \left(1 + \frac{\mathbf{w}^T \mathbf{w}}{\mathbf{u}^T \mathbf{w}}\right) \frac{\mathbf{u} \mathbf{u}^T}{\mathbf{u}^T \mathbf{w}} \nabla f(\mathbf{x}^{k+1}) \\ &\quad + \frac{(\mathbf{w} \mathbf{u}^T + \mathbf{u} \mathbf{w}^T)}{\mathbf{u}^T \mathbf{w}} \nabla f(\mathbf{x}^{k+1}) \\ &= -\nabla f(\mathbf{x}^{k+1}) + \frac{\mathbf{u} \mathbf{w}^T}{\mathbf{u}^T \mathbf{w}} \nabla f(\mathbf{x}^{k+1}) \\ &= -\nabla f(\mathbf{x}^{k+1}) + \left(\frac{\nabla f(\mathbf{x}^{k+1})^T \mathbf{w}}{\mathbf{u}^T \mathbf{w}}\right) \mathbf{u} \\ &= -\nabla f(\mathbf{x}^{k+1}) + \left(\frac{\nabla f(\mathbf{x}^{k+1})^T \mathbf{G}^T \mathbf{G} \mathbf{u}}{\mathbf{u}^T \mathbf{G}^T \mathbf{G} \mathbf{u}}\right) \mathbf{u}. \end{aligned}$$

3.1.3. Line-search

From Eq. (9) we see that in addition to the computation of \mathbf{D}^k , the update also involves a parameter $\alpha > 0$. Like many other iterative optimization procedures, standard line-search methods can be used to choose the step-size α . We omit a discussion of the same for brevity and refer the reader to ref. [8].

3.2. Convergence

In this section, we prove that our method for NNLS as described above is an exact method, i.e. it converges to the globally optimal solution of Eq. (5). The main result of this section is the following theorem.

THEOREM 1: (Convergence and Optimality). If \mathbf{G} is of full-rank and $\{\mathbf{x}^k\}$ is the sequence of points generated by Eq. (10), then every limit point of $\{\mathbf{x}^k\}$ is a stationary point of Problem Eq. (5), and hence optimal since Eq. (5) is strictly convex.

The proof of this theorem depends on several lemmas that we prove below. Our proof is structured as follows. First we show that the update Eq. (10) ensures a monotonic

descent in the objective function value (Lemma 1). Then we show that the resulting sequence of iterates $\{\mathbf{x}^k\}$ has a limit-point (Lemma 2). Finally, we show that any limit point of the sequence $\{\mathbf{x}^k\}$ is also a stationary or KKT point of Eq. (5), thereby concluding the proof.

LEMMA 1: (Descent). If \mathbf{x}^k is not a stationary point of Eq. (5), then there exists some constant $\bar{\alpha}$ such that

$$f(\mathbf{x}^{k+1}) < f(\mathbf{x}^k), \quad \forall \alpha \in (0, \bar{\alpha}],$$

where \mathbf{x}^{k+1} and $f(\mathbf{x})$ are given by Eq. (10) and Eq. (5) respectively.

PROOF: By the construction of I_+ , all components of \mathbf{y}^k satisfy:

$$\text{either } y_i^k \neq 0 \text{ or } [\nabla f(\mathbf{y}^k)]_i \leq 0.$$

Furthermore, since \mathbf{x}^k is not a stationary point, there exists at least one i such that

$$[\nabla f(\mathbf{y}^k)]_i \neq 0.$$

Thus letting $\mathbf{d} = -\bar{\mathbf{D}}^k \nabla f(\mathbf{y}^k)$, we see that

$$\nabla f(\mathbf{y}^k)^T \mathbf{d} < 0,$$

since $\bar{\mathbf{D}}^k$ is a principal submatrix of the positive definite matrix \mathbf{D}^k , and is therefore itself positive definite. This establishes the fact that \mathbf{d} is a feasible descent direction. Now let $\gamma(\alpha) = \mathbf{y}^k + \alpha \mathbf{d}$ denote a step in the direction given by \mathbf{d} and consider partitioning the *free* variables into two disjoint sets of indices such that

$$I_1 = \{i | y_i^k > 0 \text{ or } (y_i^k = 0 \text{ and } d_i \geq 0)\}, \text{ and } I_2 = \{i | y_i^k = 0 \text{ and } d_i < 0\}.$$

It is easy to see that there exists $\alpha_1 > 0$ such that $\forall i \in I_1$,

$$y_i^k + \alpha d_i \geq 0, \quad \forall \alpha \leq \alpha_1.$$

Let us define a new search direction $\bar{\mathbf{d}}$,

$$\bar{d}_i = \begin{cases} d_i, & i \in I_1, \\ 0, & \text{otherwise.} \end{cases}$$

Then we have,

$$\mathcal{P}_+[\gamma(\alpha)] = \mathbf{y}^k + \alpha \bar{\mathbf{d}}, \quad \forall \alpha \in (0, \alpha_1].$$

Since $[\nabla f(\mathbf{y}^k)]_i \leq 0$ and $d_i < 0$ for $i \in I_2$, we get $\sum_{i \in I_2} [\nabla f(\mathbf{y}^k)]_i \cdot d_i \geq q_0$. Now we can conclude that

$$\begin{aligned} \nabla f(\mathbf{y}^k)^T \bar{\mathbf{d}} &= \sum_{i \in I_1} [\nabla f(\mathbf{y}^k)]_i \cdot d_i \leq \sum_{i \in \{I_1 \cup I_2\}} [\nabla f(\mathbf{y}^k)]_i \\ &\quad \times d_i = \nabla f(\mathbf{y}^k)^T \mathbf{d} < 0. \end{aligned}$$

Hence, $\bar{\mathbf{d}}$ is also a feasible descent direction. Therefore, letting $\mathbf{y} = \mathcal{P}_+[\gamma(\alpha)]$, there exists $\bar{\alpha} \in (0, \alpha_1]$ such that

$$g^k(\mathbf{y}) < g^k(\mathbf{y}^k), \quad \forall \alpha \in (0, \bar{\alpha}]$$

where g^k is as in Eq. (12). From Eq. (10), since \mathbf{z}^k remains fixed in \mathbf{x}^{k+1} , we conclude that

$$f(\mathbf{x}^{k+1}) < f(\mathbf{x}^k), \quad \forall \alpha \in (0, \bar{\alpha}].$$

LEMMA 2: (Limit point). Let $\{\mathbf{x}^k\}$ be a sequence of points generated by Eq. (10). Then this sequence has a limit point.

