

Finding approximately rank-one submatrices with the nuclear norm and ℓ_1 -norm*

Xuan Vinh Doan[†] Stephen Vavasis[‡]

November 2010

Abstract

We propose a convex optimization formulation with the nuclear norm and ℓ_1 -norm to find a large approximately rank-one submatrix of a given nonnegative matrix. We develop optimality conditions for the formulation and characterize the properties of the optimal solutions. We establish conditions under which the optimal solution of the convex formulation has a specific sparse structure. Finally, we show that, under certain hypotheses, with high probability, the approach can recover the rank-one submatrix even when it is corrupted with random noise and inserted as a submatrix into a much larger random noise matrix.

1 Introduction

Given a nonnegative matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $a_{ij} \geq 0$ for all $i = 1, \dots, m$, $j = 1, \dots, n$, we consider the problem of finding $\mathcal{I} \subset \{1, \dots, m\}$ and $\mathcal{J} \subset \{1, \dots, n\}$ such that $\mathbf{A}(\mathcal{I}, \mathcal{J})$ is close to a rank-one matrix, and such that $\|\mathbf{A}(\mathcal{I}, \mathcal{J})\|$ is large. We shall call this problem the LAROS problem (for “large approximately rank-one submatrix”).

The main application of the LAROS problem is for finding features in data. For example, suppose \mathbf{A} represents a corpus of documents in some language. Each column of \mathbf{A} is in correspondence with

*Supported in part by the U. S. Air Force Office of Scientific Research, a Discovery Grant from the Natural Sciences and Engineering Research Council of Canada, and a grant from MITACS.

[†]Department of Combinatorics and Optimization, University of Waterloo, 200 University Avenue West, Waterloo, ON N2L 3G1, Canada, vanxuan@uwaterloo.ca.

[‡]Department of Combinatorics and Optimization, University of Waterloo, 200 University Avenue West, Waterloo, ON N2L 3G1, Canada, vavasis@uwaterloo.ca.

one document, and each row is in correspondence with a term used in the corpus. Here, “term” means a word in the language, excluding common words such as articles and prepositions. The (i, j) entry of \mathbf{A} is the number of occurrences of term i in document j , perhaps normalized. Such a matrix is called the *term-document matrix* of the underlying corpus.

In this case, an approximately rank-one submatrix of \mathbf{A} corresponds to a subset of terms and a subset of documents in which the selected terms occur with proportional frequencies in the selected documents. Such a submatrix may correspond to the intuitive notion of a topic that recurs in several documents, since a topic may manifest itself as a particular group of relevant terms that occur roughly in the same proportions.

As another example, the matrix \mathbf{A} may correspond to a database of pixelated grayscale images, where each image has the same pixel size. Each column of \mathbf{A} corresponds to one image, and each row to one pixel position. The (i, j) entry of \mathbf{A} is the intensity of the i th pixel in the j th image. In this case, the approximately rank-one submatrix corresponds to a visual feature that recurs in a certain position in some subset of the images.

If one wanted to find more than one topic in a term-document matrix or more than one feature in an image database matrix, then one could iteratively find an approximately rank-one submatrix, subtract it from \mathbf{A} (perhaps modifying the result of the subtraction to ensure that \mathbf{A} remains nonnegative), and then repeat the procedure p times. Let the submatrices discovered be denoted $(\mathcal{I}_1, \mathcal{J}_1), \dots, (\mathcal{I}_p, \mathcal{J}_p)$. Suppose $A(\mathcal{I}_i, \mathcal{J}_i) \approx \mathbf{w}_i \mathbf{h}_i^T$ for $i = 1, \dots, p$, and let $\bar{\mathbf{w}}_i, \bar{\mathbf{h}}_i$ denote the extension of $\mathbf{w}_i, \mathbf{h}_i$ to vectors of length m, n by inserting zeros for entries not in $\mathcal{I}_i, \mathcal{J}_i$ respectively.

It is known that if $\hat{\mathbf{A}}$ is a nonnegative matrix representing a submatrix of \mathbf{A} , then the minimizer $\mathbf{w} \mathbf{h}^T$ of $\|\hat{\mathbf{A}} - \mathbf{w} \mathbf{h}^T\|$ is the dominant singular vector pair in either the Frobenius or 2-norm (a consequence of the Eckart–Young theorem, Theorem 2.5.3 of [11]) and furthermore, $\mathbf{w} \geq \mathbf{0}$ and $\mathbf{h} \geq \mathbf{0}$ (a consequence of the Perron–Frobenius theorem.) Thus, without loss of generality, we may assume that each $\mathbf{w}_i \mathbf{h}_i^T$ determined by the iterative computation is nonnegative.

In this case, one has an approximate factorization

$$\mathbf{A} \approx [\bar{\mathbf{w}}_1, \dots, \bar{\mathbf{w}}_p][\bar{\mathbf{h}}_1, \dots, \bar{\mathbf{h}}_p]^T,$$

where we can write the right-hand side as $\mathbf{W} \mathbf{H}^T$ with $\mathbf{W} \geq \mathbf{0}, \mathbf{H} \geq \mathbf{0}$. This factorization is called a *nonnegative matrix factorization* of \mathbf{A} . The earliest reference known to us concerning nonnegative matrix factorization is Thomas’ solution [18] to a problem posed by A. Berman and R. Plemmons (which, according to a remark in the journal, was also solved by A. Ben-Israel). Cohen and Rothblum

[8] describe applications for NMF in probability, quantum mechanics and other fields. Lee and Seung [13] showed that NMF can find features in image databases, and Hofmann [12] showed that probabilistic latent semantic analysis, a variant of NMF, can effectively cluster documents according to their topics.

Nonnegative matrix factorization is sometimes posed as an optimization problem: find $\mathbf{W} \in \mathbb{R}^{m \times p}$ and $\mathbf{H} \in \mathbb{R}^{n \times p}$, both nonnegative, such that $\|\mathbf{A} - \mathbf{W}\mathbf{H}^T\|$ is minimized in some matrix norm. It is known that this optimization problem is NP-hard [19]. Therefore, it is not surprising that most algorithms for the problem are heuristic in the sense that they do not make guarantees about the quality of the approximation.

One class of heuristic NMF algorithms are the ‘greedy’ algorithms [2, 3, 4, 10] that follow the framework described above. In a greedy algorithm, the columns of \mathbf{W} and \mathbf{H} are generated sequentially, with each new pair of columns accounting for one feature in the original \mathbf{A} . These greedy algorithms give rise to the LAROS subproblem addressed in this paper, namely, find one pair $\mathbf{w}_i, \mathbf{h}_i$ nonnegative such that $\mathbf{w}_i \mathbf{h}_i^T$ is a good approximation for a submatrix of \mathbf{A} in the positions (i, j) where \mathbf{A} is positive.

The LAROS subproblem, however, is itself NP-hard as observed by [10]. This is because the maximum-edge biclique problem can be naturally expressed as a rank-one submatrix problem. The biclique problem takes as input a bipartite graph $G = (U, V, E)$. The output is composed of two subsets $U^* \subset U$ and $V^* \subset V$ such that $U^* \times V^* \subset E$ (i.e., all possible $|U^*| \cdot |V^*|$ edges between U^* and V^* are present in G) and such that $|U^*| \cdot |V^*|$ is maximum with this property. This problem was shown by Peeters [15] to be NP-hard.

Maximum-edge biclique can be expressed as finding a large rank-one submatrix using the following construction. Let \mathbf{A} be a $|U| \times |V|$ matrix with rows in correspondence to U and columns in correspondence to V . Entry (i, j) of \mathbf{A} for $(i, j) \in U \times V$ is 1 if $(i, j) \in E$, else this entry is 0. Then a biclique corresponds exactly of a $|U^*| \cdot |V^*|$ submatrix of all 1’s. A submatrix of all 1’s is a rank-one matrix of norm $(|U^*| \cdot |V^*|)^{1/2}$, and there is no other kind of rank-one submatrix of \mathbf{A} .

We will provide a formal definition of the LAROS problem, i.e., exactly what is the desired output in Section 2. (Some authors mentioned earlier, e.g., [4] and [10] have provided other formal definitions.) We will also propose an convex optimization problem in Section 2 that, for matrices \mathbf{A} constructed in a certain way, successfully finds large, approximately rank-one submatrices. We present two such theorems. One case is when the approximately rank-one submatrix dominates the rest of the matrix; this is presented in Section 3.

The second case is when \mathbf{A} is constructed as follows:

$$\mathbf{A} = \mathbf{A}_0 + \mathbf{R},$$

where there are two index sets \mathcal{I}, \mathcal{J} such that $\text{rank}(\mathbf{A}_0(\mathcal{I}, \mathcal{J})) = 1$, $\mathbf{A}_0(i, j) = 0$ for all $(i, j) \notin \mathcal{I} \times \mathcal{J}$, and \mathbf{R} is a random matrix representing noise. In this case, under certain assumptions, our algorithm recovers $(\mathcal{I}, \mathcal{J})$ from \mathbf{A} as proved in Section 4.

