



Technical Report No. 193

**Sparse nonnegative matrix  
approximation:  
new formulations and algorithms**

Rashish Tandon<sup>1</sup> and Suvrit Sra<sup>1</sup>

Sep. 13, 2010

**This Technical Report has been approved by:**

Director at MPIK

Postdoc at MPIK



Technical Report No. 193

**Sparse nonnegative matrix  
approximation:  
new formulations and algorithms**

Rashish Tandon<sup>1</sup> and Suvrit Sra<sup>1</sup>

Sep. 13, 2010

<sup>1</sup> MPI für biologische Kybernetik

# Sparse nonnegative matrix approximation: new formulations and algorithms

*Rashish Tandon and Suvrit Sra*

**Abstract.** We introduce several new formulations for sparse nonnegative matrix approximation. Subsequently, we solve these formulations by developing generic algorithms. Further, to help selecting a particular sparse formulation, we briefly discuss the interpretation of each formulation. Finally, preliminary experiments are presented to illustrate the behavior of our formulations and algorithms.

---

## 1 Introduction

A wide variety of applications regularly deal with inherently nonnegative data. Not surprisingly, such data often submit to modern data-analytic techniques such as nonnegative matrix approximation (NMA).<sup>1</sup> But while using NMA to analyze the data, a frequent requirement is to either incorporate prior knowledge, or to obtain models that capture additional structure such as sparsity. Indeed, sparsity has become the mainstay among data analysis requirements owing to important practical benefits such as: (i) simpler, more robust models; (ii) countering of overfitting; (iii) discovery of the most-relevant features; and (iv) potentially better modeling of prior information. Based on these benefits, as well as other motivations, sparsity continues to enjoy great interest in various models and applications [6, 20, 24, 26].

However, despite the ever growing importance of sparsity and its widespread use in numerous applications, sparsity constrained NMA has received comparatively little attention [10, 12]. We make up for this deficiency and present two new formulations for sparse NMA, each of which can be further specialized to yield a series of different sparsity constrained or penalized NMA problems. Our formulations are richer than the previously considered sparse NMA models, and therefore offer greater modeling power. Moreover, our formulations are computationally feasible and for solving them we derive generic algorithms based on traditional multiplicative updates as well as proximal-splitting ideas.

To put our work in context we briefly summarize related work below, before proceeding to the technical details.

### 1.1 Related work

The amount of related work on NMA and its applications has grown too rapidly; the number of references is thus too large to summarize here. For a broad listing of references and survey type summaries we refer the reader to [2, 22]. Below we summarize only the most directly related work.

The idea of sparsity constraints on NMA is natural, but became popular largely after the 2004 paper of Hoyer [10] who introduced a new measure of sparsity based on the ratio of the  $\ell_1$ -norm of a vector to the  $\ell_2$ -norm. This measure allows the user to intuitively tune the sparsity, and it was also used later by Heiler and Schnörr [9], who developed a cone-programming algorithm to solve the associated NMA formulation. While theoretically appealing, this second-order cone-programming formulation can be needlessly expensive. Our sparsity formulations are richer and the associated optimization algorithms simpler, and therefore potentially much faster.

Other authors have considered enforcing sparsity using  $\ell_1$ -norm (trivial choice; no definite reference possible), or even squared  $\ell_1$ -norm penalization [12]. However, to our knowledge, no author has considered sparsity measures for NMA, in the generality treated in this paper. Our sparsity measures are based on the recently popular *mixed-norms*, for which a sizable literature has grown, e.g., [11, 15, 17, 29]; we note, however, that in addition, we also permit mixed-quasinorms.

On the algorithmic front, most of the attention to sparsity based problems has been limited to convex problems, as opposed to this paper where the central problem is non-convex. But the work on convex sparse problems remains highly relevant, because the NMA subproblems that we study, are convex. Thus, we can leverage on the powerful

---

<sup>1</sup>NMA is better known as nonnegative matrix factorization, or just NMF; we prefer using the word ‘approximation’ because usually an exact factorization does not exist.

machinery developed in convex-analysis for tackling sparsity penalized subproblems. Amongst other algorithmic approaches, the most directly are two recent papers: [8, 17]. Both these papers present incremental algorithms that can compute matrix approximations; the methods of Mairal et al. [17] can be easily modified to tackle sparsity if needed. Our approach is different, since we focus on batch methods, as well are more general sparsity constraints.

## 2 NMA: Background

The basic NMA problem is formulated as follows. Let  $\mathbf{A} \in \mathbb{R}_+^{m \times n}$  be an input matrix. NMA aims to obtain a set of nonnegative “basis” vectors  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_k]$  that may be nonnegatively (conically) combined to well-approximate  $\mathbf{A}$ . In symbols, NMA seeks nonnegative  $\mathbf{b}_j$  and  $c_{jk}$  such that

$$\mathbf{a}_i \approx \sum_{j=1}^k \mathbf{b}_j c_{jk}, \quad c_{jk} \geq 0, \quad (2.1)$$

which using matrix notation may be rewritten as

$$\mathbf{A} \approx \mathbf{BC}, \quad \text{where } \mathbf{B} \in \mathbb{R}_+^{m \times k}, \mathbf{C} \in \mathbb{R}_+^{k \times n}. \quad (2.2)$$

Beginning from the formulations (2.1) and (2.2) we may now consider several different alternatives, each of which augments the basic approximation to incorporate other desirable characteristics. For example, suppose that we want the “basis” vectors  $\mathbf{b}_j$  to be sparse, or perhaps we are using NMA for overcomplete dictionary learning (e.g. [17]), and we wish to have sparse combinations, i.e., each column  $\mathbf{c}_i$  of  $\mathbf{C}$  is required to be sparse. Beyond sparsity other constraints may also be imposed on  $\mathbf{B}$  and  $\mathbf{C}$ , leading to different problems such as clustering and co-clustering—see the discussion in [22, Chapter 5] and [25, Chapter 8] for more examples.

**Computing the approximation.** To actually compute an approximation of the form (2.2), we need some quality measure that judges how well  $\mathbf{BC}$  approximates  $\mathbf{A}$ . To that end, we use the *divergence*  $D : \mathbb{R}_+^{m \times k} \times \mathbb{R}_+^{k \times n} \rightarrow \mathbb{R}_+$ , which satisfies the following properties:

1. it is continuously differentiable (at least once) in both variables;
2. it is individually convex in  $\mathbf{B}$  and  $\mathbf{C}$ ;
3. it equals 0, if and only if  $\mathbf{A} = \mathbf{BC}$ .

Notice that  $D$  is a function of  $\mathbf{B}$  and  $\mathbf{C}$ , but is parametrized by  $\mathbf{A}$ , so we choose to write  $D(\mathbf{A}; \mathbf{B}, \mathbf{C})$ . Now the NMA problem may be formulated as

$$\min \quad D(\mathbf{A}; \mathbf{B}, \mathbf{C}) \quad \mathbf{B} \geq 0, \quad \mathbf{C} \geq 0. \quad (2.3)$$

The most common special instance of (2.3) is given the *least-squares* NMA problem:

$$\min \quad D(\mathbf{A}; \mathbf{B}, \mathbf{C}) = \frac{1}{2} \|\mathbf{A} - \mathbf{BC}\|_{\mathbb{F}}^2, \quad \mathbf{B} \geq 0, \quad \mathbf{C} \geq 0. \quad (2.4)$$

Many other distance measures have also been used, including those that do not satisfy the convexity requirement—see [5], or the recent book by Cichocki et al. [3]. We restrict our attention to distance measures that satisfy the two basic properties described above, as these properties are crucial to obtaining efficient algorithms. Furthermore, we will motivate most of our discussion by using (2.4) as the quality measure, though most of the results will actually also hold for the general case.

**Sparsity.** The additional requirements on  $\mathbf{B}$  and  $\mathbf{C}$  may be (approximately) enforced by introducing additional constraints or penalty functions. Our focus in this paper is on modeling various forms of sparsity requirements that either  $\mathbf{B}$  or  $\mathbf{C}$  should satisfy. We assume that this sparsity requirement is encoded by a function  $f(\mathbf{B})$  or  $g(\mathbf{C})$  so that instead of the basic problem (2.3) we consider the more general formulation

$$\begin{aligned} \min \quad & D(\mathbf{A}; \mathbf{B}, \mathbf{C}) \\ & f(\mathbf{B}) \leq \gamma, \quad g(\mathbf{C}) \leq \delta, \quad \mathbf{B} \geq 0, \quad \mathbf{C} \geq 0. \end{aligned} \quad (2.5)$$

We also consider the penalized version of (2.5), where one or both of the sparsity constraints appear as penalizers to the objective function. But for simplicity of exposition, without loss of generality we henceforth that only  $\mathbf{B}$  is required to fulfill sparsity. Our arguments hold with trivial modifications for  $\mathbf{C}$  also, so we may safely omit them.

The most common sparsity measure is the ‘‘counting’’-norm, i.e.,  $f(\mathbf{B}) = \|\mathbf{B}\|_0$ , which counts the number of nonzeros in (support of) matrix  $\mathbf{B}$ . This particular  $f(\mathbf{B})$  is *not* a norm, and using it as a sparsity measure usually results in a combinatorial optimization problem. Therefore, for practicality we prefer the nearest convex relaxation:  $\|\mathbf{B}\|_1$  instead, which also favors sparsity. There also exist several non-convex functions that approximate  $\|\mathbf{B}\|_0$ , but we do not study them in this paper.

Beyond mere vector norms, recently there has been a lot of interest in *group-sparsity*, where instead of computing sparsity of individual vectors (or matrices), one computes a joint sparsity measure across a group of (possibly overlapping) vectors [29]. This type of sparsity will form a central part of our new formulations, but before we formally describe it, let us briefly look at some related work.

### 3 Sparse NMA

Now we are ready to introduce our new sparse nonnegative matrix approximation (SNMA) formulations. But first we need to define some notation.

**Definition 1** (Mixed-norm: Vectors). Let  $\mathbf{w} \in \mathbb{R}^d$  be split<sup>2</sup> into the set  $\{\mathbf{w}_t : \mathbf{w}_t \in \mathbb{R}^{d_t}, 1 \leq t \leq n\}$  of (column) vectors. We define the mixed  $\ell_{p,q}$ -norm (read as  $p$  norm of  $q$  norms) ( $1 \leq p, q \leq \infty$ ) for  $\mathbf{w}$  as

$$\|\mathbf{w}\|_{p,q} = \left\| \left[ \|\mathbf{w}_1\|_q; \|\mathbf{w}_2\|_q; \dots, \|\mathbf{w}_n\|_q \right] \right\|_p. \quad (3.1)$$

That is, we compute the  $\ell_p$ -norm of the vector formed by taking  $\ell_q$ -norms of the subvectors  $\mathbf{w}_t$  ( $1 \leq t \leq n$ ). We may similarly define *mixed-quasi-norms*, where  $p$  or  $q = 0$  are allowed.