PROOF: Assume that we start the iteration at \mathbf{x}^0 where $f(\mathbf{x}^0) = M$. By Lemma 1, $\{f(\mathbf{x}^k)\}$ is a monotonically decreasing sequence, whereby \mathbf{x}^0 is a maximizer of f over the M -level set of f . If a convex quadratic function f is bounded above, its M -level set is also bounded. Denote this M -level set by \mathbb{X} . Then we can choose $\mathbf{u} \in \mathbb{X}$ such that

$$\|\mathbf{u}\|_2 \geq \|\mathbf{x}\|_2, \quad \forall \mathbf{x} \in \mathbb{X}.$$

Then $\{\mathbf{x}^k\}$ is bounded as $0 \leq \|\mathbf{x}^k\|_2 \leq \|\mathbf{u}\|_2$ for all k , hence the sequence has a limit point. This concludes the proof of the lemma.

3.2.1. Gradient related condition

Let $\{\mathbf{x}^k\}$ be a sequence generated by Eq. (10). Then for any subsequence $\{\mathbf{x}^k\}_{k \in \mathcal{K}}$ that converges to a nonstationary point,

$$\limsup_{t \rightarrow \infty} \|\mathbf{x}^{k_{t+1}} - \mathbf{x}^{k_t}\| < \infty, \quad (16)$$

$$\limsup_{t \rightarrow \infty} \nabla f(\mathbf{x}^{k_t})^T (\mathbf{x}^{k_{t+1}} - \mathbf{x}^{k_t}) < 0. \quad (17)$$

This is known as the *gradient related condition* in optimization literature and plays a crucial role to prove the convergence of a number of methods. Inequality Eq. 16 follows from Lemma 2 and it can be shown that our method also satisfies condition Eq. (17) [8].

Finally we present a proof of our main Theorem 1.

PROOF: Assume $\{\mathbf{x}^k\}$ converges to a nonstationary point $\bar{\mathbf{x}}$. From Lemma 1, it can be shown that there exists some ϵ_k such that $\lim_{k \rightarrow \infty} \epsilon_k = 0$ and

$$f(\mathbf{x}^k) - f(\mathbf{x}^{k+1}) = -\nabla f(\mathbf{x}^k)^T (\mathbf{x}^{k+1} - \mathbf{x}^k) - \epsilon_k > 0,$$

Since f is continuous, $\lim_{k \rightarrow \infty} f(\mathbf{x}^k) = f(\bar{\mathbf{x}})$. Consequently,

$$\lim_{k \rightarrow \infty} f(\mathbf{x}^k) - f(\mathbf{x}^{k+1}) = 0.$$

In turn, it implies

$$\lim_{k \rightarrow \infty} \nabla f(\mathbf{x}^k)^T (\mathbf{x}^{k+1} - \mathbf{x}^k) = 0,$$

which contradicts Eq. (17).

3.3. FNMA^E: An Exact Method for NNMA

Now we extend the ideas from Section 3.1 to the matrix case. To that end, we need to redefine various quantities in terms of matrices. First, observe that the gradient matrices $\nabla_{\mathbf{C}} \mathcal{F}(\mathbf{B}; \mathbf{C})$ and $\nabla_{\mathbf{B}} \mathcal{F}(\mathbf{B}; \mathbf{C})$ are

$$\begin{aligned} \nabla_{\mathbf{C}} \mathcal{F}(\mathbf{B}; \mathbf{C}) &= \mathbf{B}^T \mathbf{B} \mathbf{C} - \mathbf{B}^T \mathbf{A}, \quad \text{and} \\ \nabla_{\mathbf{B}} \mathcal{F}(\mathbf{B}; \mathbf{C}) &= \mathbf{B} \mathbf{C} \mathbf{C}^T - \mathbf{A} \mathbf{C}^T. \end{aligned}$$

Then we redefine the *fixed set* accordingly. For example, the fixed set corresponding to \mathbf{B} is defined as:

$$I_+ = \{(i, j) | B_{ij} = 0, [\nabla_{\mathbf{B}} \mathcal{F}(\mathbf{B}; \mathbf{C})]_{ij} > 0\}.$$

Finally, we define the *zero-out operator* \mathcal{Z}_+ with respect to the fixed set I_+ so that

$$[\mathcal{Z}_+[X]]_{ij} = \begin{cases} X_{ij}, & (i, j) \notin I_+, \\ 0, & \text{otherwise.} \end{cases} \quad (18)$$

Algorithm 1 FNMA^E

Input: $\mathbf{A} \in \mathbb{R}_+^{M \times N}$, K s.t. $1 \leq K \leq \min\{M, N\}$

Output: $\mathbf{B} \in \mathbb{R}_+^{M \times K}$, $\mathbf{C} \in \mathbb{R}_+^{K \times N}$

1. Initialize $\mathbf{B}^0, \mathbf{C}^0, t = 0, \mathbf{S} = \mathbf{I}$.

repeat

2. $\mathbf{B} \leftarrow \mathbf{B}^t, \quad \mathbf{C}^{\text{old}} \leftarrow \mathbf{C}^t$.