2 Matrix norm minimization

For reasons that will become clear, we start this section by presenting the convex relaxation of the LAROS problem, and only later will we present a nonconvex exact optimization formulation. In particular, we propose the following convex optimization problem that in some cases solves the LAROS problem:

$$\begin{aligned} \min \quad & \|\mathbf{X}\|_* + \theta \|\mathbf{X}\|_1 \\ \text{s.t.} \quad & \langle \mathbf{A}, \mathbf{X} \rangle \geq 1. \end{aligned} \tag{1}$$

The matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$ is the unknown. Norm $\|\mathbf{X}\|_*$ is the *nuclear norm*, also called the *trace norm*; it is the sum of the singular values of \mathbf{X} . We use the notation $\|\mathbf{X}\|_1$ to mean the sum of the absolute values of entries of \mathbf{X} , that is, the ℓ_1^{mn} -norm applied to $\text{vec}(\mathbf{X})$, the concatenation of the columns of \mathbf{X} into a long vector. Finally, $\langle \mathbf{A}, \mathbf{X} \rangle$ means the inner product of the two matrices. We note that an objective function involving a sum of the nuclear and ℓ_1 -norms was used for a different purpose by [7]. Two other norms used extensively in this paper are $\|\mathbf{X}\|$, which is the spectral or 2-norm, i.e., $\sigma_1(\mathbf{X})$, and $\|\mathbf{X}\|_\infty$, which is the ℓ_∞^{mn} -norm applied to $\text{vec}(\mathbf{X})$, i.e., the maximum absolute entry of \mathbf{X} .

Before beginning a detailed analysis of this optimization problem, we first provide some motivation. Consider first the simplification obtained by taking $\theta = 0$:

$$\begin{aligned} \min \quad & \|\mathbf{X}\|_* \\ \text{s.t.} \quad & \langle \mathbf{A}, \mathbf{X} \rangle \geq 1. \end{aligned}$$

It follows from Proposition 1 below that the optimal solution is found using the singular value decomposition. In particular, if \mathbf{A} is factored as $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, where $\mathbf{U} \in \mathbb{R}^{m \times m}$ is orthogonal, $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$ is diagonal, and $\mathbf{V} \in \mathbb{R}^{n \times n}$ is orthogonal, then an optimizer is $\mathbf{X} = \mathbf{U}(:, 1)\mathbf{V}(:, 1)^T/\sigma_1$. Thus, when $\theta = 0$, the above formulation successfully finds the best rank-one approximation to the whole matrix \mathbf{A} .

This approximation, however, is not always well suited for identifying submatrices. Consider e.g.

the following 6×6 matrix:

$$\mathbf{A} = \begin{pmatrix} 0.8 & 0.9 & 1.1 & 0.1 & 0.2 & 0.2 \\ 0.8 & 1.1 & 0.8 & 0 & 0 & 0 \\ 1.0 & 1.0 & 0.8 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.8 & 0.9 & 1.0 \\ 0 & 0 & 0 & 0.9 & 1.0 & 0.8 \\ 0 & 0 & 0 & 1.0 & 1.1 & 0.8 \end{pmatrix}$$

It is apparent from inspection that this matrix has two 3×3 approximately rank-one blocks in positions $\{1, 2, 3\} \times \{1, 2, 3\}$ and $\{4, 5, 6\} \times \{4, 5, 6\}$. If the ‘noise’ entries in the upper right $\{1, 2, 3\} \times \{4, 5, 6\}$ block were absent, then the two dominant singular vectors would exactly identify the two diagonal blocks. Once the noise entries are inserted, however, the dominant left singular vector of the above matrix \mathbf{A} accurate to two decimal places is $[.45, .37, .37, .40, .40, .43]^T$. In other words, there is no separation at all between the rows numbered 1, 2, 3 and those numbered 4, 5, 6, so no submatrix is identified.

Armed with a preliminary understanding of the convex relaxed formulation, we now present and motivate an exact (nonconvex) formulation of LAROS, which is as follows.

$$\begin{aligned} \min \quad & \|\mathbf{X}\|_* + \theta |\mathcal{I}| |\mathcal{J}| \\ \text{s.t.} \quad & \langle \mathbf{A}, \mathbf{X} \rangle \geq 1, \\ & x_{ij} = 0, \quad \forall (i, j) \notin \mathcal{I} \times \mathcal{J}, \end{aligned}$$

where $\mathcal{I} \subset \{1, \dots, m\}$ and $\mathcal{J} \subset \{1, \dots, n\}$ are unknowns (as well as \mathbf{X}). For fixed \mathcal{I} and \mathcal{J} , the optimal solution would be the rank-one approximation of the submatrix $\mathbf{A}(\mathcal{I}, \mathcal{J})$ given by the SVD, as explained above, and the optimal value is $\|\mathbf{A}(\mathcal{I}, \mathcal{J})\|^{-1} + \theta |\mathcal{I}| |\mathcal{J}|$. The first term of the optimal value is $\|\mathbf{A}(\mathcal{I}, \mathcal{J})\|^{-1}$; therefore, for appropriate selection of θ , a large submatrix (in terms of 2-norm) will be selected. The second term is the size of the submatrix, which is a nonconvex function. Thus, the two terms balance the twin objectives of selecting a submatrix with a large first singular value and selecting a submatrix that has a relatively small number of entries.

This now motivates (1): we relax the above nonconvex formulation by replacing the cardinality term in the objective function with the ℓ_1 -norm. The relaxed term $\theta \|\mathbf{X}\|_1$ in the objective function has the well-known effect of favoring sparser matrices \mathbf{X} (those with fewer nonzero entries). Thus, the combination of the two terms seeks a low rank matrix with many entries equal to 0. For example, our formulation (1) applied to \mathbf{A} above identifies the $\{4, 5, 6\} \times \{4, 5, 6\}$ submatrix when $\theta = 0.5$. In particular, the solution \mathbf{X} has zeros in all positions except $\{4, 5, 6\} \times \{4, 5, 6\}$; in these positions it has positive entries ranging from 0.08 to 0.16.

To start our more formal analysis, let us consider the following general norm minimization problem.

$$\begin{aligned} \min \quad & \|\mathbf{X}\| \\ \text{s.t.} \quad & \langle \mathbf{A}, \mathbf{X} \rangle \geq 1, \end{aligned} \tag{2}$$

where $\|\cdot\|$ is an arbitrary norm function on $\mathbb{R}^{m \times n}$. (For example, the objective function $\|\mathbf{X}\|_* + \theta \|\mathbf{X}\|_1$ appearing in (1) is a norm.) Using the associated dual norm $\|\cdot\|^*$, we can relate Problem (2) to an equivalent problem as follows.

Lemma 1. *Consider $\mathbf{A} \neq \mathbf{0}$. Matrix \mathbf{X}^* is an optimal solution of Problem (2) if and only if $\mathbf{Y}^* = \|\mathbf{A}\|^* \mathbf{X}^*$ is an optimal solution of the following problem:*

$$\begin{aligned} \max \quad & \langle \mathbf{A}, \mathbf{Y} \rangle \\ \text{s.t.} \quad & \|\mathbf{Y}\| \leq 1. \end{aligned} \tag{3}$$

Proof. Let \mathbf{X}^* be an optimal solution of Problem (2). Clearly, $\mathbf{X}^* \neq \mathbf{0}$ and $\langle \mathbf{A}, \mathbf{X}^* \rangle = 1$. Apply the norm inequality, we have

$$\|\mathbf{X}^*\| \cdot \|\mathbf{A}\|^* \geq \langle \mathbf{A}, \mathbf{X}^* \rangle = 1 \Leftrightarrow \|\mathbf{X}^*\| \geq \frac{1}{\|\mathbf{A}\|^*}.$$

According to Boyd and Vandenberghe [5], the dual of the dual norm is the original norm and the norm inequality is tight: for any \mathbf{A} , there is always an $\mathbf{X} \neq \mathbf{0}$ such that the equality holds (for finite-dimensional vector spaces). Since \mathbf{X}^* is an optimal solution of Problem (2),

$$\|\mathbf{X}^*\| = \frac{1}{\|\mathbf{A}\|^*}.$$

Let $\mathbf{Y}^* = \|\mathbf{A}\|^* \mathbf{X}^*$, we then have: $\|\mathbf{Y}^*\| = 1$ and $\langle \mathbf{A}, \mathbf{Y}^* \rangle = \|\mathbf{A}\|^*$. We also have:

$$\begin{aligned} \|\mathbf{A}\|^* &= \max \langle \mathbf{A}, \mathbf{Y} \rangle \\ \text{s.t.} \quad & \|\mathbf{Y}\| \leq 1. \end{aligned}$$

Thus \mathbf{Y}^* is indeed an optimal solution of Problem (3). Using similar arguments, we can prove that conversely, if \mathbf{Y}^* is an optimal solution of Problem (3), then $\mathbf{X}^* = (\|\mathbf{A}\|^*)^{-1} \mathbf{Y}^*$ is an optimal solution of Problem (2). \square

The next lemma characterizes the set of all optimal solutions of Problem (3).