We alert the reader that some authors (e.g., [1]) use  $p, q$  in the order opposite to ours. Another point worth noting is that the definition (3.1) can be further generalized, e.g., if we take the  $\ell_{q_t}$ -norm  $\|\mathbf{w}_t\|_{q_t}$  of  $\mathbf{w}_t$ . This generalization comes up for example, while studying  $L_p$ -nested symmetric distributions [21]. Finally, even more generality is gained when the subvectors  $\mathbf{w}_t$  are not required to be disjoint [29]. However, for simplicity, we do not focus on the associated overlapping mixed-norms in this paper.

Recall that without loss of generality we may discuss sparsity for  $\mathbf{B}$  alone; the matrix  $\mathbf{C}$  can be treated similarly if desired. Assume that the sparsity requirement on matrix  $\mathbf{B}$  is encoded by a function  $f : \mathbb{R}_+^{m \times k} \rightarrow \mathbb{R}_+$ . We consider both *penalized* and *constrained* versions of SNMA:

<b>Penalized SNMA</b>
<p>Given a matrix <math>\mathbf{A}_{m \times n}</math>, a natural number <math>k \leq \min(m, n)</math>, and a parameter <math>\lambda &gt; 0</math>, find matrices <math>\mathbf{B}_{m \times k}</math> and <math>\mathbf{C}_{k \times n}</math> that solve</p> $\begin{aligned} &\text{minimize} && D(\mathbf{A}; \mathbf{B}, \mathbf{C}) + \lambda f(\mathbf{B}) \\ &\text{s.t.} && \mathbf{B}, \mathbf{C} \geq 0. \end{aligned} \quad (3.2)$
<b>Constrained SNMA</b>
<p>Given a matrix <math>\mathbf{A}_{m \times n}</math>, a natural number <math>k \leq \min(m, n)</math>, and a parameter <math>\gamma &gt; 0</math>, find matrices <math>\mathbf{B}_{m \times k}</math> and <math>\mathbf{C}_{k \times n}</math> that solve</p> $\begin{aligned} &\text{minimize} && D(\mathbf{A}; \mathbf{B}, \mathbf{C}) \\ &\text{subject to} && \mathbf{B}, \mathbf{C} \geq 0, \quad f(\mathbf{B}) \leq \gamma. \end{aligned} \quad (3.3)$

We specialize (3.2) and (3.3) below by describing several practical choices; our objective is to illustrate how different ‘kinds’ of sparsity can be modeled by different choices of norms for  $f(\mathbf{B})$ . Finally, we also look at a possibility for  $f(\mathbf{B})$  which seeks to promote sparsity in  $\mathbf{B}$  while maximizing diversity across its columns.

#### 3.1 Overall sparsity

The function

$$f(\mathbf{B}) = \|\mathbf{B}\|_{1,0} = \sum_i \|\mathbf{b}^i\|_0,$$

<sup>2</sup>For computational efficiency we restrict our attention where the splitting *partitions* the vector  $\mathbf{w}$  into subvectors; though generally overlapping groups may also be considered.

promotes overall sparsity in  $\mathbf{B}$ , since it exactly counts the total number of non-zero elements in  $\mathbf{B}$ . Notice that  $\|\mathbf{B}\|_{1,0} = \|\mathbf{B}^T\|_{1,0}$ , which tells us that this function does not distinguish between rows and columns.

A relaxation to the  $\ell_{1,0}$ -quasi-norm is the  $\ell_{1,1}$ -norm:

$$f(\mathbf{B}) = \|\mathbf{B}\|_{1,1} = \sum_{ij} |b_{ij}|.$$

This is the familiar (vector)  $\ell_1$ -norm, and is known to induce sparsity by shrinking small elements of  $\mathbf{B}$  toward 0. The extent of sparsity attained is, however, harder to control, though because it is a separable, convex function, the optimization burden is less demanding. Here too, notice that  $\|\mathbf{B}\|_{1,1} = \|\mathbf{B}^T\|_{1,1}$ , which means that this norm promotes overall sparsity without distinguishing between rows and columns.

Another possible choice for promoting overall sparsity in  $\mathbf{B}$  is the  $\ell_{1,2}$ -norm

$$f(\mathbf{B}) = \|\mathbf{B}\|_{1,2} = \sum_i \|\mathbf{b}^i\|_2,$$

where  $\mathbf{b}^i$  is the  $i$ -th row of  $\mathbf{B}$ . Observe that this norm favors making the  $\ell_2$ -norm of each row as small as possible, rather than merely focusing on individual elements. Thus, even though this choice does promote overall sparsity, there is a relatively greater stress on rows. Note that in this case,  $\|\mathbf{B}\|_{1,2} \neq \|\mathbf{B}^T\|_{1,2}$  in general.

### 3.2 Group sparsity: Maximum Row-wise Budget

The penalty function

$$f(\mathbf{B}) = \|\mathbf{B}\|_{\infty,0} = \max_{1 \leq i \leq m} \{\|\mathbf{b}^i\|_0\},$$

counts the maximum number of non-zero elements over all the rows. Thus, with this choice of  $f(\mathbf{B})$ , we can impose an upper bound on the maximum number of non-zero elements that any row of  $\mathbf{B}$  is allowed to have. If needed, one could also subject the columns of  $\mathbf{B}$  to a similar constraint by using  $f(\mathbf{B}) = \|\mathbf{B}^T\|_{\infty,0}$ . The great benefit of this function is its explicit ability to enforce a “sparsity-budget” on either the rows or columns of  $\mathbf{B}$ .

One disadvantage of the above choice of  $f(\mathbf{B})$  is its non-convexity. We can instead use the convex relaxation

$$f(\mathbf{B}) = \|\mathbf{B}\|_{\infty,1} = \max_{1 \leq i \leq m} \{\|\mathbf{b}^i\|_1\},$$

which is a natural relaxation to the  $\ell_{\infty,0}$ -quasi-norm. This choice also imposes a maximum budget on the rows (or columns), but now in terms of the  $\ell_1$ -norm.

### 3.3 Group sparsity: Zeroing out insignificant rows

Observe that  $\|\mathbf{B}\|_{0,0} = \|\mathbf{B}\|_{0,1} = \|\mathbf{B}\|_{0,\infty}$  counts the total number of *non-zero rows* in  $\mathbf{B}$ . Therefore, the penalty function  $f(\mathbf{B}) = \|\mathbf{B}\|_{0,0}$  would have a tendency to favor setting the less important rows in  $\mathbf{B}$  to zero. When the rows of  $\mathbf{B}$  correspond to features, this choice of  $f(\mathbf{B})$  allows performing *feature selection* by eliminating less significant features from the matrix  $\mathbf{B}$ .

Since the above choice of  $f(\mathbf{B})$  is non-convex, we might prefer to use

$$f(\mathbf{B}) = \|\mathbf{B}\|_{1,\infty} = \sum_i \|\mathbf{b}^i\|_\infty$$

instead, which may be viewed as a convex relaxation to the  $\ell_{0,0}$ -quasinorm. This choice of  $f$  favors reducing  $\ell_\infty$ -norms of the rows of  $\mathbf{B}$ ; if for some row  $\mathbf{b}^i$ , the norm  $\|\mathbf{b}^i\|_\infty$  is set to 0, then the entire row must be zero, thereby achieving the aforementioned feature selection effect.

### 3.4 Sparsity with Diversity

Sparsity with diversity seeks columns of  $\mathbf{B}$  that have non-zero elements at different locations in each column. For example,  $\mathbf{B}$  with orthogonal columns is maximally diverse, though not necessarily with sparse columns.

Two simple functions to enforce diversity could be

$$f(\mathbf{B}) = \|\mathbf{B}^T \mathbf{B} - \text{Diag}(\mathbf{B}^T \mathbf{B})\|_{\mathbb{F}}^2, \quad \text{or} \quad f(\mathbf{B}) = \|\mathbf{B}^T \mathbf{B} - \mathbf{I}\|_{\mathbb{F}}^2.$$

An advantage of using the latter as compared to the former is that the former permits solutions where the inner product of a columns with itself can be close to zero, while the latter favors where this inner product is close to one. Both these diversity promoting functions can be combined with any of the abovementioned sparsity measures to jointly enforce sparsity with diversity. We note, however, that using these diversity measures seems to be harder for the constrained formulation (3.3), so we shall only look at solving (3.2) with these measures.

## 4 Sparse NMA: Algorithms

The NMA problem is a difficult non-convex optimization problem. Indeed, Vavasis [28] showed that it is NP-Hard to find the global optimum for NMA, and so most algorithms can at best hope to guarantee locally optimum solutions. Our Sparse NMA formulations that utilize a true norm (like  $\|\mathbf{B}\|_{1,1}$ ,  $\|\mathbf{B}\|_{1,2}$ ,  $\|\mathbf{B}\|_{1,\infty}$ ,  $\|\mathbf{B}\|_{\infty,1}$ ) can be reduced to the standard NMA problem and vice-versa. Thus, our SNMA formulations are essentially equivalent to the basic NMA problem. The NMA problem can also be reduced to our SNMA formulations that use quasi-norms (like  $\|\mathbf{B}\|_{1,0}$ ,  $\|\mathbf{B}\|_{0,0}$ ,  $\|\mathbf{B}\|_{\infty,0}$ ), showing that these problems are potentially harder than basic NMA. In either case, it follows that it is NP-Hard to find a global optimum for our SNMA formulations as well.

Although non-convex, and NP-Hard, the SNMA problems do have nice structure which allows the use of *alternating-descent* procedures. Here, one iteratively minimizes (or merely descends) with respect to one of the matrices  $\mathbf{B}$  or  $\mathbf{C}$ , at a time. Such an alternating approach is appealing because upon fixing one of the matrices, the optimization problems becomes convex in the other.