repeat

3.1. Compute the gradient matrix $\nabla_{\mathbf{C}} \mathcal{F}(\mathbf{B}; \mathbf{C}^{\text{old}})$.

3.2. Compute fixed set I_+ for \mathbf{C}^{old} .

3.3. Compute the step length vector α using line-search.

3.4. Update \mathbf{C}^{old} as

$$\begin{aligned} \mathbf{U} &\leftarrow \mathcal{Z}_+[\nabla_{\mathbf{C}} \mathcal{F}(\mathbf{B}; \mathbf{C}^{\text{old}})]; \quad \mathbf{U} \leftarrow \mathcal{Z}_+[\mathbf{S}\mathbf{U}]; \\ \mathbf{C}^{\text{new}} &\leftarrow \mathcal{P}[\mathbf{C}^{\text{old}} - \mathbf{U} \cdot \text{diag}(\alpha)]. \end{aligned}$$

3.5. $\mathbf{C}^{\text{old}} \leftarrow \mathbf{C}^{\text{new}}$.

3.6. Update \mathbf{S} if necessary.

until \mathbf{C}^{old} converges

4. $\mathbf{C}^{t+1} \leftarrow \mathbf{C}^{\text{old}}$.

5. Repeat steps similar to Step 2–4 to obtain \mathbf{B}^{t+1} .

6. $t \leftarrow t + 1$.

until Stopping criteria are met

Now we have all the pieces to describe the overall algorithm for solving the NNMA problem Eq. (2). Algorithm 1 presents our proposed method which we name fast non-negative matrix approximation—exact, i.e. FNMA^E.

In Step 3.4 of Algorithm 1, the first $\mathcal{Z}_+[\cdot]$ eliminates the “fixed” gradient information from the search direction, the second $\mathcal{Z}_+[\cdot]$ ensures that the fixed set remains fixed, and the projection $\mathcal{P}_+[\cdot]$ maintains feasibility of the next iterate.

Note that we maintain only one gradient scaling matrix \mathbf{D} at each alternating step even though our algorithm for NNLS suggests that each column should have its own gradient scaling matrix. We justify this strategy as follows. In Problem Eq. (5), a series of BFGS updates for \mathbf{D} aim at estimating the inverse of the Hessian. However, the true Hessian for Problem Eq. (5) is $\mathbf{G}^T \mathbf{G}$ which is a constant matrix. Thus in Problem Eq. (2), a matrix-wise extension of NNLS, every column shares the same true Hessian, whereby each column can also share the approximation of the inverse Hessian, namely \mathbf{D} . As long as we retain the positive definiteness of the matrix \mathbf{D} , this shared \mathbf{D} provides an effective gradient scaling, and it does not impede convergence of the algorithm. Also note that since the Hessian is of size $K \times K$, its exact inverse can be used if K is not too large at a computational cost of $O(K^2(M + N) + K^3)$ operations. This strategy is included in our second algorithm, FNMA^I, in the next section.

THEOREM 2: (Convergence of FNMA^E). If \mathbf{B}^t and \mathbf{C}^t retain full-rank, then the sequence $\{\mathbf{B}^t, \mathbf{C}^t\}$ generated by Algorithm FNMA^E converges to a stationary point of Problem Eq. (2).

PROOF: Algorithm 1 essentially performs the following alternating minimization at each outer iteration

$$\begin{aligned} \mathbf{C}^{t+1} &\leftarrow \underset{\mathbf{C} \geq 0}{\text{argmin}} \|\mathbf{A} - \mathbf{B}^t \mathbf{C}\|_{\mathbb{F}}^2, \quad \text{and } \mathbf{B}^{t+1} \\ &\leftarrow \underset{\mathbf{B} \geq 0}{\text{argmin}} \|\mathbf{A} - \mathbf{B} \mathbf{C}^{t+1}\|_{\mathbb{F}}^2. \end{aligned}$$

Similar to the argument in Lemma 2, the domain of Problem Eq. (2) can be considered to be compact. Since $\{\mathcal{F}(\mathbf{B}^t; \mathbf{C}^t)\}$ is monotone decreasing and bounded below, it has a limit point, $(\overline{\mathbf{B}}, \overline{\mathbf{C}})$, i.e.

$$\lim_{t \rightarrow \infty} \mathcal{F}(\mathbf{B}^t; \mathbf{C}^t) = \mathcal{F}(\overline{\mathbf{B}}; \overline{\mathbf{C}}).$$

Since \mathcal{F} is continuous, we have

$$\lim_{t \rightarrow \infty} \mathbf{B}^t = \overline{\mathbf{B}}, \text{ and } \lim_{t \rightarrow \infty} \mathbf{C}^t = \overline{\mathbf{C}}.$$

Now we can invoke the proof of the two-block Gauss-Seidel method [17] to conclude our claim.

3.4. FNMA^I: An *Inexact* Method for NNMA

In this section we present an *inexact* version of our approach. This method has the same underlying framework as FNMA^E, but uses some heuristics to reduce computational effort at each iteration.

Algorithm 2 gives the pseudocode for FNMA^I and it differs from the exact method in three main aspects. First, it uses the inverse of the Hessian as the nondiagonal gradient scaling matrix \mathbf{D} . Whenever the rank K of the factor matrices \mathbf{B} and \mathbf{C} is small, using the inverse Hessian can be advantageous for problems where $\mathbf{O}(K^3)$ costs are acceptable. Second, the step-size α is made an input parameter, and FNMA^I guarantees monotonic descent on the objective function for a sufficiently small α . Third, FNMA^I accepts the number of iterations for each alternating step as an input parameter. This modification permits premature termination of each alternating step, which naturally translates into large computational savings by trading-off accuracy for speed.

THEOREM 3: (Monotonicity of FNMA^I). If \mathbf{B}^t and \mathbf{C}^t retain full-rank, then $FNMA^I$ decreases its objective function monotonically for sufficiently small α .