Lemma 2. *The set of all optimal solutions of Problem (3) with $\mathbf{A} \neq \mathbf{0}$ is the subgradient of the dual norm function $\|\cdot\|^*$ at \mathbf{A} , $\partial \|\mathbf{A}\|^*$.*

Proof. Let \mathbf{Y}^* be an optimal solution of Problem (3), we have $\|\mathbf{Y}^*\| = 1$ since $\mathbf{A} \neq \mathbf{0}$. Thus we have: $\|\mathbf{A}\|^* = \langle \mathbf{A}, \mathbf{Y}^* \rangle$. For an arbitrary matrix $\mathbf{B} \in \mathbb{R}^{m \times n}$,

$$\|\mathbf{A} + \mathbf{B}\|^* \geq \langle \mathbf{A} + \mathbf{B}, \mathbf{Y}^* \rangle = \|\mathbf{A}\|^* + \langle \mathbf{B}, \mathbf{Y}^* \rangle.$$

Thus $\mathbf{Y}^* \in \partial\|\mathbf{A}\|^*$.

Now consider $\mathbf{Y} \in \partial\|\mathbf{A}\|^*$:

$$\|\mathbf{B}\|^* \geq \|\mathbf{A}\|^* + \langle \mathbf{B} - \mathbf{A}, \mathbf{Y} \rangle \Leftrightarrow \langle \mathbf{A}, \mathbf{Y} \rangle - \|\mathbf{A}\|^* \geq \langle \mathbf{B}, \mathbf{Y} \rangle - \|\mathbf{B}\|^*, \quad \forall \mathbf{B} \in \mathbb{R}^{m \times n}.$$

With $\mathbf{B} = \mathbf{0}$ and $\mathbf{B} = 2\mathbf{A}$, we obtain the equality $\langle \mathbf{A}, \mathbf{Y} \rangle = \|\mathbf{A}\|^* > 0$. We have:

$$\langle \mathbf{A}, \mathbf{Y} \rangle = \|\mathbf{A}\|^* \leq \|\mathbf{A}\|^* \|\mathbf{Y}\| \Leftrightarrow (\|\mathbf{Y}\| - 1) \|\mathbf{A}\|^* \geq 0 \Rightarrow \|\mathbf{Y}\|_* \geq 1.$$

In addition, $\langle \mathbf{B}, \mathbf{Y} \rangle - \|\mathbf{B}\|^* \leq 0$ for all $\mathbf{B} \in \mathbb{R}^{m \times n}$. The norm inequality $\langle \mathbf{B}, \mathbf{Y} \rangle \leq \|\mathbf{B}\|^* \|\mathbf{Y}\|$ is tight and $\mathbf{Y} \neq \mathbf{0}$ ($\|\mathbf{Y}\| \geq 1$); therefore, there exists $\mathbf{B} \neq \mathbf{0}$ such that $\langle \mathbf{B}, \mathbf{Y} \rangle = \|\mathbf{B}\|^* \|\mathbf{Y}\|$. Thus we have:

$$\|\mathbf{B}\|^* \|\mathbf{Y}\| - \|\mathbf{B}\|^* \leq 0 \Leftrightarrow (\|\mathbf{Y}\| - 1) \|\mathbf{B}\|^* \leq 0 \Rightarrow \|\mathbf{Y}\| \leq 1.$$

Thus $\|\mathbf{Y}\| = 1$ and $\langle \mathbf{A}, \mathbf{Y} \rangle = \|\mathbf{A}\|^*$, the optimal value of Problem (3), which means \mathbf{Y} is an optimal solution of Problem (3). \square

Lemma 1 and 2 show that the set of all optimal solutions of Problem (2) is $(\|\mathbf{A}\|^*)^{-1} \partial\|\mathbf{A}\|^*$. The uniqueness of the optimal solution of Problem (2) is equivalent to the differentiability of the dual norm function $\|\cdot\|^*$ at \mathbf{A} . These results are summarized in the following theorem.

Theorem 1. *Consider $\mathbf{A} \neq \mathbf{0}$. The following statements are true:*

(i) *The set of optimal solutions of Problem (2) is $(\|\mathbf{A}\|^*)^{-1} \partial\|\mathbf{A}\|^*$.*

(ii) *Problem (2) has a unique optimal solution if and only if the dual norm function $\|\cdot\|^*$ is differentiable at \mathbf{A} .*

If the norm is set to be the nuclear norm, we obtain the following minimization problem, which has been used [9, 16, 6, 1] as a relaxation of rank minimization optimization problems:

$$\begin{aligned} \min \quad & \|\mathbf{X}\|_* \\ \text{s.t.} \quad & \langle \mathbf{A}, \mathbf{X} \rangle \geq 1. \end{aligned} \tag{4}$$

The dual norm of the nuclear norm is the spectral norm. According to Ziętak [20], if $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ is a singular value decomposition of \mathbf{A} and s is the multiplicity of the largest singular value of \mathbf{A} , the subgradient $\partial \|\mathbf{A}\|$ is written as follows:

$$\partial \|\mathbf{A}\| = \left\{ \mathbf{U} \begin{bmatrix} \mathbf{S} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}^T : \mathbf{S} \in \mathcal{S}_+^s, \|\mathbf{S}\|_* = 1 \right\}.$$

Clearly, the largest rank-one approximation of \mathbf{A} , $\mathbf{u}_1\mathbf{v}_1^T$, always belongs to the subgradient $\partial \|\mathbf{A}\|$. The description of the subgradient $\partial \|\mathbf{A}\|$ shows that the maximum possible rank of an optimal solution of Problem (4) is the multiplicity of the largest singular value of \mathbf{A} . In addition, the spectral norm function $\|\cdot\|$ is not differentiable in general. The uniqueness of the optimal solution of Problem (4) is equivalent to the differentiability of the spectral norm function $\|\cdot\|$ at \mathbf{A} . The necessary and sufficient condition is $s = 1$ or equivalently, $\sigma_1(\mathbf{A}) > \sigma_2(\mathbf{A})$. In the case of unique optimal solution, we obtain the largest rank-one approximation of \mathbf{A} (up to the scaling factor $\|\mathbf{A}\|^{-1}$). These results are stated in the following proposition:

Proposition 1. *Consider $\mathbf{A} \neq \mathbf{0}$. The following statements are true:*

- (i) *The set of optimal solutions of Problem (4) is $\|\mathbf{A}\|^{-1} \partial \|\mathbf{A}\|$.*
- (ii) *The largest rank-one approximation of \mathbf{A} is an optimal solution of Problem (4) and it is the unique solution if and only if $\sigma_1(\mathbf{A}) > \sigma_2(\mathbf{A})$.*

Similar to low-rank minimization problems with nuclear norm approximation, sparse optimization problems can be approximately handled by the (vector) ℓ_1 -norm function $\|\cdot\|_1$. Let us consider the following problem

$$\begin{aligned} \min \quad & \|\mathbf{X}\|_1 \\ \text{s.t.} \quad & \langle \mathbf{A}, \mathbf{X} \rangle \geq 1. \end{aligned} \tag{5}$$

The dual norm of ℓ_1 -norm is the (vector) infinity norm $\|\cdot\|_\infty$, i.e., the maximum absolute entry of the matrix, and the subgradient $\partial \|\mathbf{A}\|_\infty$ can be written as follows,

$$\partial \|\mathbf{A}\|_\infty = \text{conv} \left\{ \text{sgn}(a_{ij}) \mathbf{E}_{ij} \mid (i, j) \in \arg \max_{(k,l)} |a_{kl}| \right\},$$

where \mathbf{E}_{ij} is the unit matrix in $\mathbb{R}^{m \times n}$ with $\mathbf{E}_{ij}(i, j) = 1$. The sparsity of the optimal solution \mathbf{X}^* of Problem (5) is clearly related to the multiplicity of the maximum absolute value of elements of \mathbf{A} . Applying Theorem 1 for this particular ℓ_1 -norm, we obtain the following results:

Proposition 2. Consider $\mathbf{A} \neq \mathbf{0}$. The following statements are true:

(i) The set of optimal solutions of Problem (5) is $\|\mathbf{A}\|_\infty^{-1} \partial \|\mathbf{A}\|_\infty$.

(ii) The matrix $\text{sgn}(a_{ij})\mathbf{E}_{ij}$, where $(i, j) \in \arg \max_{(k,l)} |a_{kl}|$, and $\text{sgn}(\cdot)$ is the usual sign function, is an optimal solution of Problem (5) and it is the unique solution if and only if $|a_{ij}| > |a_{kl}|$ for all $(k, l) \neq (i, j)$.

As mentioned above, finding a low-rank submatrix clearly involves both low-rank and sparse optimization (with a specific sparse structure). Let us return to the parametric optimization problem (1) proposed at the beginning of this section

$$\begin{aligned} \min \quad & \|\mathbf{X}\|_* + \theta \|\mathbf{X}\|_1 \\ \text{s.t.} \quad & \langle \mathbf{A}, \mathbf{X} \rangle \geq 1, \end{aligned}$$

where $\theta \geq 0$. Clearly, if $\theta = 0$, we obtain Problem (4) and when $\theta \rightarrow \infty$, we approach Problem (5). This optimization problem clearly addresses both low-rank and sparse requirements of the solution \mathbf{X} . We now would like to characterize the set of optimal solutions of the problem.