Formally, if  $\mathbf{B}_t$  and  $\mathbf{C}_t$  are the matrices at the  $t$ -th iteration, then we can obtain  $\mathbf{B}_{t+1}$  and  $\mathbf{C}_{t+1}$  as follows:

$$\mathbf{B}_{t+1} \leftarrow \operatorname{argmin}_{\mathbf{B} \geq 0} D(\mathbf{A}; \mathbf{B}, \mathbf{C}_t) + \lambda f(\mathbf{B}), \quad \text{for (3.2)} \quad (4.1)$$

$$\mathbf{B}_{t+1} \leftarrow \operatorname{argmin}_{\mathbf{B} \geq 0, f(\mathbf{B}) \leq \gamma} D(\mathbf{A}; \mathbf{B}, \mathbf{C}_t), \quad \text{for (3.3)} \quad (4.2)$$

$$\mathbf{C}_{t+1} \leftarrow \operatorname{argmin}_{\mathbf{C} \geq 0} D(\mathbf{A}; \mathbf{B}_{t+1}, \mathbf{C}). \quad (4.3)$$

Thanks to the *argmin*, descent is guaranteed, i.e.,

$$D(\mathbf{A}; \mathbf{B}_t, \mathbf{C}_t) + \lambda f(\mathbf{B}_t) \geq D(\mathbf{A}; \mathbf{B}_{t+1}, \mathbf{C}_{t+1}) + \lambda f(\mathbf{B}_{t+1}), \quad \text{for (3.2)} \quad (4.4)$$

$$D(\mathbf{A}; \mathbf{B}_t, \mathbf{C}_t) \geq D(\mathbf{A}; \mathbf{B}_{t+1}, \mathbf{C}_{t+1}), \quad \text{for (3.3)} \quad (4.5)$$

Under appropriate assumptions one can then show (e.g., using [7]) that the above procedure converges to a local optimum for the overall problem.

However, in general exact minimization in the subproblems might be either overkill, or perhaps complicated and computationally expensive. In such a case it is more useful to replace the *argmin* which merely a guaranteed descent. Although doing so destroys the theoretical guarantees of convergence to a local optimum, empirically the approach may still work well. In principle, algorithms performing mere descent can be augmented via appropriate line-searches (that guarantee sufficient descent) to ensure convergence to a stationary point.

Let us now look at methods that can be used to solve the *nonsmooth* optimization problems (4.1) and (4.2). We prefer particularly simple algorithms, which can also be stopped before reaching the minimum, in case a desired level of descent has been achieved.

### 4.1 Forward-Backward Splitting Algorithms

Recall that Forward-Backward Splitting (FBS) algorithms [4] are iterative procedures for solving the following optimization problem:

$$\operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} f_1(\mathbf{x}) + f_2(\mathbf{x}), \quad (4.6)$$

where  $f_1(\mathbf{x})$  is a smooth, convex function, while  $f_2$  is convex, but not necessarily smooth. To solve (4.6), FBS algorithms depend on proximity operators, which are defined as:

**Definition 2.** Let  $f_2 : \mathbb{R}^n \rightarrow \mathbb{R}$  be a lower semi-continuous convex function. The *proximity-operator* for  $f_2$  is denoted by  $\operatorname{prox}_{f_2}$ , and its application at a point  $\mathbf{y}$  is defined by

$$\operatorname{prox}_{f_2}(\mathbf{y}) \equiv \operatorname{argmin}_{\mathbf{x} \in \operatorname{dom} f_2} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + f_2(\mathbf{x}). \quad (4.7)$$

It can be shown that if  $f_1(\mathbf{x})$  is a convex differentiable function with a  $\beta$ -Lipschitz continuous gradient  $\nabla f_1$ , and  $f_2(\mathbf{x})$  is a convex lower semi-continuous function, then (4.6) admits at least one solution; this solution is characterized by the following fixed point equation (for any  $\alpha > 0$ ) (see [4] for details):

$$\mathbf{x} \leftarrow \operatorname{prox}_{\alpha f_2}(\mathbf{x} - \alpha \nabla f_1(\mathbf{x}))$$

This fixed-point equation suggests the following *forward-backward splitting* procedure

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \theta_t (\text{prox}_{\alpha_t f_2}(\mathbf{x}_t - \alpha_t \nabla f_1(\mathbf{x}_t)) - \mathbf{x}_t), \quad (4.8)$$

where  $\alpha_t \in [\epsilon, 2/\beta - \epsilon]$  and  $\theta_t \in [\epsilon, 1]$  for a fixed  $\epsilon \in (0, \min\{1, 1/\beta\})$ . It can be shown (see [4]) that the sequence  $\mathbf{x}_t$  generated via (4.8) converges to a solution of (4.6). A common implementation choice for (4.8) is  $\theta_t = 1$ , for which the iteration simplifies to

$$\mathbf{x}_{t+1} = \text{prox}_{\alpha_t f_2}(\mathbf{x}_t - \alpha_t \nabla f_1(\mathbf{x}_t)). \quad (4.9)$$

We use iteration (4.9) to minimize over  $\mathbf{B}$ , noting that a similar choice could also be made for  $\mathbf{C}$ . Moreover, to save time, we need not run (4.9) to convergence, and may accept an intermediate iterate satisfying descent.

## 4.2 Details for (3.2)

For the penalized formulation (3.2), descending along  $\mathbf{B}$  requires solving

$$\underset{\mathbf{B} \geq 0}{\text{argmin}} \quad D(\mathbf{A}; \mathbf{B}, \mathbf{C}_t) + \lambda f(\mathbf{B}). \quad (4.10)$$

First, define the indicator function for a set  $\mathcal{B}$

$$\mathbb{1}_{\mathcal{B}}(\mathbf{B}) = \begin{cases} 0, & \mathbf{B} \in \mathcal{B} \\ \infty, & \text{otherwise.} \end{cases} \quad (4.11)$$

Now, with  $\mathcal{B} = \mathbb{R}_+^{m \times k}$  (i.e., the nonnegativity constraints  $\mathbf{B} \geq 0$ ), we can rewrite (4.10) as,

$$\underset{\mathbf{B}}{\text{argmin}} \quad \underbrace{D(\mathbf{A}; \mathbf{B}, \mathbf{C}_t)}_{f_1} + \underbrace{\lambda f(\mathbf{B}) + \mathbb{1}_{\geq 0}(\mathbf{B})}_{f_2}. \quad (4.12)$$

Now we may invoke (4.9) to solve (4.12), using the indicated choices for  $f_1$  and  $f_2$ . Notice that  $f_2 = \lambda f + \mathbb{1}_{\geq 0}$  is a sum of two functions. Normally, computing the proximity operator for a sum of two functions is not easy. However, below we show how one can compute  $\text{prox}_{f_2}$  efficiently. Another point worth mentioning is that to invoke (4.9), the divergence  $f_1(\mathbf{B}) = D(\mathbf{A}; \mathbf{B}, \mathbf{C}_t)$  must be a convex, differentiable function with a  $\beta$ -Lipschitz continuous, while  $f_2$  must a convex lower semi-continuous function. However, not all choices of  $f$  that we have previously mentioned lead to a lower semi-continuous  $f_2$ . But we do not exclude such  $f$ , because we can compute the corresponding proximity operators efficiently.

## 4.3 Details for (3.3)

For the constrained formulation (3.3), descending along  $\mathbf{B}$  requires solving

$$\underset{\mathbf{B} \geq 0, f(\mathbf{B}) \leq \gamma}{\text{argmin}} \quad D(\mathbf{A}; \mathbf{B}, \mathbf{C}_t). \quad (4.13)$$

We move both constraints into the objective, and obtain the problem

$$\underset{\mathbf{B}}{\text{argmin}} \quad \underbrace{D(\mathbf{A}; \mathbf{B}, \mathbf{C}_t)}_{f_1} + \underbrace{\mathbb{1}_{\geq 0}(\mathbf{B}) + \mathbb{1}_{f \leq \gamma}(\mathbf{B})}_{f_2}. \quad (4.14)$$

Here, we need to compute  $\text{prox}_{f_2} = \text{prox}_{\mathbb{1}_{\geq 0} + \mathbb{1}_{f \leq \gamma}}$ , where as alluded to above, not all choices of  $f$  are convex and lower semi-continuous, although the proximity operator can still be computed efficiently.

Since we solve both (4.10) and (4.13) via the FBS iteration (4.9), we need to efficiently compute the associated proximity operators. We discuss this computation below.

## 4.4 Computing the proximity operators

Computing  $\text{prox}_{f_2}$  requires solving one of the following optimization problems:

$$\begin{aligned} & \underset{\mathbf{B} \geq 0}{\text{argmin}} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_2^2 + \lambda f(\mathbf{B}) \\ & \underset{\mathbf{B} \geq 0}{\text{argmin}} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_2^2, \quad \text{s.t.} \quad f(\mathbf{B}) \leq \gamma. \end{aligned}$$

To simplify solving these optimization problems, we derive two simple lemmas below.



**Lemma 3** (Signs). *Let  $f : \mathbb{R}^{m \times k} \rightarrow \mathbb{R}$  be an elementwise even function,  $\mathbf{B}_1^* = \text{prox}_{\lambda f(\mathbf{B})}(\mathbf{B}')$ , and  $\mathbf{B}_2^* = \text{prox}_{\mathbb{1}_{f \leq \gamma}}(\mathbf{B}')$ . Then, we must have*

$$\left. \begin{aligned} \text{sgn}(\mathbf{B}_1^*(i, j)) &= \text{sgn}(\mathbf{B}'(i, j)) \text{ or } 0, \\ \text{sgn}(\mathbf{B}_2^*(i, j)) &= \text{sgn}(\mathbf{B}'(i, j)) \text{ or } 0, \end{aligned} \right\} \quad i = 1, \dots, m, \text{ and } j = 1, \dots, k. \quad (4.15)$$

*Proof.* Since  $\mathbf{B}_1^* = \text{prox}_{\lambda f}(\mathbf{B}')$ , it is the optimal solution to

$$\underset{\mathbf{B}}{\text{argmin}} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_2^2 + \lambda f(\mathbf{B}).$$

To prove the lemma, we show that for all  $i, j$ , the inequality  $\mathbf{B}_1^*(i, j)\mathbf{B}'(i, j) \geq 0$  holds. Thus, suppose by way of contradiction that there exist  $p, q$  such that  $\mathbf{B}_1^*(p, q)\mathbf{B}'(p, q) < 0$ . Then, we can construct  $\hat{\mathbf{B}}$  such that,

$$\hat{\mathbf{B}}(i, j) = \begin{cases} \mathbf{B}_1^*(i, j) & \text{for } (i, j) \neq (p, q), \\ -\mathbf{B}_1^*(p, q) & \text{otherwise.} \end{cases} \quad (4.16)$$