PROOF: It is enough to consider Steps 2–3 from FNMA^I (the argument for Steps 4–5 is similar). Since \mathbf{B} is assumed to be full-rank at every step, $(\mathbf{B}^T \mathbf{B})^{-1}$ is positive definite. For a sufficiently small α that satisfies

$$\alpha \leq \min\{\alpha_i, i = 1, \dots, K\},$$

Algorithm 2 FNMA^I

Input: $\mathbf{A} \in \mathbb{R}_+^{M \times N}$, $K, \tau \in \mathbb{N}$, $\alpha \in \mathbb{R}_+$.

Output: $\mathbf{B} \in \mathbb{R}_+^{M \times K}$, $\mathbf{C} \in \mathbb{R}_+^{K \times N}$

1. Initialize $\mathbf{B}^0, \mathbf{C}^0, t = 0$.

repeat

2. $\mathbf{B} \leftarrow \mathbf{B}^t, \mathbf{C}^{\text{old}} \leftarrow \mathbf{C}^t$.

for $i = 1$ to τ **do**

3.1. Compute the gradient matrix $\nabla_{\mathbf{C}} \mathcal{F}(\mathbf{B}; \mathbf{C}^{\text{old}})$.

3.2. Compute fixed set I_+ for \mathbf{C}^{old} .

3.3. Update \mathbf{C}^{old} as:

$$\mathbf{U} \leftarrow \mathcal{Z}_+[\nabla_{\mathbf{C}} \mathcal{F}(\mathbf{B}; \mathbf{C}^{\text{old}})]; \quad \mathbf{U} \leftarrow \mathcal{Z}_+[(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{U}];$$

$$\mathbf{C}^{\text{new}} \leftarrow \mathcal{P}[\mathbf{C}^{\text{old}} - \alpha \mathbf{U}].$$

3.4. $\mathbf{C}^{\text{old}} \leftarrow \mathbf{C}^{\text{new}}$.

end for

4. $\mathbf{C}^{t+1} \leftarrow \mathbf{C}^{\text{old}}$.

5. Repeat steps similar to Step 2–4 to obtain \mathbf{B}^{t+1} .

6. $t \leftarrow t + 1$.

until Stopping criteria are met

where the α_i are computed by Step 3.3 from FNMA^E, it can be shown that Steps 2-3 decrease the objective monotonically by arguments similar to the ones in the proof of Lemma 1.

REMARK 2: If any α_i is zero, then the inner loop for the current alternating step should be terminated to guarantee monotonicity.

REMARK 3: A sufficiently small α is important to guarantee monotonicity of FNMA^I, but too small a value will hurt the computational benefit by slowing down convergence. On the other hand, if α is too large, it can push the search direction out of the feasible region, or introduce too many zeros into the current iterate, resulting in a singular or ill-conditioned Hessian for the next iterate.

To overcome these subtleties and to find a proper α in practice, the following simple heuristic can be used. Writing Step 3.3 from FNMA^I in the form:

$$\mathbf{W} \leftarrow \mathcal{P}_+[\mathbf{W} - \alpha \mathbf{U}],$$

1. Let number of *inner* iterations be small ($\tau = 2$ or 3),
2. Start with a large scaling λ (typically 0.1) and compute

$$\alpha = \lambda \frac{\|\mathbf{W}\|_F}{\|\mathbf{U}\|_F},$$

for each alternating step,

3. Decrease λ until it passes the inner steps without error,
4. Increase the number of iterations (typically $\tau = 10$).

3.5. Extensions to Handle Regularization

The regularized NNMA problem Eq. (3) can be solved by suitably modifying the FNMA^E and FNMA^I procedures. Essentially the gradient and Hessian get redefined. For example, the gradient

$$\nabla_C \mathcal{F}(\mathbf{B}; \mathbf{C}) = (\mathbf{B}^T \mathbf{B} + \lambda \mathbf{I})\mathbf{C} - \mathbf{B}^T \mathbf{A},$$

and the Hessian

$$\nabla_C^2 \mathcal{F}(\mathbf{B}; \mathbf{C}) = (\mathbf{B}^T \mathbf{B} + \lambda \mathbf{I}),$$

are suitably modified to include the contribution of the regularization term. We just use these updated values in the algorithms FNMA^E and FNMA^I to handle regularization. Notice that regularization provides the benefit of ensuring that the Hessian remains positive-definite. All the convergence results carry over without any additional work.

3.6. Handling Box-Constraints

FNMA^E and FNMA^I can be easily extended to handle box-constraints, i.e. constraints of the form $\mathbf{p} \leq \mathbf{x} \leq \mathbf{q}$. We motivate the details by first looking at the box-constrained version of Eq. (5), which is also known as *bounded least squares* (BLS) (Björck, 1996),

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \frac{1}{2} \|\mathbf{G}\mathbf{x} - \mathbf{h}\|^2, \\ & \text{subject to} && \mathbf{p} \leq \mathbf{x} \leq \mathbf{q}. \end{aligned} \tag{19}$$

Problem Eq. (19) can be solved just as we solved Eq. (2.1). We need to modify the definition of the *fixed-set* Eq. (8) so that

$$I_+ = \{i | (x_i^k = p_i, [\nabla f(\mathbf{x}^k)]_i > 0) \text{ or } (x_i^k = q_i, [\nabla f(\mathbf{x}^k)]_i < 0)\},$$