The objective function $\|\mathbf{X}\|_* + \theta \|\mathbf{X}\|_1$ is a norm function since $\theta \geq 0$. Denote $\|\mathbf{X}\|_\theta$ to be this parametric norm of \mathbf{X} ,

$$\|\mathbf{X}\|_\theta := \|\mathbf{X}\|_* + \theta \|\mathbf{X}\|_1$$

and consider its dual norm function $\|\cdot\|_\theta^*$. Clearly, Problem (1) is a special case of Problem (2). The set of optimal solutions of Problem (1) can therefore be characterized as follows:

Proposition 3. Consider $\mathbf{A} \neq \mathbf{0}$. The following statements are true:

(i) The set of optimal solutions of Problem (1) is $(\|\mathbf{A}\|_\theta^*)^{-1} \partial \|\mathbf{A}\|_\theta^*$.

(ii) There is a unique solution if and only if the dual norm function $\|\cdot\|_\theta^*$ is differentiable at \mathbf{A} .

We now focus on deriving some properties of the dual norm $\|\cdot\|_\theta^*$. We have:

$$\begin{aligned} \|\mathbf{A}\|_\theta^* = \max \quad & \langle \mathbf{A}, \mathbf{X} \rangle \\ \text{s.t.} \quad & \|\mathbf{X}\|_\theta \leq 1. \end{aligned} \tag{6}$$

We will use the gauge function and its dual polar function (see Rockafellar [17] for more details) to compute this dual norm.

Proposition 4. *The dual norm $\|\mathbf{A}\|_\theta^*$ with $\theta > 0$ is the optimal value of the following optimization problem:*

$$\begin{aligned} \|\mathbf{A}\|_\theta^* &= \min \max \{ \|\mathbf{Y}\|, \theta^{-1} \|\mathbf{Z}\|_\infty \} \\ \text{s.t. } & \mathbf{Y} + \mathbf{Z} = \mathbf{A}. \end{aligned} \quad (7)$$

Proof. Consider the closed unit ball $\mathcal{C}_* = \{\mathbf{X} \in \mathbb{R}^{m \times n} \mid \|\mathbf{X}\|_* \leq 1\}$ with respect to the nuclear norm and similarly, the unit ball \mathcal{C}_1 with respect to the ℓ_1 -norm $\|\cdot\|_1$. We have the polar of \mathcal{C}_* is the closed unit ball with respect to the spectral norm, $\mathcal{C}_*^\circ = \mathcal{C}$. Similarly, we have: $\mathcal{C}_1^\circ = \mathcal{C}_\infty$, the unit ball with respect to the infinity norm.

Using the definition of gauge functions, we have: $\|\mathbf{X}\|_* = \gamma_{\mathcal{C}_*}(\mathbf{X}) = \min\{\lambda \geq 0 \mid \mathbf{X} \in \lambda\mathcal{C}_*\}$. In addition, the support function $\sigma_{\mathcal{S}}(\mathbf{X}) = \max\{\langle \mathbf{X}, \mathbf{Y} \rangle \mid \mathbf{Y} \in \mathcal{S}\}$ is the gauge function of \mathcal{S}° for all symmetric closed bounded convex set with $\mathbf{0} \in \text{int}(\mathcal{S})$. All unit balls satisfy these conditions; therefore, we obtain the well-known results $\|\mathbf{X}\|_* = \gamma_{\mathcal{C}_*}(\mathbf{X}) = \sigma_{\mathcal{C}_*^\circ}(\mathbf{X}) = \sigma_{\mathcal{C}}(\mathbf{X})$ and $\|\mathbf{X}\|_1 = \gamma_{\mathcal{C}_1}(\mathbf{X}) = \sigma_{\mathcal{C}_1^\circ}(\mathbf{X}) = \sigma_{\mathcal{C}_\infty}(\mathbf{X})$.

Now consider the unit ball $\mathcal{C}_\theta = \{\mathbf{X} \in \mathbb{R}^{m \times n} \mid \|\mathbf{X}\|_* + \theta \|\mathbf{X}\|_1 \leq 1\}$, we have:

$$\mathcal{C}_\theta = \{\mathbf{X} \in \mathbb{R}^{m \times n} \mid \sigma_{\mathcal{C}}(\mathbf{X}) + \theta \sigma_{\mathcal{C}_\infty}(\mathbf{X}) \leq 1\}.$$

Applying the definition of support functions, we have: $\sigma_{\mathcal{C}}(\mathbf{X}) + \theta \sigma_{\mathcal{C}_\infty}(\mathbf{X}) = \sigma_{\mathcal{C} + \theta\mathcal{C}_\infty}(\mathbf{X})$, where $\mathcal{C} + \theta\mathcal{C}_\infty$ is the Minkowski sum of two sets, \mathcal{C} and $\theta\mathcal{C}_\infty$. This set satisfies all the conditions above; therefore, $\sigma_{\mathcal{C} + \theta\mathcal{C}_\infty}(\mathbf{X}) = \gamma_{(\mathcal{C} + \theta\mathcal{C}_\infty)^\circ}(\mathbf{X})$. Thus

$$\mathcal{C}_\theta = \{\mathbf{X} \in \mathbb{R}^{m \times n} \mid \gamma_{(\mathcal{C} + \theta\mathcal{C}_\infty)^\circ}(\mathbf{X}) \leq 1\} = (\mathcal{C} + \theta\mathcal{C}_\infty)^\circ.$$

We also have: $\|\mathbf{A}\|_\theta^* = \sigma_{\mathcal{C}_\theta}(\mathbf{A})$. Thus

$$\|\mathbf{A}\|_\theta^* = \gamma_{\mathcal{C}_\theta^\circ}(\mathbf{A}) = \gamma_{\mathcal{C} + \theta\mathcal{C}_\infty}(\mathbf{A}).$$

We have: $\gamma_{\mathcal{C} + \theta\mathcal{C}_\infty}(\mathbf{A}) = \min\{\lambda \geq 0 \mid \mathbf{A} \in \lambda(\mathcal{C} + \theta\mathcal{C}_\infty)\}$ or equivalently,

$$\begin{aligned} \gamma_{\mathcal{C} + \theta\mathcal{C}_\infty}(\mathbf{A}) &= \min \lambda \\ \text{s.t. } & \mathbf{A} = \mathbf{Y} + \mathbf{Z}, \\ & \|\mathbf{Y}\| \leq \lambda, \\ & \theta^{-1} \|\mathbf{Z}\|_\infty \leq \lambda, \\ & \lambda \geq 0. \end{aligned}$$

Rewriting the minimization problem above, we obtain the final result as shown in (7):

$$\begin{aligned} \|\mathbf{A}\|_{\theta}^* &= \min \max \{ \|\mathbf{Y}\|, \theta^{-1} \|\mathbf{Z}\|_{\infty} \} \\ \text{s.t. } & \mathbf{Y} + \mathbf{Z} = \mathbf{A}. \end{aligned}$$

□

We can now derive the optimality conditions for both problems (6) and (7):

Lemma 3. *Nonzero feasible solutions \mathbf{X} and (\mathbf{Y}, \mathbf{Z}) are optimal for Problem (6) and (7) respectively if and only if they satisfy the conditions below:*

(i) $\|\mathbf{Y}\| = \theta^{-1} \|\mathbf{Z}\|_{\infty}$,

(ii) $\mathbf{X} \in \alpha \partial \|\mathbf{Y}\|$, $\alpha \geq 0$,

(iii) $\mathbf{X} \in \beta \partial \|\mathbf{Z}\|_{\infty}$, $\beta \geq 0$, and

(iv) $\alpha + \theta\beta = 1$.

Proof. We first prove the weak duality result. Consider feasible solutions \mathbf{X} and (\mathbf{Y}, \mathbf{Z}) for Problem (6) and (7) respectively, we have:

$$\begin{aligned} \langle \mathbf{X}, \mathbf{A} \rangle &= \langle \mathbf{X}, \mathbf{Y} \rangle + \langle \mathbf{X}, \mathbf{Z} \rangle \\ &\leq \|\mathbf{X}\|_* \|\mathbf{Y}\| + \|\mathbf{X}\|_1 \|\mathbf{Z}\|_{\infty} \\ &\leq \|\mathbf{X}\|_* \max\{\|\mathbf{Y}\|, \theta^{-1} \|\mathbf{Z}\|_{\infty}\} + \theta \|\mathbf{X}\|_1 \max\{\|\mathbf{Y}\|, \theta^{-1} \|\mathbf{Z}\|_{\infty}\} \\ &= (\|\mathbf{X}\|_* + \theta \|\mathbf{X}\|_1) \max\{\|\mathbf{Y}\|, \theta^{-1} \|\mathbf{Z}\|_{\infty}\} \\ &\leq \max\{\|\mathbf{Y}\|, \theta^{-1} \|\mathbf{Z}\|_{\infty}\}. \end{aligned}$$

The strong duality result shows that \mathbf{X} and (\mathbf{Y}, \mathbf{Z}) are the optimal solutions if and only if $\langle \mathbf{X}, \mathbf{A} \rangle = \max\{\|\mathbf{Y}\|, \theta^{-1} \|\mathbf{Z}\|_{\infty}\}$. This happens if and only if all the conditions below are satisfied:

(i) $\langle \mathbf{X}, \mathbf{Y} \rangle = \|\mathbf{X}\|_* \|\mathbf{Y}\|$ and $\langle \mathbf{X}, \mathbf{Z} \rangle = \|\mathbf{X}\|_1 \|\mathbf{Z}\|_{\infty}$,

(ii) $\|\mathbf{Y}\| = \max\{\|\mathbf{Y}\|, \theta^{-1} \|\mathbf{Z}\|_{\infty}\} = \theta^{-1} \|\mathbf{Z}\|_{\infty}$, and

(iii) $\|\mathbf{X}\|_* + \theta \|\mathbf{X}\|_1 = 1$.