Since we assumed  $f$  to be elementwise even, it holds that  $f(\mathbf{B}_1^*) = f(\hat{\mathbf{B}})$ . Now, consider the difference

$$\begin{aligned} & \frac{1}{2} \|\mathbf{B}_1^* - \mathbf{B}'\|_2^2 + \lambda f(\mathbf{B}_1^*) - \frac{1}{2} \|\hat{\mathbf{B}} - \mathbf{B}'\|_2^2 - \lambda f(\hat{\mathbf{B}}) \\ &= \frac{1}{2} (\mathbf{B}_1^*(p, q) - \mathbf{B}'(p, q))^2 - \frac{1}{2} (\mathbf{B}_1^*(p, q) + \mathbf{B}'(p, q))^2 \\ &= -2(\mathbf{B}_1^*(p, q)\mathbf{B}'(p, q)) > 0, \end{aligned}$$

which contradicts the optimality of  $\mathbf{B}_1^*$ . Similarly, one can also prove the inequality  $\mathbf{B}_2^*(i, j)\mathbf{B}'(i, j) \geq 0$ .  $\square$

**Lemma 4.** *Let  $f : \mathbb{R}^{m \times k} \rightarrow \mathbb{R}$  be an elementwise function for which*

$$f(\mathbf{B}_1, \dots, \underbrace{0}_{i^{\text{th}} \text{ pos}}, \dots, \mathbf{B}_{mn}) \leq f(\mathbf{B}_1, \dots, \mathbf{B}_i, \dots, \mathbf{B}_{mn}), \quad i \in [mn], \mathbf{B}_i \in \mathbb{R}_+, \text{ and } \mathbf{B}_j \in \mathbb{R} \text{ for } j \neq i.$$

*Then for  $\mathbf{B}_1^* = \text{prox}_{\lambda f + \mathbb{1}_{\geq 0}}(\mathbf{B}')$  and  $\mathbf{B}_2^* = \text{prox}_{\mathbb{1}_{\geq 0} + \mathbb{1}_{f \leq \gamma}}(\mathbf{B}')$ , we must have:*

$$\left. \begin{aligned} \text{sgn}(\mathbf{B}_1^*(p, q)) &= 0 \quad \text{if } \text{sgn}(\mathbf{B}'(p, q)) = -1 \quad \text{or } 0, \\ \text{sgn}(\mathbf{B}_2^*(p, q)) &= 0 \quad \text{if } \text{sgn}(\mathbf{B}'(p, q)) = -1 \quad \text{or } 0, \end{aligned} \right\} \quad p \in [m], q \in [k].$$

*Proof.* Since  $\mathbf{B}_1^* = \text{prox}_{\lambda f + \mathbb{1}_{\geq 0}}(\mathbf{B}')$ , it is the optimal solution to

$$\underset{\mathbf{B}_{\geq 0}}{\text{argmin}} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_2^2 + \lambda f(\mathbf{B}).$$

We again proceed by contradiction; thus, let there be some  $p, q$  such that  $\mathbf{B}'(p, q) \leq 0$ , and  $\mathbf{B}_1^*(p, q) > 0$ . Construct  $\hat{\mathbf{B}}$  such that,

$$\hat{\mathbf{B}}(i, j) = \begin{cases} \mathbf{B}_1^*(i, j) & \text{if } (i, j) \neq (p, q), \\ 0 & \text{otherwise.} \end{cases}$$

Thus,  $\hat{\mathbf{B}} \geq 0$  and  $f(\mathbf{B}_1^*) \geq f(\hat{\mathbf{B}})$ . Now, consider the difference

$$\begin{aligned} & \frac{1}{2} \|\mathbf{B}_1^* - \mathbf{B}'\|_2^2 + \lambda f(\mathbf{B}_1^*) - \frac{1}{2} \|\hat{\mathbf{B}} - \mathbf{B}'\|_2^2 - \lambda f(\hat{\mathbf{B}}) \\ &= \frac{1}{2} (\mathbf{B}_1^*(p, q) - \mathbf{B}'(p, q))^2 - \frac{1}{2} (\mathbf{B}'(p, q))^2 + \lambda (f(\mathbf{B}_1^*) - f(\hat{\mathbf{B}})) \\ &= \frac{1}{2} (\mathbf{B}_1^*(p, q))^2 - (\mathbf{B}_1^*(p, q)\mathbf{B}'(p, q)) + \lambda (f(\mathbf{B}_1^*) - f(\hat{\mathbf{B}})) > 0, \end{aligned}$$

which contradicts the optimality of  $\mathbf{B}_1^*$ . The proof for  $\mathbf{B}_2^*$  follows similarly.  $\square$

**Corollary 5.** *Let  $f$  be as in Lemma 4, and let  $P_+(\mathbf{X})$  denote projection onto the nonnegative orthant. Then,*

$$\text{prox}_{\lambda f + \mathbb{1}_{\geq 0}}(\mathbf{B}') = \text{prox}_{\lambda f}(P_+(\mathbf{B}')), \quad (4.17)$$

$$\text{prox}_{\mathbb{1}_{\geq 0} + \mathbb{1}_{f \leq \gamma}}(\mathbf{B}') = \text{prox}_{\mathbb{1}_{f \leq \gamma}}(P_+(\mathbf{B}')). \quad (4.18)$$

*Proof.* By Lemma 4, we see that replacing the negative entries of  $\mathbf{B}'$  by 0, does not affect  $\text{prox}_{\lambda f + \mathbb{1}_{\geq 0}}(\mathbf{B}')$ . Thus,

$$\text{prox}_{\lambda f + \mathbb{1}_{\geq 0}}(\mathbf{B}') = \text{prox}_{\lambda f + \mathbb{1}_{\geq 0}}(P_+(\mathbf{B}')).$$

Now, since  $P_+(\mathbf{B}') \geq 0$ , by Lemma 3, it follows that  $\text{prox}_{\lambda f}(P_+(\mathbf{B}'))$  must have all non-negative entries. Thus, enforcing  $\mathbf{B} \geq 0$  does not change the solution, and (4.17) follows; a similar argument holds for (4.18).  $\square$

All the norms and quasi-norms (as choices for  $f(\mathbf{B})$ ) introduced above satisfy the conditions set forth in Corollary 5. Thus, we can compute the corresponding proximity operators with the constraint  $\mathbf{B} \geq 0$ , by merely replacing  $\mathbf{B}'$  by  $P_+(\mathbf{B}')$ . Let us look at details of specific choices for  $f(\mathbf{B})$  below.

#### 4.4.1 Proximity and projection for $\ell_{1,0}$ -quasinorm

Let  $f(\mathbf{B}) = \|\mathbf{B}\|_{1,0}$ . Computing  $\text{prox}_{\lambda f}(\mathbf{B}')$  requires solving

$$\underset{\mathbf{B}}{\text{argmin}} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\text{F}}^2 + \lambda \|\mathbf{B}\|_{1,0}.$$

This problem has the following simple closed-form solution:

$$\mathbf{B}^*(i, j) = \begin{cases} \mathbf{B}'(i, j), & \text{if } |\mathbf{B}'(i, j)| > \sqrt{2\lambda}, \\ 0 & \text{otherwise.} \end{cases} \quad (4.19)$$

Thus, for an  $m \times k$  matrix  $\mathbf{B}'$ , the proximity operator can be computed in  $O(mk)$ , i.e., linear time.

The proximity (actually, merely projection)  $\text{prox}_{\mathbb{1}_{f \leq \gamma}}(\mathbf{B}')$  requires solving

$$\underset{\|\mathbf{B}\|_{1,0} \leq \gamma}{\text{argmin}} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\text{F}}^2.$$

This problem can be solved by simply picking the top  $\lfloor \gamma \rfloor$  elements from  $\mathbf{B}'$ , and using them as the corresponding entries in  $\mathbf{B}$ ; the remainder of the elements of  $\mathbf{B}$  are set to 0. Using an order-statistic algorithm (median of medians procedure), this solution can be computed in  $O(mk)$  time.

*Remark:* Note that  $\|\mathbf{B}\|_{1,0}$  is not a convex function and  $\{\mathbf{B} \mid \|\mathbf{B}\|_{1,0} \leq \gamma\}$  is not a convex set. Hence, convergence guarantees of a FBS procedure based on using the corresponding proximity operator do not hold. However, convergence guarantees hold for the relaxed versions of  $\|\mathbf{B}\|_{1,0}$ , namely the norms:  $\|\mathbf{B}\|_{1,1}$  and  $\|\mathbf{B}\|_{1,2}$ .

#### 4.4.2 Proximity and projection for $\ell_{1,1}$ -norm

Let  $f(\mathbf{B}) = \|\mathbf{B}\|_{1,1}$ . Here we must solve

$$\underset{\mathbf{B}}{\text{argmin}} \quad \|\mathbf{B} - \mathbf{B}'\|_{\text{F}}^2 + \lambda \|\mathbf{B}\|_{1,1}.$$

The solution to this problem is well-known, and is called *soft-thresholding*:

$$\mathbf{B}^*(i, j) = \text{sgn}(\mathbf{B}'(i, j)) \max(|\mathbf{B}'(i, j)| - \lambda, 0). \quad (4.20)$$

This operation can obviously be computed in  $O(mk)$  time.

The corresponding projection problem, i.e.,  $\text{prox}_{\mathbb{1}_{f \leq \gamma}}(\mathbf{B}')$  requires solving

$$\underset{\|\mathbf{B}\|_{1,1} \leq \gamma}{\min} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\text{F}}^2.$$

This projection is also very well-known, e.g., [13, 16, 18], and can be solved in linear time.

#### 4.4.3 Proximity and projection for $\ell_{1,2}$ -norm

Let  $f(\mathbf{B}) = \|\mathbf{B}\|_{1,2}$ . Here we must solve

$$\underset{\mathbf{B}}{\min} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\text{F}}^2 + \lambda \|\mathbf{B}\|_{1,2}.$$

Using separability of  $\|\cdot\|_{\text{F}}^2$  and  $\|\cdot\|_{1,2}$ , the above problems can be broken into  $m$  independent optimization problems of the form:

$$\underset{\mathbf{b}}{\min} \quad \|\mathbf{b} - \mathbf{b}'\|_2^2 + \lambda \|\mathbf{b}\|_2,$$

which has the well-known closed-form solution

$$\mathbf{b}^* = \max \{1 - \lambda / \|\mathbf{b}'\|_2, 0\} \mathbf{b}'. \quad (4.21)$$

Thus, the  $\ell_{1,2}$ -norm proximity requires only linear  $O(mk)$  time.