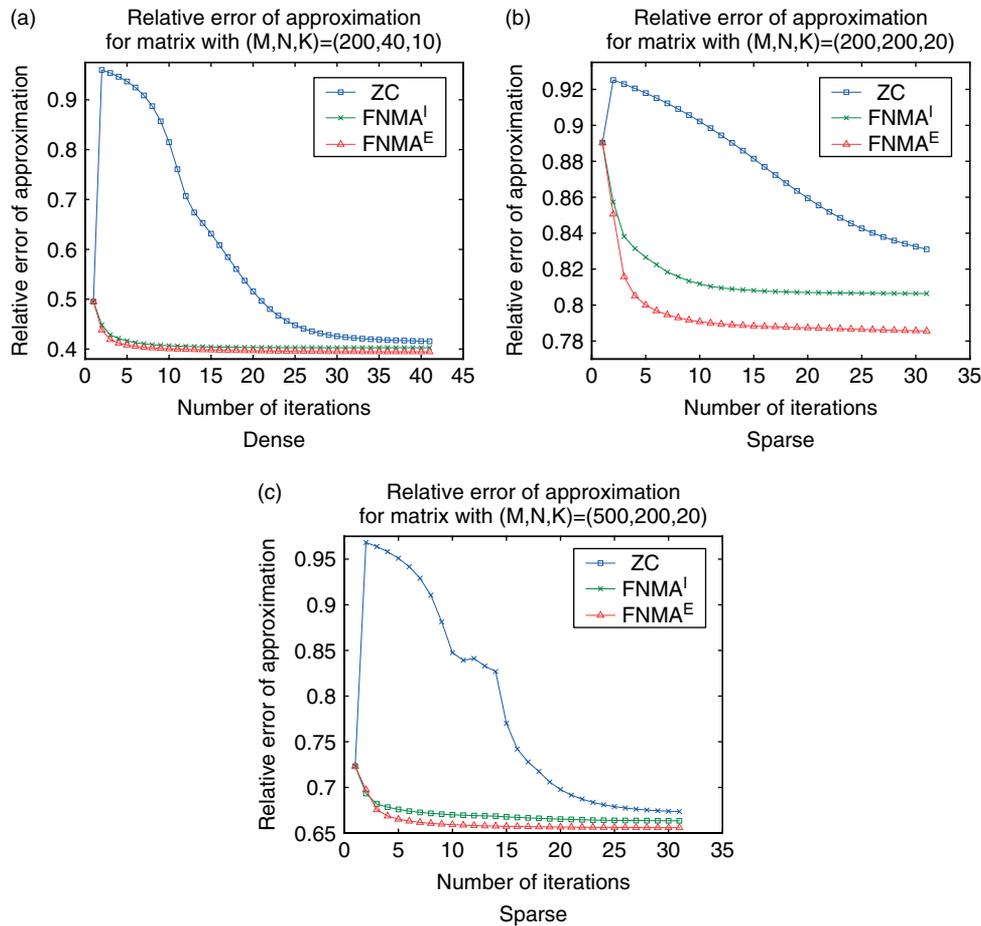


Fig. 2 Relative error of approximation against iteration count for ZC, FNMA^I, and FNMA^E. The relative errors achieved by both FNMA^I and FNMA^E are lower than ZC. Note that ZC does not decrease the errors monotonically.

and to replace the $\mathcal{P}_+[\cdot]$ projection by $\mathcal{P}_\Omega[\cdot]$, where

$$[\mathcal{P}_\Omega[\mathbf{x}]]_i = \begin{cases} p_i & : x_i \leq p_i \\ x_i & : p_i < x_i < q_i \\ q_i & : q_i \leq x_i \end{cases} \quad (20)$$

Given these definitions, it can be verified that Lemma 1 holds without significant modification and Theorem 1 also follows. The fact that the domain of Eq. (19) is a compact set obviates the need for Lemma 2 in this case.

Given the above method for BLS we can appropriately modify FNMA^E and FNMA^I for solving the bounded matrix approximation (BMA) problem Eq. (4). We omit the details for brevity, noting that the modifications needed are minor, for example, the fixed set for \mathbf{B} is redefined as

$$I_\Omega = \{(i, j) | (B_{ij} = P_{ij}, [\nabla_{\mathbf{B}} \mathcal{F}(\mathbf{B}; \mathbf{C})]_{ij} > 0), \text{ or} \\ (B_{ij} = Q_{ij}, [\nabla_{\mathbf{B}} \mathcal{F}(\mathbf{B}; \mathbf{C})]_{ij} < 0)\}.$$

By taking a projection step similar to Eq. (20) we can construct the desired method.

4. EXPERIMENTS

We now present experimental results to demonstrate the performance of our FNMA^E and FNMA^I methods. We give numerical results to assess the performance of our methods as compared to the standard LS method [4], ZC method [15], and the ALS approach [1] for solving the least-squares NNMA problem. Our experiments show that FNMA^E and FNMA^I produce better quality approximations than the LS, ZC, and the ALS procedures. We implemented LS, ALS, FNMA^E, and FNMA^I in MATLAB, while the ZC method was available in the NMFLAB toolbox [18]. We present results for the ZC method only with small matrices as the implementation available in NMFLAB was unable to run on larger matrices.

Since NNMA enjoys a vast number of applications [19], all of them stand to benefit from our new methods, especially because our methods achieve better objective function values and come with theoretical guarantees. As an illustration, we include some simple results on image processing in Section 4.2.

4.1. Error of Approximation

For our experiments, we initialize all the methods randomly or with one step of LS. Our results below show plots of the relative error of approximation, i.e. $\|\mathbf{A} - \mathbf{BC}\|_F / \|\mathbf{A}\|_F$ against the number of iterations. However, a word of caution is in order—iterations of these different

methods are not strictly comparable to each other, since some methods do more work than others in one iteration. A more interesting plot would have been “time” on the X-axis; however, at present we are unable to conduct such experiments since different implementations of each of the methods can change the running time substantially, for example, implementations that use BLAS3 versus those that do not. To perform timing comparisons, we intend to compare C/C++ implementations of these methods in the future.