The first two conditions are equivalent to the fact that $\mathbf{X} = \alpha \partial \|\mathbf{Y}\|$, where $\alpha = \|\mathbf{X}\|_*$, and $\mathbf{X} = \beta \partial \|\mathbf{Z}\|_{\infty}$, where $\beta = \|\mathbf{X}\|_1$. The second condition is simply $\|\mathbf{Y}\| = \theta^{-1} \|\mathbf{Z}\|_{\infty}$ and the third condition is equivalent to $\alpha + \theta\beta = 1$. Thus we have proved the necessary and sufficient optimality conditions for Problem (6) and (7). □

Using these optimality conditions, we can obtain simple sufficient conditions for the uniqueness of the optimal solution \mathbf{X} :

Proposition 5. *Consider the feasible solution \mathbf{X} of Problem (6). If there exists (\mathbf{Y}, \mathbf{Z}) that satisfies the conditions below,*

(i) $\mathbf{Y} + \mathbf{Z} = \mathbf{A}$ and $\|\mathbf{Y}\| = \theta^{-1} \|\mathbf{Z}\|_\infty$,

(ii) $\mathbf{X} \in \alpha \partial \|\mathbf{Y}\|$, $\alpha \geq 0$,

(iii) $\mathbf{X} \in \beta \partial \|\mathbf{Z}\|_\infty$, $\beta \geq 0$,

(iv) $\alpha + \theta\beta = 1$, and

(v) $\|\cdot\|$ is differentiable at \mathbf{Y} or $\|\cdot\|_\infty$ is differentiable at \mathbf{Z} ,

then \mathbf{X} is the unique optimal solution of Problem (6).

Proof. Using the first four conditions, we can prove that \mathbf{X} is an optimal solution of Problem (6) and (\mathbf{Y}, \mathbf{Z}) is an optimal solution of Problem (7). Now assume that $\|\cdot\|$ is differentiable at \mathbf{Y} , we have: $\partial \|\mathbf{Y}\|$ is a singleton, $\partial \|\mathbf{Y}\| = \{\mathbf{V}\}$. Thus we have:

$$\|\mathbf{X}\|_1 = \alpha \|\mathbf{V}\|_1 = \beta \Rightarrow \alpha(1 + \theta \|\mathbf{V}\|_1) = 0 \Rightarrow \alpha = \frac{1}{1 + \theta \|\mathbf{V}\|_1}.$$

Assume there is another optimal solution $\bar{\mathbf{X}} \neq \mathbf{X}$ of Problem (6). Applying Lemma 3, we will have $\bar{\mathbf{X}} \in \bar{\alpha} \partial \|\mathbf{Y}\|$ and similarly $\bar{\mathbf{X}} \in \bar{\beta} \partial \|\mathbf{Z}\|_\infty$ with $\bar{\alpha} + \theta \bar{\beta} = 1$. Same calculation results in $\bar{\alpha} = \alpha$ (contradiction). Thus \mathbf{X} is the unique optimal solution of Problem (6). Similar arguments can be used to prove the uniqueness of \mathbf{X} if $\|\cdot\|_\infty$ is differentiable at \mathbf{Z} . \square

Proposition 5 relies on dual solutions \mathbf{Y} and \mathbf{Z} to show the uniqueness of the primal solution \mathbf{X} . Next, we will focus on the low-rank and sparse property of the optimal solution \mathbf{X} for different values of θ . The following theorem provides the sufficient conditions on matrix \mathbf{A} for the rank-one property (and uniqueness) of the optimal solution \mathbf{X} when θ is small enough.

Theorem 2. *If \mathbf{A} satisfies the condition $\sigma_1(\mathbf{A}) > \sigma_2(\mathbf{A})$, then Problem (6) has a (unique) rank-one optimal solution \mathbf{X} for all $0 \leq \theta < \theta_A$, where $\theta_A = \frac{1}{\sqrt{mn}} \left(\frac{\sigma_1(\mathbf{A}) - \sigma_2(\mathbf{A})}{3\sigma_1(\mathbf{A}) - \sigma_2(\mathbf{A})} \right)$.*

Proof. The optimality conditions in Lemma 3 show that there exist \mathbf{Y} and \mathbf{Z} such that $\mathbf{A} = \mathbf{Y} + \mathbf{Z}$, $\|\mathbf{Z}\|_\infty = \theta \|\mathbf{Y}\|$ and $\mathbf{X} \in \alpha \partial \|\mathbf{Y}\|$. Applying a standard perturbation theorem of singular values (see Cor. 8.6.2 of [11]), we have:

$$|\sigma_i(\mathbf{A}) - \sigma_i(\mathbf{Y})| \leq \|\mathbf{Z}\|, \quad i = 1, 2.$$

We also have: $\|\mathbf{Z}\| \leq \sqrt{mn} \|\mathbf{Z}\|_\infty$. Thus

$$\|\mathbf{Z}\| \leq \sqrt{mn} (\theta \|\mathbf{Y}\|) = \sqrt{mn} (\theta \sigma_1(\mathbf{Y})).$$

For all $0 \leq \theta < \theta_A$, we have:

$$\sigma_1(\mathbf{Y}) \leq \sigma_1(\mathbf{A}) + \|\mathbf{Z}\| \leq \sigma_1(\mathbf{A}) + \sqrt{mn} (\theta \sigma_1(\mathbf{Y})) < \sigma_1(\mathbf{A}) + \sqrt{mn} (\theta_A \sigma_1(\mathbf{Y})).$$

This implies

$$(1 - \theta_A \sqrt{mn}) \sigma_1(\mathbf{Y}) < \sigma_1(\mathbf{A}) \Leftrightarrow \frac{2\sigma_1(\mathbf{A})}{3\sigma_1(\mathbf{A}) - \sigma_2(\mathbf{A})} \sigma_1(\mathbf{Y}) < \sigma_1(\mathbf{A}) \Leftrightarrow \sigma_1(\mathbf{Y}) < \frac{1}{2}(3\sigma_1(\mathbf{A}) - \sigma_2(\mathbf{A})).$$

We then have:

$$\|\mathbf{Z}\| \leq \sqrt{mn} (\theta \sigma_1(\mathbf{Y})) < \sqrt{mn} (\theta_A \sigma_1(\mathbf{Y})) < \frac{1}{2}(\sigma_1(\mathbf{A}) - \sigma_2(\mathbf{A})).$$

Thus

$$\sigma_1(\mathbf{Y}) \geq \sigma_1(\mathbf{A}) - \|\mathbf{Z}\| > \frac{1}{2}(\sigma_1(\mathbf{A}) + \sigma_2(\mathbf{A})) > \sigma_2(\mathbf{A}) + \|\mathbf{Z}\| \geq \sigma_2(\mathbf{Y}).$$

We have $\sigma_1(\mathbf{Y}) > \sigma_2(\mathbf{Y})$; therefore, $\|\cdot\|$ is differentiable at \mathbf{Y} . According to Proposition 5, we have \mathbf{X} is the unique rank-one optimal solution of Problem (6). \square

The last result of this section concerns the nonnegativity of \mathbf{X} . If \mathbf{A} is nonnegative, then one might expect \mathbf{X} to be nonnegative. For $\theta = 0$ or $\theta = \infty$, this is certainly true by preceding results in this section. It is not always necessarily true for intermediate values of θ . The following theorem shows that, at least in the rank-one case, nonnegativity is assured.

Theorem 3. *Consider the set of optimal solutions of Problem (6) when $\mathbf{A} \geq \mathbf{0}$. We have:*

(i) *If Problem (6) has a rank-one optimal solution, then there exists a nonnegative rank-one optimal solution.*

(ii) *If $\theta > 1$, then all optimal solutions of Problem (6) are nonnegative.*

Proof.