The corresponding projection problem is

$$\min_{\|\mathbf{B}\|_{1,2} \leq \gamma} \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\text{F}}^2,$$

which can be solved in linear time using an algorithm of van den Berg et al. [27].

#### 4.4.4 Proximity and projection for $\ell_{0,0}$ -quasinorm

Let  $f(\mathbf{B}) = \|\mathbf{B}\|_{0,0}$ . Here we must solve

$$\min_{\mathbf{B}} \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\text{F}}^2 + \lambda \|\mathbf{B}\|_{0,0}.$$

One can obtain the optimal solution by noting that the  $i$ -th row  $\mathbf{b}_i^*$ , of the optimal solution  $\mathbf{B}^*$  is given by

$$\mathbf{b}_i^* = \begin{cases} \mathbf{b}'_i, & \text{if } \|\mathbf{b}'_i\|_2 > \sqrt{2\lambda}, \\ 0 & \text{otherwise.} \end{cases} \quad (4.22)$$

Thus this proximity operator too can be computed in linear time.

The corresponding projection requires solving

$$\min_{\|\mathbf{B}\|_{0,0} \leq \gamma} \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\text{F}}^2.$$

The optimal solution is comprised of the top  $\lceil \gamma \rceil$  rows (ordered by  $\|\cdot\|_2$ ) of  $\mathbf{B}'$  copied over into  $\mathbf{B}$ ; the other rows of  $\mathbf{B}$  are set to zero. Again, using a linear-time median of medians algorithm, the top rows can be determined in linear time.

*Remark:* Due to non-convexity of  $\|\mathbf{B}\|_{0,0}$ , the associated convergence guarantees for FBS schemes do not hold. Convergence does hold, however, with the convex relaxation  $\|\mathbf{B}\|_{1,\infty}$ .

#### 4.4.5 Proximity and projection for $\ell_{1,\infty}$ -norm

Let  $f(\mathbf{B}) = \|\mathbf{B}\|_{1,\infty}$ . Here we must solve

$$\min_{\mathbf{B}} \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\text{F}}^2 + \lambda \|\mathbf{B}\|_{1,\infty}.$$

This problem can be again decomposed into  $m$  independent optimization problems of the form

$$\min_{\mathbf{b}} \frac{1}{2} \|\mathbf{b} - \mathbf{b}'\|_{\text{F}}^2 + \lambda \|\mathbf{b}\|_{\infty}.$$

Recalling that  $\|\mathbf{b}\|_1$  is the dual-norm to  $\|\mathbf{b}\|_{\infty}$ , we can express the optimal solution to the above problem as

$$\mathbf{b}^* = \mathbf{b}' - \text{prox}_{\mathbb{1}_{\|\cdot\|_1 \leq \lambda}}(\mathbf{b}'),$$

which requires merely solving an  $\ell_1$ -norm projection, which can be done in linear time as previously described.

The corresponding projection task is

$$\min_{\|\mathbf{B}\|_{1,\infty} \leq \gamma} \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\text{F}}^2.$$

Quattoni et al. [19] provided an algorithm for performing this projection in  $O(mk \log k)$  time; though empirically the techniques described in [23] result in much faster computation.

#### 4.4.6 Proximity and projection for $\ell_{\infty,0}$ -quasinorm

Let  $f(\mathbf{B}) = \|\mathbf{B}\|_{\infty,0}$ . Here we must solve

$$\min_{\mathbf{B}} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\mathbb{F}}^2 + \lambda \|\mathbf{B}\|_{\infty,0}.$$

Let  $\mathbf{b}'_i{}^l$  represent the  $l$ -largest element in the  $i$ -th row,  $\mathbf{b}'_i$  of  $\mathbf{B}'$ . To solve the above problem, we need to find the appropriate rank  $p$ , such that

$$\sum_{i=1}^m (\mathbf{b}'_i{}^p)^2 \geq 2\lambda, \quad \text{and} \quad \sum_{i=1}^m (\mathbf{b}'_i{}^{p+1})^2 < 2\lambda.$$

Once such a  $p$  has been found, the  $i^{\text{th}}$  row of the optimal solution is obtained by

$$\mathbf{b}_i^*(j) = \begin{cases} \mathbf{b}'_i(j), & \text{if } \mathbf{b}'_i(j) \geq \mathbf{b}'_i{}^p, \\ 0 & \text{otherwise.} \end{cases} \quad (4.23)$$

If the rows of the matrix  $\mathbf{B}'$  are sorted, then the required  $p$  can be found in linear time. Thus, sorting dominates the overall runtime, implying an  $O(mk \log k)$  runtime.

The corresponding projection task is

$$\min_{\|\mathbf{B}\|_{\infty,0} \leq \gamma} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\mathbb{F}}^2.$$

This computation decomposes into  $m$  independent optimization problems of the form

$$\min_{\|\mathbf{b}\|_0 \leq \gamma} \quad \frac{1}{2} \|\mathbf{b} - \mathbf{b}'\|_2^2.$$

Each of these problems can be solved easily; the optimal solution is obtained by picking the top  $\lfloor \gamma \rfloor$  elements from  $\mathbf{b}'$ , and setting the rest to zero. As before, a median of medians algorithm allows doing this in linear time, thus the overall runtime is  $O(mk)$ .

*Remark:* Since  $\|\mathbf{B}\|_{\infty,0}$  is non-convex, it may be desirable to replace it with its convex relaxation,  $\|\mathbf{B}\|_{\infty,1}$ .

#### 4.4.7 Proximity and projection for $\ell_{\infty,1}$ -norm

Let  $f(\mathbf{B}) = \|\mathbf{B}\|_{\infty,1}$ . Here we must solve

$$\min_{\mathbf{B}} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\mathbb{F}}^2 + \lambda \|\mathbf{B}\|_{\infty,1}.$$

Since the  $\ell_{\infty,1}$ -norm is dual to the  $\ell_{1,\infty}$ -norm, the optimal solution  $\mathbf{B}^*$  to the above problem can be expressed as

$$\mathbf{B}^* = \mathbf{B}' - \text{prox}_{\lambda \|\cdot\|_{1,\infty}}(\mathbf{B}').$$

We already described how to efficiently compute the  $\ell_{1,\infty}$ -norm projection involved above.

The corresponding projection problem is

$$\min_{\|\mathbf{B}\|_{\infty,1} \leq \gamma} \quad \frac{1}{2} \|\mathbf{B} - \mathbf{B}'\|_{\mathbb{F}}^2,$$

which can be decomposed into  $m$  independent  $\ell_1$ -norm ball projection problems, each solvable in linear time. Thus, the overall projection can be computed in  $O(mk)$  time.

### 4.5 Alternative: Multiplicative Update Algorithms

A popular alternative for solving NMA subproblems are multiplicative update algorithms. These algorithms ensure descent for the alternating sub-problems over  $\mathbf{B}$  or  $\mathbf{C}$ ; but they do not obtain convergence to the minimum of each separate alternate minimization, as they usually perform only one descent iteration at a time.

When the regularizer  $f(\mathbf{B}) = 0$ , several multiplicative algorithms are known depending on the choice of the divergence measure  $D(\mathbf{A}; \mathbf{B}, \mathbf{C})$ ; the most well-known amongst these are the updates of Lee and Seung [14].

For  $D(\mathbf{A}; \mathbf{B}, \mathbf{C})$  being the least squares distance<sup>3</sup>, we provide two new update rules that guarantee descent on our penalized formulation (3.2), when  $f(\mathbf{B})$  is either  $\|\mathbf{B}\|_{1,1}$  or  $\|\mathbf{B}\|_{1,2}$ . Since we penalize only  $\mathbf{B}$ , the updates for  $\mathbf{C}$  remain the usual ones.

<sup>3</sup>Although we do not explicitly state the update rules, yet it is fairly straightforward to derive update rules even when  $D(\mathbf{A}; \mathbf{B}, \mathbf{C})$  is a Bregman Divergence

#### 4.5.1 Overall sparsity: $\ell_{1,1}$ -norm

Here  $f(\mathbf{B}) = \|\mathbf{B}\|_{1,1}$ . The subproblem to be considered is

$$\mathbf{B}_{t+1} = \operatorname{argmin}_{\mathbf{B} \geq 0} \frac{1}{2} \|\mathbf{A} - \mathbf{B}\mathbf{C}_t\|_F^2 + \lambda \|\mathbf{B}\|_{1,1}. \quad (4.24)$$

We shall show that the following update rule

$$\mathbf{B}_{t+1}(i, j) \leftarrow \max\left(0, \mathbf{B}_t(i, j) \frac{[\mathbf{A}\mathbf{C}_t^T]_{ij} - \lambda}{[\mathbf{B}_t\mathbf{C}_t\mathbf{C}_t^T]_{ij}}\right), \quad (4.25)$$

ensures descent on the objective function in (4.24). Note that update (4.25) respects the non-negativity of  $\mathbf{B}$ .

Since (4.24) decomposes over rows of  $\mathbf{B}$ , it suffices to derive the updates for an arbitrary row  $\mathbf{b}$ , i.e., by ensuring descent for

$$\min_{\mathbf{b} \geq 0} \frac{1}{2} \|\mathbf{a} - \mathbf{b}\mathbf{C}\|_2^2 + \lambda \|\mathbf{b}\|_1. \quad (4.26)$$

The following proposition derives the update for (4.26).

**Proposition 6.** *The function  $h(\mathbf{b}) = \frac{1}{2} \|\mathbf{a} - \mathbf{b}\mathbf{C}\|_2^2 + \lambda \|\mathbf{b}\|_1$  is non-increasing under the update rule*

$$\mathbf{b}(p) \leftarrow \max\left(0, \mathbf{b}(p) \frac{[\mathbf{a}\mathbf{C}^T]_p - \lambda}{[\mathbf{b}\mathbf{C}\mathbf{C}^T]_p}\right), \quad p \in [k].$$

*Proof.* We prove the proposition by illustrating an appropriate auxiliary function technique for  $h$ . That is, first we construct a function  $g(\mathbf{b}, \mathbf{b}')$  such that  $h(\mathbf{b}) \leq g(\mathbf{b}, \mathbf{b}')$ , and  $h(\mathbf{b}) = g(\mathbf{b}, \mathbf{b})$ ; then, we show that the update is obtained by solving

$$\mathbf{b}^{t+1} \leftarrow \operatorname{argmin}_{\mathbf{b}} g(\mathbf{b}, \mathbf{b}^t). \quad (4.27)$$

Following the above steps guarantees descent, as  $h(\mathbf{b}^t) = g(\mathbf{b}^t, \mathbf{b}^t) \geq g(\mathbf{b}^{t+1}, \mathbf{b}^t) \geq h(\mathbf{b}^{t+1})$ .