4.1.1. Comparisons against ZC

The first experiment compares FNMA^E and FNMA^I against ZC on three data matrices. The results are reported in Fig. 2. As previously noted, the data matrices used are fairly small since ZC (NMFLAB) seems to be unable to cope with larger matrices. We initialized all methods using one iteration of LS, which itself was initialized randomly. However, in the figures we do not report the relative error for the random initialization as it is too large to display properly. Figure 2(a) indicates that our methods outperform ZC. The differences between the three algorithms are sharper in Fig. 2(b). Also note that ZC actually increases the approximation errors after the first iteration.

4.1.2. Comparisons against LS and ALS

On a larger matrix, Fig. 3 shows a comparison of the approximation errors for LS, ALS, FNMA^E, and FNMA^I.

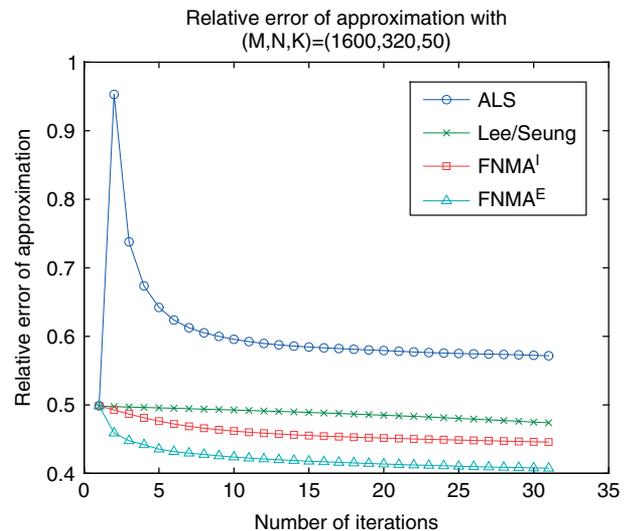


Fig. 3 Relative error values against iteration count for a random dense matrix of size 1600×320 for a rank 50 approximation. All methods other than ALS show a monotonic decrease when initialized with one step of LS.

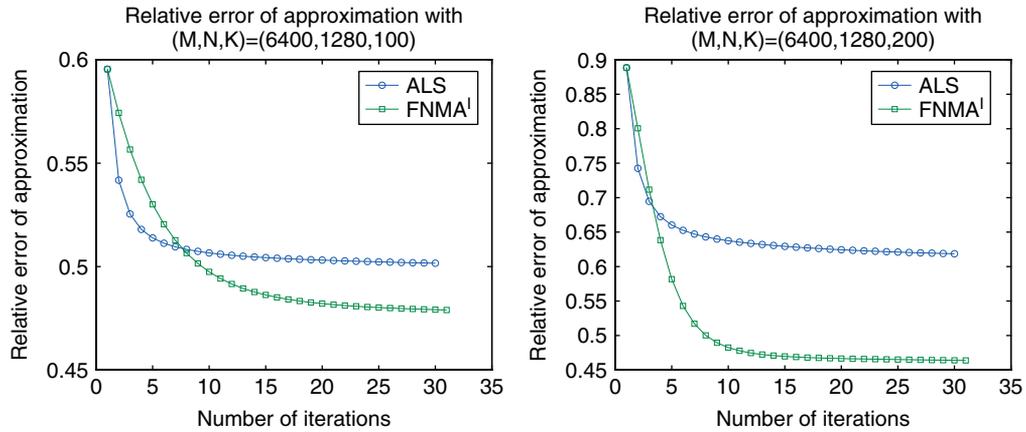


Fig. 4 Relative error values against iteration count for a random dense matrix of size 6400×1280 for a rank 100 (top) and rank 200 approximation (bottom).

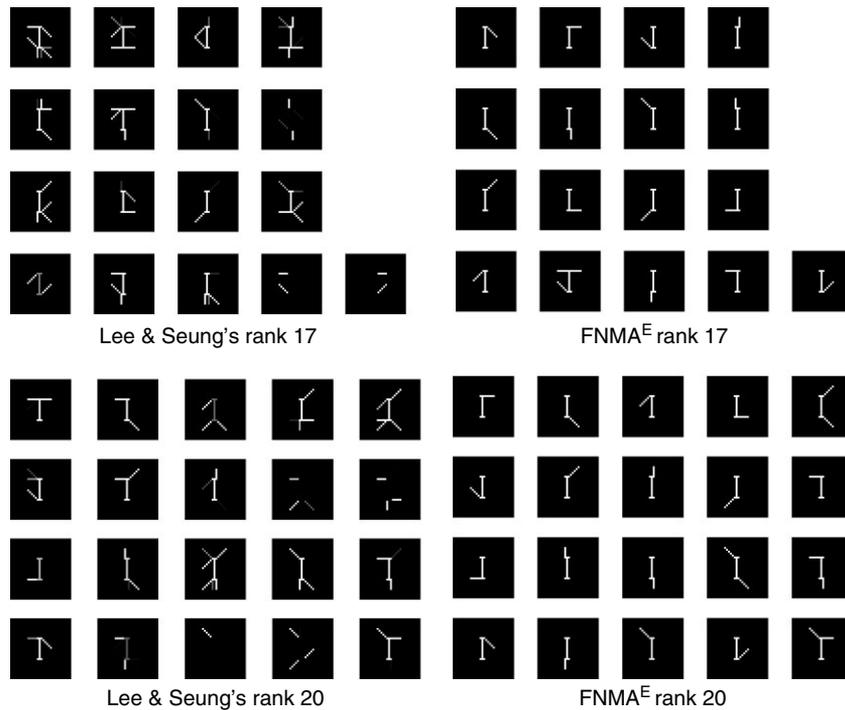


Fig. 5 Recovered nonnegative factors (matrix \mathbf{B}) from the swimmer dataset. The panels on the left are generated by Lee and Seung's algorithm while ones on the right by FNMA^E. From top to bottom, the approximation is performed for rank 17 and 20, respectively.