(i) Consider a rank-one optimal solution \mathbf{X} , $\mathbf{X} = \sigma \mathbf{u} \mathbf{v}^T$, of Problem (6). We prove that $|\mathbf{X}| = \sigma |\mathbf{u}| |\mathbf{v}|^T \geq \mathbf{0}$ is also an optimal solution. Let $\tilde{\mathbf{X}}$ denote $|\mathbf{X}|$. We have:

$$\|\tilde{\mathbf{X}}\|_\theta = \|\tilde{\mathbf{X}}\|_* + \theta \|\tilde{\mathbf{X}}\|_1 = \|\mathbf{X}\|_* + \theta \|\mathbf{X}\|_1.$$

In addition, $\langle \mathbf{A}, \tilde{\mathbf{X}} \rangle \geq \langle \mathbf{A}, \mathbf{X} \rangle$ since $\mathbf{A} \geq \mathbf{0}$. Thus clearly $\tilde{\mathbf{X}}$ is also an optimal solution.

(ii) Assume that there exists an optimal solution \mathbf{X} of Problem (6) is not nonnegative. Without loss of generality, assume $x_{11} < 0$. Consider $\mathbf{X}(\epsilon) = \mathbf{X} + \epsilon \mathbf{E}_{11}$, where $\epsilon > 0$ and \mathbf{E}_{11} is the matrix of all zeros except the element $\mathbf{E}_{11}(1, 1) = 1$, we have:

$$\|\mathbf{X}(\epsilon)\|_* \leq \|\mathbf{X}\|_* + \epsilon \|\mathbf{E}_{11}\|_* = \|\mathbf{X}\|_* + \epsilon.$$

In addition, $\|\mathbf{X}(\epsilon)\|_1 = \|\mathbf{X}\|_1 - \epsilon$ if $\epsilon \leq |x_{11}|$. Therefore, we have:

$$\|\mathbf{X}(\epsilon)\|_\theta \leq \|\mathbf{X}\| + (1 - \theta)\epsilon = 1 + (1 - \theta)\epsilon < 1, \quad \forall 0 < \epsilon \leq |x_{11}|.$$

Here we assume that $\mathbf{A} \neq \mathbf{0}$; therefore, $\|\mathbf{X}\|_\theta = 1$. We also have

$$\langle \mathbf{A}, \mathbf{X}(\epsilon) \rangle = \langle \mathbf{A}, \mathbf{X} \rangle + \epsilon a_{11} \geq \langle \mathbf{A}, \mathbf{X} \rangle.$$

Now consider $\bar{\mathbf{X}} = \frac{1}{1 + (1 - \theta)\epsilon} \mathbf{X}(\epsilon)$. Clearly, $\|\bar{\mathbf{X}}\|_\theta = 1$ and $\langle \mathbf{A}, \bar{\mathbf{X}} \rangle > \langle \mathbf{A}, \mathbf{X} \rangle > 0$ (contradiction). Thus all optimal solutions of Problem (6) are nonnegative if $\theta > 1$.

□

3 Sparsity

As mentioned in the introduction, the penalty term $\theta \|\mathbf{X}\|_1$ in the objective function of (1) is intended to promote sparsity of \mathbf{X} . For some very simple convex optimization problems with an ℓ_1 penalty term, e.g., the unconstrained problem of minimizing $\|\mathbf{x} - \mathbf{c}\|_2 + \theta \|\mathbf{x}\|_1$ for a given vector \mathbf{c} , it is known that sparsity increases monotonically with θ (i.e., if \mathbf{x}_1^* is the optimizer for θ_1 and \mathbf{x}_2^* is the optimizer for θ_2 with $\theta_1 \leq \theta_2$, then the indices of nonzeros of \mathbf{x}_2^* are a subset of the indices of nonzeros of \mathbf{x}_1^*).

For a more complicated problem such as (1), monotonicity does not hold in general. But nonetheless, some weaker statements about the relationship between θ and sparsity are possible. Two such results are derived in this section. We start with a lemma that leads to a sparsity result.

Lemma 4. *Assume $\mathbf{X} = \sigma \mathbf{u} \mathbf{v}^T$, where $\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1$, $\mathbf{u} \geq \mathbf{0}$, and $\mathbf{v} \geq \mathbf{0}$, is the optimal solution of Problem (6). If $u_i > u_j = 0$ then*

$$\|A\|_\theta^* = \frac{\mathbf{a}_i^T \mathbf{v}}{\theta \|\mathbf{v}\|_1 + u_i} \geq \frac{\mathbf{a}_j^T \mathbf{v}}{\theta \|\mathbf{v}\|_1}, \quad (8)$$

where \mathbf{a}_i and \mathbf{a}_j are the i th and j th row of \mathbf{A} .

Proof. We again assume here $\mathbf{A} \neq \mathbf{0}$, which means $\|\mathbf{X}\|_\theta = 1$. We have:

$$\|\mathbf{X}\|_\theta = \sigma + \theta\sigma \|\mathbf{u}\|_1 \|\mathbf{v}\|_1 = 1 \Leftrightarrow \sigma = \frac{1}{1 + \theta \|\mathbf{u}\|_1 \|\mathbf{v}\|_1}.$$

Consider $\mathbf{X}(\epsilon) = \sigma(\mathbf{u} + \epsilon \mathbf{e}_i) \mathbf{v}^T$, where $\epsilon \geq 0$ and \mathbf{e}_i is the i th unit vector, we have:

$$\|\mathbf{u} + \epsilon \mathbf{e}_i\|_2^2 = \|\mathbf{u}\|_2^2 + (u_i + \epsilon)^2 - u_i^2 = 1 + 2u_i\epsilon + \epsilon^2.$$

Thus $\|\mathbf{X}(\epsilon)\|_* = \sigma\sqrt{1 + 2u_i\epsilon + \epsilon^2}$. On the other hand, $\|\mathbf{X}(\epsilon)\|_1 = \sigma(\|\mathbf{u}\|_1 + \epsilon) \|\mathbf{v}\|_1$. So we have:

$$\|\mathbf{X}(\epsilon)\|_\theta = 1 + \sigma \left(\sqrt{1 + 2u_i\epsilon + \epsilon^2} - 1 + \theta\epsilon \|\mathbf{v}\|_1 \right) = 1 + \sigma\epsilon \left(\frac{2u_i + \epsilon}{\sqrt{1 + 2u_i\epsilon + \epsilon^2} + 1} + \theta \|\mathbf{v}\|_1 \right).$$

Let $\alpha = \frac{2u_i + \epsilon}{\sqrt{1 + 2u_i\epsilon + \epsilon^2} + 1} + \theta \|\mathbf{v}\|_1$ and consider $\bar{\mathbf{X}} = \frac{\mathbf{X}(\epsilon)}{1 + \sigma\epsilon\alpha}$, we have: $\|\bar{\mathbf{X}}\|_\theta = 1$ and

$$\langle \mathbf{A}, \bar{\mathbf{X}} \rangle = \frac{\langle \mathbf{A}, \mathbf{X}(\epsilon) \rangle}{1 + \sigma\epsilon\alpha} = \frac{\|\mathbf{A}\|_\theta^* + \sigma\epsilon \mathbf{a}_i^T \mathbf{v}}{1 + \sigma\epsilon\alpha} \leq \|\mathbf{A}\|_\theta^*.$$

With $\sigma > 0$ and $\epsilon > 0$, we obtain the following inequality

$$\mathbf{a}_i^T \mathbf{v} \leq \left(\frac{2u_i + \epsilon}{\sqrt{1 + 2u_i\epsilon + \epsilon^2} + 1} + \theta \|\mathbf{v}\|_1 \right) \|\mathbf{A}\|_\theta^*.$$

Taking the limit $\epsilon \rightarrow 0^+$, we have:

$$\|\mathbf{A}\|_\theta^* \geq \frac{\mathbf{a}_i^T \mathbf{v}}{\theta \|\mathbf{v}\|_1 + u_i}.$$

Now consider the case in which $u_i > 0$ and set $\mathbf{X}(\epsilon) = \sigma(\mathbf{u} - \epsilon \mathbf{e}_i) \mathbf{v}^T$, where $0 \leq \epsilon \leq u_i$. Similarly, we have:

$$\|\mathbf{X}(\epsilon)\|_\theta = 1 - \sigma\epsilon \left(\frac{2u_i - \epsilon}{\sqrt{1 - 2u_i\epsilon + \epsilon^2} + 1} + \theta \|\mathbf{v}\|_1 \right).$$

This implies the following inequality

$$\mathbf{a}_i^T \mathbf{v} \geq \left(\frac{2u_i - \epsilon}{\sqrt{1 - 2u_i\epsilon + \epsilon^2} + 1} + \theta \|\mathbf{v}\|_1 \right) \|\mathbf{A}\|_\theta^*.$$

Again, taking the limit $\epsilon \rightarrow 0^+$, we have:

$$\|\mathbf{A}\|_\theta^* \leq \frac{\mathbf{a}_i^T \mathbf{v}}{\theta \|\mathbf{v}\|_1 + u_i}.$$

From these two results, we can see that if $u_i > u_j = 0$, then

$$\|\mathbf{A}\|_\theta^* = \frac{\mathbf{a}_i^T \mathbf{v}}{\theta \|\mathbf{v}\|_1 + u_i} \geq \frac{\mathbf{a}_j^T \mathbf{v}}{\theta \|\mathbf{v}\|_1}.$$

□

Since the roles of columns and rows are interchangeable, we also have the following result. If $v_k > v_l = 0$, then

$$\|A\|_{\theta}^* = \frac{\mathbf{u}^T \mathbf{A}_k}{\theta \|\mathbf{u}\|_1 + v_k} \geq \frac{\mathbf{u}^T \mathbf{A}_l}{\theta \|\mathbf{u}\|_1}, \quad (9)$$

where \mathbf{A}_k and \mathbf{A}_l are k -th and l -th column of \mathbf{A} .