Now consider,

$$\begin{aligned} h(\mathbf{b}) &= \frac{1}{2} \|\mathbf{a} - \mathbf{b}\mathbf{C}\|_2^2 + \lambda \|\mathbf{b}\|_1, \\ &\leq \frac{1}{2} \|\mathbf{a}\|_2^2 - \sum_i \mathbf{a}_i \sum_j \mathbf{b}_j \mathbf{c}_{ij} + \frac{1}{2} \sum_{ij} \frac{(\mathbf{b}_j \mathbf{c}_{ij})^2}{\lambda_{ij}} + \lambda \sum_j |\mathbf{b}_j| = g(\mathbf{b}, \mathbf{b}'), \end{aligned} \quad (4.28)$$

where  $\lambda_{ij} = \mathbf{b}'_j \mathbf{c}_{ij} / \mathbf{b}'_i \mathbf{c}_i$ ,  $\mathbf{c}_i$  is the  $i^{\text{th}}$  column of  $\mathbf{C}$  and  $\mathbf{c}_{ij}$  is its  $j^{\text{th}}$  term,  $\mathbf{b}_j$  is the  $j^{\text{th}}$  term of the vector  $\mathbf{b}$ . Observe that  $g(\mathbf{b}, \mathbf{b}) = h(\mathbf{b})$  holds too; thus  $g$  is a valid auxiliary function.

Now, differentiate  $g$  to obtain

$$\begin{aligned} \frac{\partial g}{\partial \mathbf{b}_p} &= -\sum_i \mathbf{a}_i \mathbf{c}_{ip} + \sum_{ij} \frac{\mathbf{b}_p \mathbf{c}_{ip} \mathbf{b}'_i \mathbf{c}_i}{\mathbf{b}'_p} + \lambda \frac{\mathbf{b}_p}{|\mathbf{b}_p|} \\ &= -[\mathbf{a}\mathbf{C}^T]_p + \frac{\mathbf{b}_p}{\mathbf{b}'_p} [\mathbf{b}'\mathbf{C}\mathbf{C}^T]_p + \lambda \frac{\mathbf{b}_p}{|\mathbf{b}_p|}. \end{aligned} \quad (4.29)$$

We must now solve  $\frac{\partial g}{\partial \mathbf{b}_p} = 0$ . Here we have the following two cases:

1.  $[\mathbf{a}\mathbf{C}^T]_p > \lambda$ : A positive root exists, and is given by  $\mathbf{b}_p = \mathbf{b}'_p \left( \frac{[\mathbf{a}\mathbf{C}^T]_p - \lambda}{[\mathbf{b}'\mathbf{C}\mathbf{C}^T]_p} \right)$ . No other roots exist.
2.  $[\mathbf{a}\mathbf{C}^T]_p < \lambda$ : No roots exist. However, for  $\mathbf{b}_p > 0$  we see that  $\frac{\partial g}{\partial \mathbf{b}_p} > 0$ . Thus, decreasing  $\mathbf{b}_p$  at any positive value will lead to a decrease in  $g$ . Also, for  $\mathbf{b}_p < 0$  we see that  $\frac{\partial g}{\partial \mathbf{b}_p} < 0$ . So, increasing  $\mathbf{b}_p$  at any of its negative values will lead to a decrease in  $g$ . So, setting  $\mathbf{b}_p = 0$  will lead to the smallest value of  $g$ .

Thus,  $\mathbf{b}(p) = \max\left(0, \mathbf{b}(p) \frac{[\mathbf{a}\mathbf{C}^T]_p - \lambda}{[\mathbf{b}\mathbf{C}\mathbf{C}^T]_p}\right)$  corresponds to (4.27), completing the proof.  $\square$

#### 4.5.2 Group-sparsity: $\ell_{1,2}$ -norm

For  $f(\mathbf{B}) = \|\mathbf{B}\|_{1,2}$ , we shall show that the following update

$$\mathbf{B}_{t+1}(i, j) = \mathbf{B}_t(i, j) \left( \frac{[\mathbf{A}\mathbf{C}_t^T]_{ij}}{[\mathbf{B}_t\mathbf{C}\mathbf{C}^T]_{ij} + \lambda[\mathbf{B}_t]_{ij}\|\mathbf{b}_i\|_2^{-1}} \right),$$

ensures descent. As before, we consider the matrix  $\mathbf{B}$  row-wise; the associated claim is formalized below.

**Proposition 7.** *The objective function  $h(\mathbf{b}) = \frac{1}{2}\|\mathbf{a} - \mathbf{b}\mathbf{C}\|_2^2 + \lambda\|\mathbf{b}\|_2$  is non-increasing under the update rule*

$$\mathbf{b}(p) = \mathbf{b}(p) \left( \frac{[\mathbf{a}\mathbf{C}^T]_p}{[\mathbf{b}\mathbf{C}\mathbf{C}^T]_p + \lambda\mathbf{b}(p)\|\mathbf{b}\|_2^{-1}} \right).$$

*Proof.* Consider the inequality

$$\begin{aligned} h(\mathbf{b}) &= \frac{1}{2}\|\mathbf{a} - \mathbf{b}\mathbf{C}\|_2^2 + \lambda\|\mathbf{b}\|_2 \\ &\leq \frac{1}{2}\|\mathbf{a}\|_2^2 - \sum_i \mathbf{a}_i \left( \sum_j \mathbf{b}_j \mathbf{c}_{ij} \right) + \frac{1}{2} \sum_{ij} \frac{(\mathbf{b}_j \mathbf{c}_{ij})^2}{\lambda_{ij}} + \lambda \sqrt{\sum_j \mathbf{b}_j^2}, \end{aligned} \quad (4.30)$$

where  $\lambda_{ij} = \frac{\mathbf{b}'_j \mathbf{c}_{ij}}{\mathbf{b}'_i \mathbf{c}_i}$  as before. We now bound  $\|\mathbf{b}\|_2$  by invoking the inequality  $\sqrt{x} \leq \frac{\sqrt{y}}{2} + \frac{x}{2\sqrt{y}}$ , whereby

$$f(\mathbf{b}) \leq \frac{1}{2}\|\mathbf{a}\|_2^2 - \sum_i \mathbf{a}_i \left( \sum_j \mathbf{b}_j \mathbf{c}_{ij} \right) + \frac{1}{2} \sum_{ij} \frac{(\mathbf{b}_j \mathbf{c}_{ij})^2}{\lambda_{ij}} + \left( \frac{1}{2}\|\mathbf{b}'\|_2 + \frac{1}{2\|\mathbf{b}'\|_2} \sum_j \mathbf{b}_j^2 \right) = g(\mathbf{b}, \mathbf{b}'). \quad (4.31)$$

Additionally, one can easily verify that  $g(\mathbf{b}, \mathbf{b}) = f(\mathbf{b})$ , thus establishing  $g$  to be a valid auxiliary function.

Finally, to obtain the update, we first compute the derivative

$$\begin{aligned} \frac{\partial g}{\partial \mathbf{b}_p} &= - \sum_i \mathbf{a}_i \mathbf{c}_{ip} + \sum_{ij} \frac{\mathbf{b}_p \mathbf{c}_{ip} \mathbf{b}'_i \mathbf{c}_i}{\mathbf{b}'_i} + \lambda \frac{\mathbf{b}_p}{\|\mathbf{b}'\|_2} \\ &= -[\mathbf{a}\mathbf{C}^T]_p + \frac{\mathbf{b}_p}{\mathbf{b}'_p} [\mathbf{b}'\mathbf{C}\mathbf{C}^T]_p + \lambda \frac{\mathbf{b}_p}{\|\mathbf{b}'\|_2}. \end{aligned}$$

Then, we solve  $\frac{\partial g}{\partial \mathbf{b}_p} = 0$ ; this solution immediately yields the desired update.  $\square$

## 5 Extension to Multifactor Matrix Approximation

We shall digress here from our main problem and see that the methods described for solving our problem can be applied to the more general setting of Multifactor Matrix Approximation, wherein a matrix  $\mathbf{A}$  may be required to be decomposed as a product of more than 2 matrices.

As an example, consider 3-factor decomposition with non-negativity constraints. Here we wish to decompose  $\mathbf{A}$  approximately into 3 matrices,  $\mathbf{R}$ ,  $\mathbf{B}$  and  $\mathbf{C}$ , with appropriately chosen dimensions, by solving

$$\min_{\mathbf{R} \geq 0, \mathbf{B} \geq 0, \mathbf{C} \geq 0} D(\mathbf{A}; \mathbf{R}, \mathbf{B}, \mathbf{C}). \quad (5.1)$$

Problem (5.1) can also be solved using the obvious Alternating Minimization approach:

$$\begin{aligned} \mathbf{R}_{t+1} &\leftarrow \operatorname{argmin}_{\mathbf{R} \geq 0} D(\mathbf{A}; \mathbf{R}, \mathbf{B}_t, \mathbf{C}_t), \\ \mathbf{B}_{t+1} &\leftarrow \operatorname{argmin}_{\mathbf{B} \geq 0} D(\mathbf{A}; \mathbf{R}_{t+1}, \mathbf{B}, \mathbf{C}_t), \\ \mathbf{C}_{t+1} &\leftarrow \operatorname{argmin}_{\mathbf{C} \geq 0} D(\mathbf{A}; \mathbf{R}_{t+1}, \mathbf{B}_{t+1}, \mathbf{C}). \end{aligned}$$

From this iteration, it is evident that any sparsity constraints imposed on  $\mathbf{R}$ ,  $\mathbf{B}$  or  $\mathbf{C}$ , may be easily incorporated by suitably modifying the corresponding minimization subproblem. Thus, our previous derivations extend to handle sparsity constrained multifactor approximation too.

## 6 Preliminary experiments

We took a random  $10 \times 10$  matrix  $A$ , with values in the range of  $0 - 100$ , and ran our Sparse Formulations (only **Formulation 1**) on this matrix, for a random initial choice of  $B$  and  $C$ . The following figures illustrate the descent in the objective function, for different formulations.