We see that FNMA^E achieves the best objective function values of all the methods presented. However, FNMA^E can take more running time than the other methods because of its *exact* nature. Therefore, FNMA^E is to be preferred when reconstruction accuracy is more important, while FNMA^I is recommended when running time is more important. We now present two more experiments to highlight the advantages of FNMA^I over ALS, which owing to its *ad-hoc* nature leads to inferior accuracies (see also Section 2.1).

Figure 4 compares the relative errors of approximation achieved by ALS and FNMA^I for a dense random matrix of size 6400×1280 . We emphasize again that the number of iterations is merely used as an indicator of progress of the algorithms, and is not to be taken as an indicator of time. From these figures one sees the interesting trend that as the rank of approximation increases, ALS becomes less and less competitive in terms of the objective function value achieved. For a rank-200 approximation (see Fig. 4),

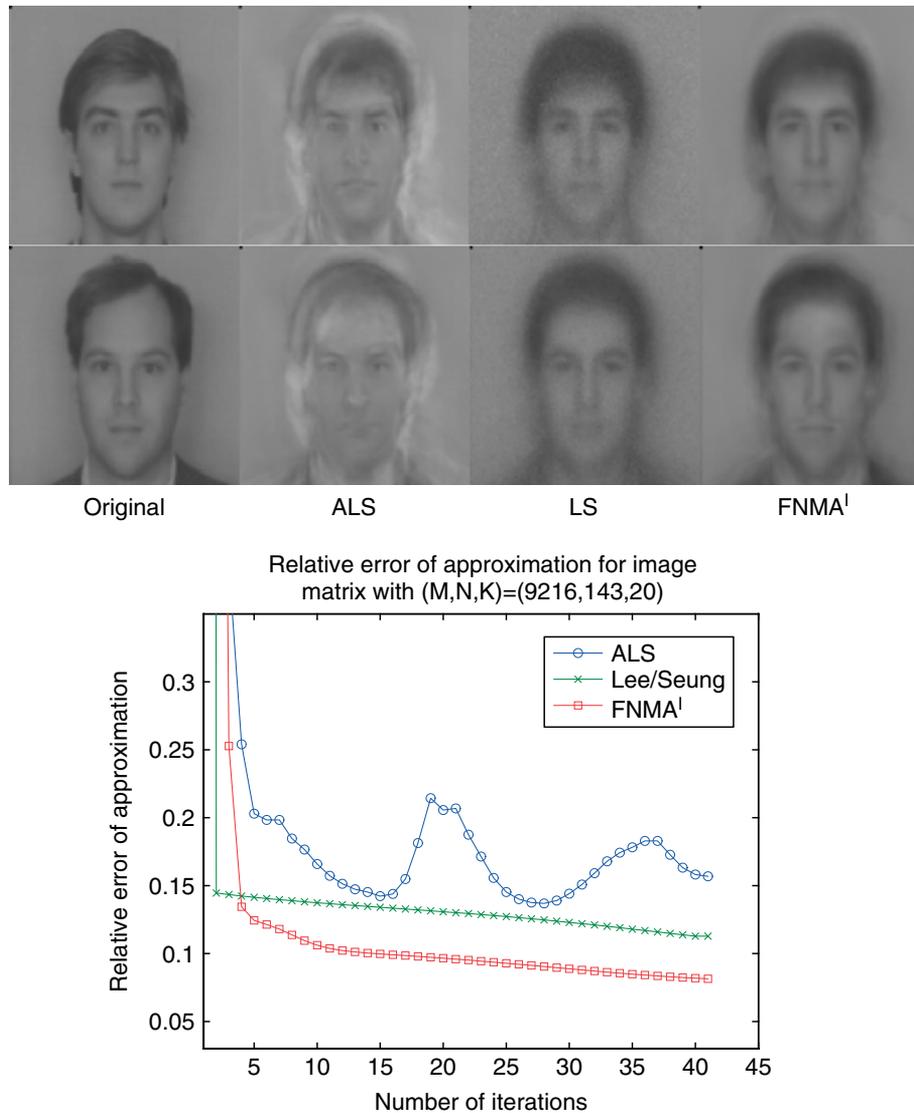


Fig. 6 Image reconstruction as obtained by the ALS, LS, and FNMA^I procedures. The figure illustrates two randomly chosen images out of the 143 reconstructed images, each with 96×96 pixels. The reconstruction was computed from a rank-20 approximation to the input image matrix, which was of size 9216×143 . The first image in each row is the original, followed by reconstructions obtained via ALS, LS, and FNMA^I. From the images above, FNMA^I is seen to obtain the best reconstruction and the relative errors as plotted on the right attest to this observation. Observe how ALS leads to a nonmonotonic change in the objective function value (as explained in Section 2.1). Note that the error values for the initialization are not fully shown as they are too large to display properly without obscuring the rest of the plot.

the accuracy achieved by FNMA^I is 25% higher than that achieved by ALS.

4.2. Application to Image Processing

NNMA was originally motivated by Lee and Seung [2] using an image processing application. Many other authors have also considered NNMA for image processing, graphics, or face recognition applications. Figure 5 compares LS method to FNMA^E for the swimmer dataset [20]. In the

figure, the effect of differences in final objective function values is more apparent. In our experiment, we run LS method up to 3000 iterations and FNMA^E up to 20 iterations. We also set the minimal threshold in the objective function value to be 10^{-5} ; hence both methods either stop after the maximum number of iterations or if progress is below the threshold. From Table 1, we can see that FNMA^E generally outperforms LS method, both in terms of the elapsed CPU time and the final objective function value achieved.

Table 1. Results on the swimmer dataset.