The sparsity structure of $\mathbf{X} = \sigma \mathbf{u} \mathbf{v}^T$ depends on the sparsity structure of \mathbf{u} and \mathbf{v} . The results obtained above help us derive some conditions under which a row (or column) of \mathbf{X} is zero.

Corollary 1. *Consider two rows \mathbf{a}_i^T and \mathbf{a}_j^T of matrix \mathbf{A} . If $\min_k \{a_{ik}\} \geq \alpha \max_k \{a_{jk}\}$, where $\alpha > 1$, then for every $\theta > \frac{1}{\alpha - 1}$ and for every nonnegative rank-one optimal solution \mathbf{X} of Problem (6), the j th row of \mathbf{X} is zero.*

Proof. Assume that $\mathbf{X} = \sigma \mathbf{u} \mathbf{v}^T$, where $\mathbf{u} \geq \mathbf{0}$ and $\mathbf{v} \geq \mathbf{0}$, and $\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1$. We have:

$$\frac{\mathbf{a}_i^T \mathbf{v}}{\theta \|\mathbf{v}\|_1 + u_i} \geq \frac{\min_k \{a_{ik}\} \|\mathbf{v}\|_1}{\theta \|\mathbf{v}\|_1 + 1} \geq \frac{\alpha \max_k \{a_{jk}\} \|\mathbf{v}\|_1}{\theta \|\mathbf{v}\|_1 + 1} = \frac{\alpha \max_k \{a_{jk}\}}{\theta + \frac{1}{\|\mathbf{v}\|_1}} \geq \frac{\alpha \max_k \{a_{jk}\}}{\theta + 1},$$

since $\|\mathbf{v}\|_1 \geq \|\mathbf{v}\|_2 = 1$. On the other hand, we also have the following inequality

$$\frac{\mathbf{a}_j^T \mathbf{v}}{\theta \|\mathbf{v}\|_1 + u_j} \leq \frac{\max_k \{a_{jk}\} \|\mathbf{v}\|_1}{\theta \|\mathbf{v}\|_1} = \frac{\max_k \{a_{jk}\}}{\theta}.$$

We have:

$$\frac{\alpha \max_k \{a_{jk}\}}{\theta + 1} - \frac{\max_k \{a_{jk}\}}{\theta} = \max_k \{a_{kj}\} \frac{(\alpha - 1)\theta - 1}{\theta(\theta + 1)} > 0, \quad \forall \theta > \frac{1}{\alpha - 1}.$$

Thus we have:

$$\frac{\mathbf{a}_i^T \mathbf{v}}{\theta \|\mathbf{v}\|_1 + u_i} > \frac{\mathbf{a}_j^T \mathbf{v}}{\theta \|\mathbf{v}\|_1 + u_j}, \quad \forall \theta > \frac{1}{\alpha - 1},$$

which means $u_j = 0$ according to Lemma 4. Thus the j -th row of \mathbf{X} is zero. \square

We would like to use these results to build up results for columns and rows simultaneously. More exactly, consider a subset $\mathcal{I} \subset \{1, \dots, m\}$ and $\mathcal{J} \subset \{1, \dots, n\}$, we would like to obtain conditions on magnitudes of elements of $\mathbf{A}(\mathcal{I}, \mathcal{J})$ as compared to those of the remaining elements of \mathbf{A} to guarantee that all rows and columns that are not in \mathcal{I} and \mathcal{J} have to be zero in the nonnegative rank-one optimal matrix \mathbf{X} of Problem (6) for $\theta \geq \theta_0$. One of the difficulties here is that under these conditions, there is a coupling relationship between rows and columns. More exactly, in order to prove the rows that are not in \mathcal{I} are zero, we need to prove the columns that are not in \mathcal{J} are small or zero at the same time.

Lemma 4 and Corollary 1 are based on local optimality conditions with respect to rows or columns. We can obtain additional results on the sparsity of the optimal solution \mathbf{X} using the global optimality conditions.

The following theorem states that if the weight of nonnegative matrix \mathbf{A} is concentrated in a particular subblock then for θ sufficiently large, the optimal solution \mathbf{X} will have nonzero entries only in that subblock. ‘‘Concentration of weight’’ in this sense means that the average of those entries dominates all the other entries of the matrix.

As a special case, this theorem implies that if the maximum entry of \mathbf{A} is unique, then for θ sufficiently large, \mathbf{X} will be a singleton matrix whose unique nonzero entry corresponds to the maximum entry of \mathbf{A} .

Theorem 4. *Assume $\mathbf{A} \geq \mathbf{0}$. Let \mathcal{I} and \mathcal{J} be subsets of $\{1, \dots, m\}$ and $\{1, \dots, n\}$, respectively; $|\mathcal{I}| = M$ and $|\mathcal{J}| = N$. Define $\bar{a}(\mathcal{I}, \mathcal{J}) = \frac{1}{MN} \sum_{(i,j) \in (\mathcal{I}, \mathcal{J})} a_{ij}$ and $a_{\max}(\overline{\mathcal{I}}, \overline{\mathcal{J}}) = \max_{(i,j) \notin (\mathcal{I}, \mathcal{J})} a_{ij}$. If $\bar{a}(\mathcal{I}, \mathcal{J}) > a_{\max}(\overline{\mathcal{I}}, \overline{\mathcal{J}})$ then all optimal solutions \mathbf{X} of Problem (6) are sparse, $x_{ij} = 0$ for all $(i, j) \notin (\mathcal{I}, \mathcal{J})$, for all $\theta > \theta_B$, where $\theta_B = \frac{1}{\sqrt{MN}} \left(\frac{\bar{a}(\mathcal{I}, \mathcal{J})\sqrt{MN} + a_{\max}(\overline{\mathcal{I}}, \overline{\mathcal{J}})}{\bar{a}(\mathcal{I}, \mathcal{J}) - a_{\max}(\overline{\mathcal{I}}, \overline{\mathcal{J}})} \right)$.*

Proof. Assume there exists an optimal solution \mathbf{X} of Problem (6) such that $x_{ij} \neq 0$ for some $(i, j) \notin (\mathcal{I}, \mathcal{J})$ when $\theta > \theta_B$. We have: $\theta_B > 1$; therefore, according to Theorem 3, $\mathbf{X} \geq \mathbf{0}$. Thus $x_{ij} > 0$. We also have: $\mathbf{A} \neq \mathbf{0}$; therefore $\mathbf{X} \neq \mathbf{0}$ and $\|\mathbf{X}\|_\theta = 1$. Consider two cases, $a_{ij} = 0$ and $a_{ij} > 0$.

If $a_{ij} = 0$, let $\mathbf{X}_0 = \mathbf{X} - x_{ij}\mathbf{E}_{ij}$, where \mathbf{E}_{ij} is the matrix of all zeros but the element $\mathbf{E}_{ij}(i, j) = 1$. We have: $\|\mathbf{X}_0\| \leq \|\mathbf{X}\| + x_{ij}$ and $\|\mathbf{X}_0\|_1 = \|\mathbf{X}\|_1 - x_{ij}$. Thus

$$\|\mathbf{X}_0\|_\theta \leq \|\mathbf{X}\|_\theta + (1 - \theta)x_{ij} < 1, \quad \forall \theta > \theta_B > 1.$$

We also have: since $a_{ij} = 0$, $\langle \mathbf{A}, \mathbf{X}_0 \rangle = \langle \mathbf{A}, \mathbf{X} \rangle = \|\mathbf{A}\|_\theta^* > 0$. Thus $\mathbf{X}_0 \neq \mathbf{0}$ or $\|\mathbf{X}_0\|_\theta > 0$. Define $\mathbf{X}_0^s = \frac{1}{\|\mathbf{X}_0\|_\theta} \mathbf{X}_0$, we have: \mathbf{X}_0^s is a feasible solution of Problem (6) with the objective value $\langle \mathbf{A}, \mathbf{X}_0^s \rangle = \frac{\|\mathbf{A}\|_\theta^*}{\|\mathbf{X}_0\|_\theta} > \|\mathbf{A}\|_\theta^*$ (contradiction).