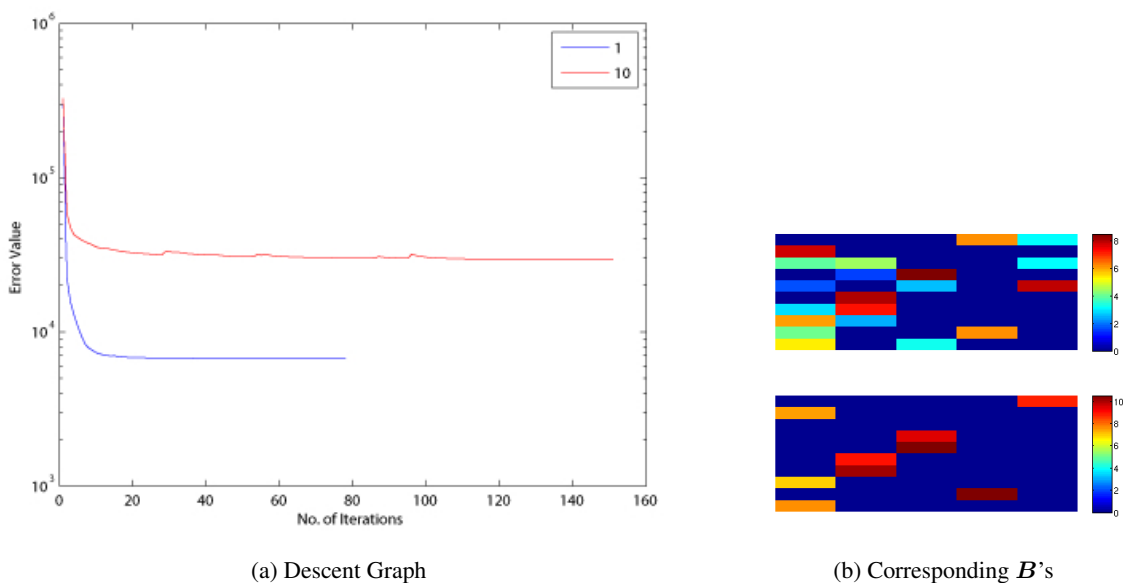


Figure 1:  $f(B) = \|B\|_{1,0}$ , with  $\lambda = 1$  and  $10$ . The number of non-zeros in the matrix  $B$  obtained finally were — 20 for  $\lambda = 1$  and 9 for  $\lambda = 10$ . The increase in error on increasing sparsity is evident here

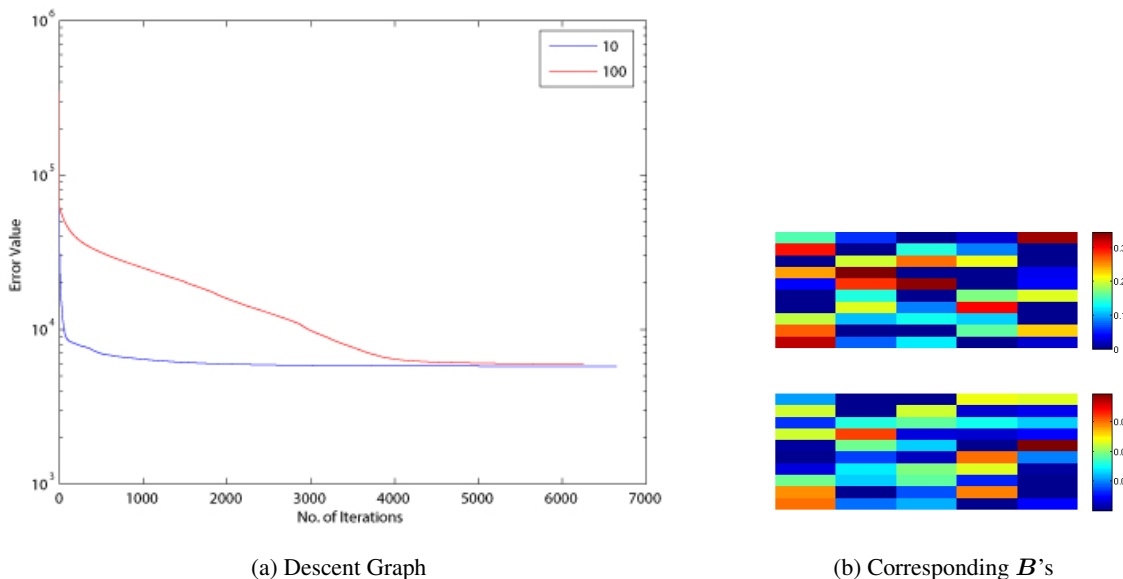
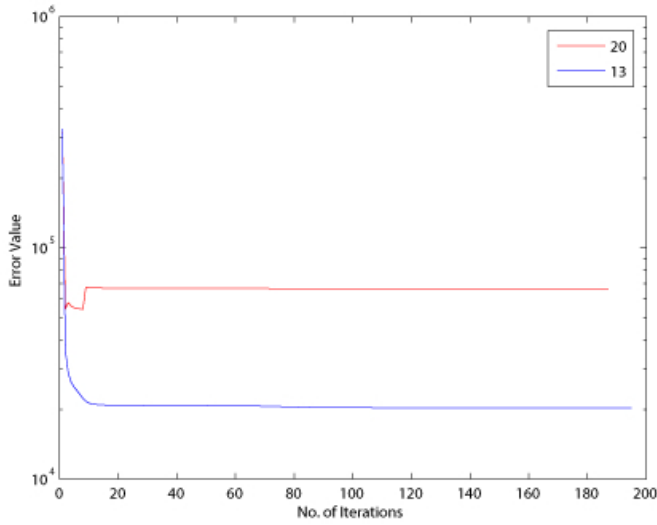
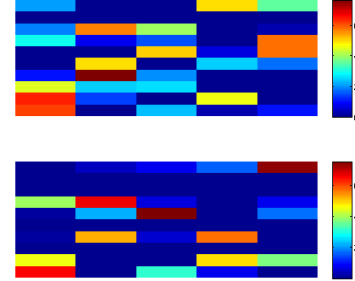


Figure 2:  $f(B) = \|B\|_{1,1}$ , with  $\lambda = 10$  and  $100$ . The number of non-zeros in the matrix  $B$  obtained finally were almost the same in both cases — 36 for  $\lambda = 10$  and 42 for  $\lambda = 100$ . The reason for this ambiguity is that since  $\|\frac{B}{\alpha}\|_{1,1} = \frac{1}{\alpha}\|B\|_{1,1}$ , the  $(1,1)$  norm can be made arbitrarily small by dividing  $B$  by an appropriate  $\alpha$ . The objective function can be kept the same by multiplying  $C$  with the same  $\alpha$ . So, increasing  $\lambda$  here seems to result in every element becoming comparatively smaller. One way around this could be to threshold the value of the elements to 0 below a particular point, and then rescale  $B$  and  $C$  in the end. Another possible way out could be to impose some reasonable sparsity constraints on  $C$  as well.

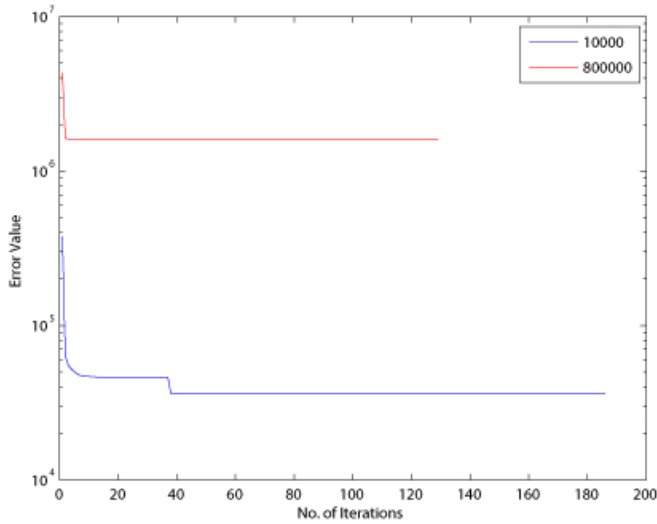


(a) Descent Graph

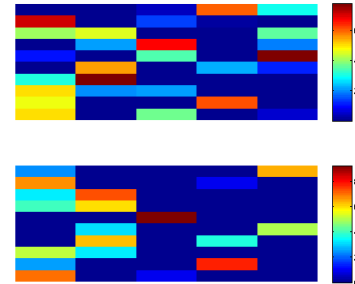


(b) Corresponding  $B$ 's

Figure 3:  $f(\mathbf{B}) = \|\mathbf{B}\|_{0,0}$ , with  $\lambda = 13$  and  $20$ . The number of zero rows in the matrix  $\mathbf{B}$  obtained finally was — 1 for  $\lambda = 13$  and 4 for  $\lambda = 20$ , with the total number of zero elements being 20 and 28 respectively.



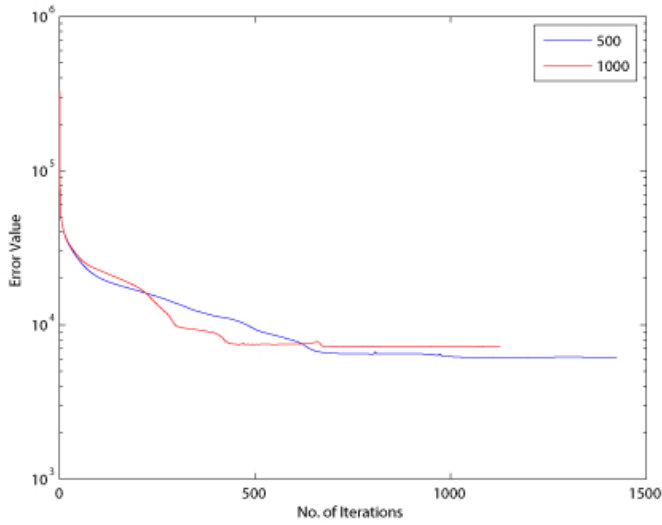
(a) Descent Graph



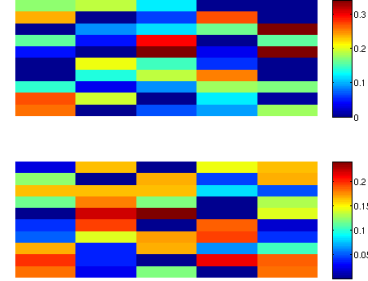
(b) Corresponding  $B$ 's

Figure 5:  $f(\mathbf{B}) = \|\mathbf{B}\|_{\infty,0}$ , with  $\lambda = 10000$  and  $800000$ . Every row of  $\mathbf{B}$  had exactly 2 zeros for  $\lambda = 10000$  and exactly 3 zeros for  $\lambda = 800000$ .