	Lee and Seung's	FNMA ^E	
Rank 17	182.24	62.29	Elapsed CPU time
	2.41×10^7	6.85×10^{-4}	Objective function value
Rank 20	156.18	41.93	Elapsed CPU time
	5.61×10^5	4.71×10^3	Objective function value

FNMA^E also produces sparser images (factors)—each image contains less number of component (limb) images. The difference becomes even more pronounced as the approximated rank approaches 17.²

Since, the quality of the reconstruction achieved by NNMA is important to many image processing applications, we provide a comparison of the various NNMA methods in terms of reconstruction accuracy—sample results are reported in Fig. 6, which shows accuracies for a rank-20 approximation to a 9216×143 matrix of face images.³ All methods were initialized with the same random \mathbf{B} and \mathbf{C} values.

This image dataset is an example of a real-world dense matrix for which ALS fails to decrease the objective function monotonically, resulting in a corresponding poorer reconstruction accuracy. FNMA^I achieves the best objective values of all three algorithms compared, and a corresponding better reconstruction is observed (see Fig. 6).

5. CONCLUSIONS

In this paper, we have presented new and improved Newton-type methods for the least-squares NNMA problem. By employing a nondiagonal gradient scaling scheme, our algorithms use curvature information and thus overcome deficiencies of gradient descent-based methods. Our methods also rectify serious draw-backs in existing methods such as ALS and ZC quasi-Newton heuristic. We provide

² In ref. [20], Donoho presents the result of a rank 16 estimation. Since we are not aware of an efficient algorithm to compute the nonnegative rank of matrix, we are not able to confirm the true nonnegative rank of the dataset. However, it obviously lies between the matrix rank 13 and the number of component images 17. We conjecture that the nonnegative rank is either 16 or 17 and our experiment confirms this conjecture by achieving almost zero error for a rank 17 approximation.

³ We preprocessed a publicly available face image database to create a subset of 143 grey-scale images of dimension 96×96 for our experiments.

convergence guarantees for our algorithms and verify their performance on real-life data from applications.

We provide two implementations based on the same algorithmic framework. Our *exact* method FNMA^E, which shows good performance in terms of approximation accuracy, is suitable for applications that require superior accuracy. Our *inexact* implementation FNMA^I is more suitable for applications that are more constrained by computational efficiency rather than accuracy.

ACKNOWLEDGEMENTS

This research was supported by NSF grant CCF-0431257, NSF Career Award ACI-0093404, and NSF-ITR award IIS-0325116.

REFERENCES

- [1] M. Berry, M. Browne, A. Langville, P. Pauca, and R. J. Plemmons, Algorithms and applications for approximation nonnegative matrix factorization, *Comput Stat Data Anal* (2006), Preprint.
- [2] D. D. Lee and H. S. Seung, Learning the parts of objects by nonnegative matrix factorization, *Nature* 401 (1999), 788–791.
- [3] P. Paatero and U. Tapper, Positive matrix factorization: A nonnegative factor model with optimal utilization of error estimates of data values, *Environmetrics* 5 (1994), 111–126.
- [4] D. D. Lee and H. S. Seung, Algorithms for Nonnegative Matrix Factorization. In *Neural Information Processing Systems*, 2000, 556–562.
- [5] C. L. Lawson and R. J. Hanson, *Solving Least Squares Problems*, Prentice-Hall, 1974.
- [6] R. Bro and S. D. Jong, A Fast Non-negativity-constrained Least Squares Algorithm, *J Chemomet* 11(5) (1997), 393–401.
- [7] Åke Björck, *Numerical Methods for Least Squares Problems*, SIAM, 1996.
- [8] D. Kim, S. Sra, and I. S. Dhillon, A New Projected Quasi-Newton Approach for the Non-negative Least Squares Problem. Technical Report TR-06-54, Computer Sciences, The University of Texas at Austin, 2006.
- [9] P. Paatero, Least-squares formulation of robust nonnegative factor analysis, *Chemomet Intell Lab Syst* 37 (1997), 23–35.
- [10] P. Paatero, The multilinear engine—a table-driven least squares program for solving multilinear problems, including the N-way parallel factor analysis model, *J Comput Graphical Statist* 8(4) (1999), 854–888.
- [11] M. Bierlaire, P. L. Toint, and D. Tuytens, On iterative algorithms for linear least squares problems with bound constraints, *Linear Algebra Appl* 143 (1991), 111–143.
- [12] C. Lin, Projected Gradient Methods for Non-negative Matrix Factorization. Technical Report ISSTECH-95-013, National Taiwan University, 2005.
- [13] M. Merritt and Y. Zhang, Interior-point gradient method for large-scale totally nonnegative least squares problems, *J Optim Theory Appl* 126(1) (2005), 191–202.

- [14] E. F. Gonzalez and Y. Zhang, Accelerating The Lee-Seung Algorithm for Nonnegative Matrix Factorization. Technical Report TR-05-02, Rice University, 2005.
- [15] R. Zdunek and A. Cichocki, Non-Negative Matrix Factorization with Quasi-Newton Optimization. In *Eighth International Conference on Artificial Intelligence and Soft Computing, ICAISC*, 2006, 870–879.
- [16] D. P. Bertsekas, Projected Newton methods for optimization problems with simple constraints, *SIAM J Control Optimizat* 20(2) (1982), 221–246.
- [17] L. Grippo and M. Sciandrone, On the convergence of the block nonlinear gauss-seidel method under convex constraints, *Operat Res Lett* 26 (2000), 127–136.
- [18] A. Cichocki and R. Zdunek, NMFLAB–MATLAB Toolbox for Non-Negative Matrix Factorization, Online, 2006.
- [19] S. Sra and I. S. Dhillon, Nonnegative Matrix Approximation: Algorithms and Applications. Technical Report Tr-06-27, Computer Sciences, University of Texas at Austin, 2006.
- [20] D. Donoho and V. Stodden, When Does Nonnegative Matrix Factorization Give a Correct Decomposition into Parts? In *Neural Information Processing Systems*, 2003.