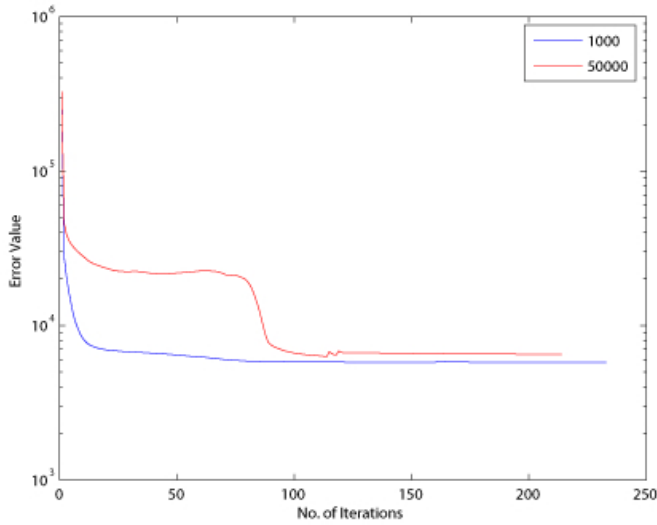


(a) Descent Graph

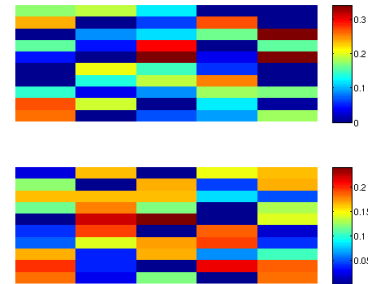


(b) Corresponding  $B$ 's

Figure 4:  $f(\mathbf{B}) = \|\mathbf{B}\|_{1,\infty}$ , with  $\lambda = 500$  and  $1000$ . No complete zero rows were obtained in the matrix  $\mathbf{B}$  here, but for  $\lambda = 1000$ , more number of rows had all their elements close to 0 than for  $\lambda = 500$ . Again, a thresholding scheme might be useful here since we have  $\|\frac{\mathbf{B}}{\alpha}\|_{1,\infty} = \frac{1}{\alpha}\|\mathbf{B}\|_{1,\infty}$



(a) Descent Graph



(b) Corresponding  $B$ 's

Figure 6:  $f(\mathbf{B}) = \|\mathbf{B}\|_{\infty,1}$ , with  $\lambda = 1000$  and  $50000$ . In  $\mathbf{B}$ , more number of elements per row were closer to 0 for  $\lambda = 50000$  than for  $\lambda = 1000$ . Imposing a thresholding scheme or sparsity constraints on  $\mathbf{C}$  will be meaningful here as well since  $\|\frac{\mathbf{B}}{\alpha}\|_{\infty,1} = \frac{1}{\alpha}\|\mathbf{B}\|_{\infty,1}$

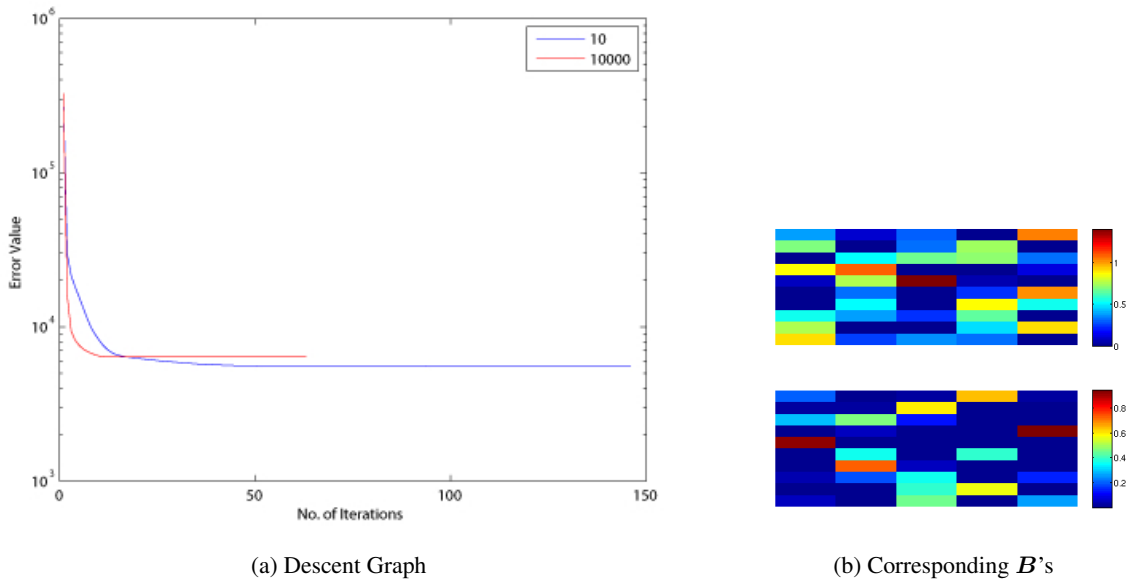


Figure 7:  $f(\mathbf{B}) = \frac{1}{4}\|\mathbf{B}^T\mathbf{B} - \mathbf{I}\|_F^2$ , with  $\lambda = 10$  and 10000.  $f(\mathbf{B})$  was 12.5941 for  $\lambda = 10$  and 0.1016 for  $\lambda = 10000$ . Qualitatively as well,  $\mathbf{B}$  obtained for  $\lambda = 10000$  seemed more diverse than that obtained for  $\lambda = 10$ .

## 6.1 On a Real Dataset

We have tested out our Sparse NMA formulations on the ORL face image datasets (<http://www.uk.research.att.com/facedatabase.html>), also used by Hoyer for testing sparsity on NMF[10]. The preliminary results are described below. However, more extensive experiments, requiring a careful tuning of the parameters involved, is still to be done.

## 7 Discussion

### References

- [1] A. Argyriou, T. Evgeniou, and M. Pontil. Convex multi-task feature learning. *Mach. Learn.*, 73(3):243–272, 2008. 3
- [2] M. Berry, M. Browne, A. Langville, P. Pauca, and R. J. Plemmons. Algorithms and applications for approximation nonnegative matrix factorization. *Computational Statistics and Data Analysis*, 52:155–173, 2007. 1
- [3] A. Cichocki, R. Zdunek, A. H. Phan, and S. ichi Amari. *Nonnegative matrix and tensor factorizations: applications to exploratory multi-way data analysis and blind source separation*. John Wiley & Sons, Ltd., 2009. 2
- [4] P. L. Combettes and J.-C. Pesquet. Proximal splitting methods in signal processing. *arXiv*, (0912.3522), 2010. 5, 6
- [5] I. S. Dhillon and S. Sra. Generalized Nonnegative Matrix Approximations with Bregman Divergences. In *NIPS 18*, Vancouver, Canada, 2006. 2
- [6] D. L. Donoho. Compressed sensing. *IEEE Tran. Inf. Theory*, 52(4):1289–1306, 2006. 1
- [7] L. Grippo and M. Sciandrone. On The Convergence of The Block Nonlinear Gauss-Seidel Method under Convex Constraints. *Operations Research Letters*, 26:127–136, 2000. 5
- [8] S. Harmeling, M. Hirsch, S. Sra, and B. Schölkopf. Online Blind Deconvolution for Astronomy. In *IEEE Int. Conf. on Computational Photography*, Apr. 2009. 2
- [9] M. Heiler and C. Schnörr. Learning Sparse Representations by Non-Negative Matrix Factorization and Sequential Cone Programming. *JMLR*, 7:1385–1407, 2006. 1

- [10] P. O. Hoyer. Non-negative matrix factorization with sparseness constraints. *Journal of Machine Learning Research*, 5:1457–1469, 2004. 1, 16
- [11] D. Kim, S. Sra, and I. S. Dhillon. A scalable trust-region algorithm with application to mixed-norm regression. In *Int. Conf. Machine Learning (ICML)*, 2010. 1
- [12] H. Kim and H. Park. Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis. *Bioinformatics*, 23(12):1495–1502, 2007. 1
- [13] K. Kiwiel. On linear-time algorithms for the continuous quadratic knapsack problem. *Journal of Optimization Theory and Applications*, 134:549–554, 2007. 8
- [14] D. D. Lee and H. S. Seung. Algorithms for Non-negative Matrix Factorization. In *NIPS*, pages 556–562, 2000. 10
- [15] H. Liu, M. Palatucci, and J. Zhang. Blockwise Coordinate Descent Procedures for the Multi-task Lasso, with Applications to Neural Semantic Basis Discovery. In *ICML*, Jun. 2009. 1
- [16] N. Maculan and G. G. de Paula. A linear-time median-finding algorithm for projecting a vector on the simplex of  $\mathbb{R}^n$ . *Operations Research Letters*, 8(4):219 – 222, 1989. 8
- [17] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online Learning for Matrix Factorization and Sparse Coding. *JMLR*, 11:10–60, 2010. 1, 2
- [18] M. Patriksson. A survey on a classic core problem in operations research. Technical Report 2005:33, Chalmers University of Technology and Göteborg University, Oct. 2005. 8
- [19] A. Quattoni, X. Carreras, M. Collins, and T. Darrell. An efficient projection for  $\ell_{1,\infty}$  regularization. In *ICML*, page 108, 2009. 9
- [20] I. Rish and G. Grabarnik. Sparse modeling: ICML 2010 tutorial. Online, Jun. 2010. 1
- [21] F. Sinz, E. Simoncelli, and M. Bethge. Hierarchical modeling of local image features through lp-nested symmetric distributions. In *NIPS*, 2009. 3
- [22] S. Sra. *Matrix Nearness Problems in Data Mining*. PhD thesis, Univ. of Texas at Austin, Aug. 2007. 1, 2
- [23] S. Sra. Generalized proximity and projection with norms and mixed-norms. Technical Report 192, Max Planck Institute for Biological Cybernetics, May 2010. 9
- [24] R. Tibshirani. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996. 1
- [25] J. A. Tropp. *Topics in Sparse Approximation*. PhD thesis, Univ. of Texas at Austin, 2004. 2
- [26] R. University. Compressive sensing resources. <http://dsp.rice.edu/cs>, 2010. 1
- [27] E. van den Berg, M. Schmidt, M. Friedlander, and K. Murphy. Group Sparsity via Linear-Time Projection. *Technical Report, Department of Computer Science, University of British Columbia*, 2008. 9
- [28] S. A. Vavasis. On the complexity of nonnegative matrix factorization. *SIAM Journal on Optimization*, 20(3): 1364–1377, 2009. 5
- [29] P. Zhao, G. Rocha, and B. Yu. The composite absolute penalties family for grouped and hierarchical variable selection. *Annals of Statistics*, 37(6A):3468–3497, 2009. 1, 3