Now consider the case when $a_{ij} > 0$. Define $\mathbf{D} = \mathbf{e}_\mathcal{I} \mathbf{e}_\mathcal{J}^T - MNr \mathbf{e}_i \mathbf{e}_j^T$, where $r = \frac{\bar{a}(\mathcal{I}, \mathcal{J})}{a_{ij}}$, $\mathbf{e}_\mathcal{I} = \sum_{i \in \mathcal{I}} \mathbf{e}_i$, $\mathbf{e}_i \in \mathbb{R}^m$ is the i th unit vector in \mathbb{R}^m , and similarly, $\mathbf{e}_\mathcal{J} = \sum_{j \in \mathcal{J}} \mathbf{e}_j$, $\mathbf{e}_j \in \mathbb{R}^n$ is the j -th unit vector in \mathbb{R}^n . We have r is well-defined since $\mathbf{A} > \mathbf{0}$ and $r > 1$. We now consider a new solution $\mathbf{X}_\alpha = \mathbf{X} + \alpha \mathbf{D}$, where $0 < \alpha \leq \frac{x_{ij}}{MNr}$. Clearly, $\mathbf{X}_\alpha \geq \mathbf{0}$. Thus we have: $\|\mathbf{X}_\alpha\|_1 = \|\mathbf{X}\|_1 + \alpha MN(1 - r)$. Applying the triangle inequality, we can bound $\|\mathbf{X}_\alpha\|_*$ as follows:

$$\|\mathbf{X}_\alpha\|_* \leq \|\mathbf{X}\|_* + \alpha \|\mathbf{D}\|_* \leq \|\mathbf{X}\|_* + \alpha(\sqrt{MN} + MNr).$$

Thus we have:

$$\|\mathbf{X}_\alpha\|_\theta \leq \|\mathbf{X}\|_\theta + \alpha \sqrt{MN} \left[(1 + r\sqrt{MN}) + \theta \sqrt{MN}(1 - r) \right].$$

Since $\theta > \theta_B = \frac{1}{\sqrt{MN}} \left(\frac{\bar{a}(\mathcal{I}, \mathcal{J})\sqrt{MM} + a_{\max}(\overline{\mathcal{I}}, \overline{\mathcal{J}})}{\bar{a}(\mathcal{I}, \mathcal{J}) - a_{\max}(\overline{\mathcal{I}}, \overline{\mathcal{J}})} \right)$ and $0 < a_{ij} \leq a_{\max}(\overline{\mathcal{I}}, \overline{\mathcal{J}}) < \bar{a}(\mathcal{I}, \mathcal{J})$, we have:

$$\theta > \frac{1}{\sqrt{MN}} \left(\frac{1 + r\sqrt{MN}}{r - 1} \right).$$

This implies that $\|\mathbf{X}_\alpha\|_\theta < \|\mathbf{X}\|_\theta = 1$ for all $0 < \alpha \leq \frac{x_{ij}}{MNr}$. Now consider the scaled solution $\mathbf{X}_\alpha^s = \frac{1}{\|\mathbf{X}_\alpha\|_\theta} \mathbf{X}_\alpha$, which is also a feasible solution of Problem (6). In terms of the objective, we have: $\langle \mathbf{A}, \mathbf{D} \rangle = \mathbf{e}_{\mathcal{I}}^T \mathbf{A} \mathbf{e}_{\mathcal{J}} - MNra_{ij} = 0$. Thus $\langle \mathbf{A}, \mathbf{X}_\alpha \rangle = \langle \mathbf{A}, \mathbf{X} \rangle = \|\mathbf{A}\|_\theta^*$ for all α . We then have:

$$\langle \mathbf{A}, \mathbf{X}_\alpha^s \rangle = \frac{\|\mathbf{A}\|_\theta^*}{\|\mathbf{X}_\alpha\|_\theta} > \|\mathbf{A}\|_\theta^*, \quad \forall \alpha \in \left(0, \frac{x_{ij}}{MNr} \right],$$

which is a contradiction because $\|\mathbf{A}\|_\theta^*$ is the optimal value of Problem (6).

Thus we can conclude that if $\mathbf{A} \geq \mathbf{0}$, all optimal solutions \mathbf{X} of Problem (6) are sparse with $x_{ij} = 0$ for all $(i, j) \notin (\mathcal{I}, \mathcal{J})$ when $\theta > \theta_B$. \square

4 Random noise

The main technical result of this article is that the proposed algorithm can find a large rank-one submatrix hidden in a substantial amount of noise. The noise takes two forms: the rank-one submatrix itself has random noise added to it (so that its rank is no longer 1), and the entries outside the rank-one submatrix are generated by a random process.

First, we recall the following definition: a random variable x is *b-subgaussian* if its mean is zero, and if there exists a $b > 0$ such that for all $t \geq 0$,

$$\mathbb{P}(|x| \geq t) \leq \exp(-t^2/(2b^2)). \quad (10)$$

For example, a normally distributed variable or any mean-zero variable with a discrete distribution is subgaussian.

The result of this section is the following bound. For this entire section, we adopt the notation that \mathbf{e}_M denotes the vector of all 1's of length M , and similarly for \mathbf{e}_N .

Theorem 5. *Let \mathbf{A} be an $m \times n$ matrix defined as follows.*

$$\mathbf{A} = \begin{pmatrix} \sigma_0 \mathbf{u}_0 \mathbf{v}_0^T & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} + \begin{pmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{R}_{21} & \mathbf{R}_{22} \end{pmatrix}, \quad (11)$$

where $\sigma_0 > 0$, $\mathbf{u}_0 \in \mathbb{R}^M$, $\mathbf{u}_0 \geq \mathbf{0}$, $M < m$, and $\mathbf{v}_0 \in \mathbb{R}^N$, $\mathbf{v}_0 \geq \mathbf{0}$, $N < n$. Furthermore, assume that $\mathbf{u}_0 = \mathbf{e}_M + \mathbf{p}$ with $\|\mathbf{p}\|_2 \leq c_1\sqrt{M}$, and $\mathbf{v}_0 = \mathbf{e}_N + \mathbf{q}$ with $\|\mathbf{q}\|_2 \leq c_2\sqrt{N}$. The matrix \mathbf{R} is a random matrix with i.i.d. nonnegative elements r_{ij} with mean $c_3\sigma_0$, where $c_3 > 0$ is a constant, such that $r_{ij}/\sigma_0 - c_3$ is b -subgaussian. Here c_1, c_2, c_3 are positive constants. Assume that these scalar constants c_1, c_2, c_3 satisfy the following relations

$$c_5 \leq 1/3, \quad c_3 + c_5 < 1. \quad (12)$$

where c_5 is chosen to satisfy

$$c_5 > c_1 + c_2 + c_1c_2. \quad (13)$$

Under these hypotheses concerning \mathbf{A} , and assuming θ satisfies

$$\theta \leq \min\left(\frac{1}{c_3 + c_5}, \frac{1 + c_3 - 3c_5}{2c_5}\right) \cdot \frac{1}{\sqrt{MN}}, \quad (14)$$

$$\theta \geq \frac{2c_3}{1 - c_3 - c_5} \cdot \frac{1}{\sqrt{MN}}, \quad (15)$$

the solution \mathbf{X} to problem (1) is a rank-one matrix with positive entries in positions that are indexed by $\{1, \dots, M\} \times \{1, \dots, N\}$ and zeros elsewhere with probability exponentially close to 1 (i.e., of the form $1 - \exp(-(M + N)^{\text{const}})$) provided that $MN \geq \Omega((M + N)^{4/3})$ and $MN \geq \Omega(m + n)$. Here, the constants implicit in the $\Omega(\cdot)$ notation depend on b and c_5 . See (32)–(35) below for a detailed presentation of these constants.

Remarks.

1. Naturally, the theorem also applies if the MN distinguished entries occur as any $M \times N$ submatrix of \mathbf{A} ; we have numbered the distinguished submatrix first in order to simplify notation.
2. It is not enough to assume simply that $\mathbf{u}_0 > \mathbf{0}$ and $\mathbf{v}_0 > \mathbf{0}$ because if these vectors have very small entries, then they cannot be distinguished from the noise.
3. This theorem is not a consequence of Theorem 4 because the hypotheses do not force entries outside the distinguished block to be smaller than the average of the distinguished block's entries.
4. The relationships among the constants as well as (14), (15) can all be satisfied provided c_3, c_5 are sufficiently small.
5. The result holds with probability exponentially close to 1 as long as $M \sim N$ and $M \geq \Omega(m^{1/2})$, $N \geq \Omega(n^{1/2})$. Thus, the rank-one submatrix can be much smaller than the entire matrix \mathbf{A} .

Before beginning the proof of the theorem, we require the following key lemma regarding matrices constructed from independent b -subgaussian random variables.

Lemma 5. *Let $\mathbf{B} \in \mathbb{R}^{m \times n}$ be a random matrix, where b_{ij} are independent b -subgaussian random variables for all $i = 1, \dots, m$, and $j = 1, \dots, n$. Then for any $u > 0$,*

$$(i) \mathbb{P}(\|\mathbf{B}\| \geq u) \leq \exp\left(-\left(\frac{8u^2}{81b^2} - (\log 7)(m+n)\right)\right)$$

$$(ii) \mathbb{P}(\|\mathbf{CB}\| \geq u) \leq \exp\left(-\left(\frac{8u^2}{81b^2 \|\mathbf{C}\|^2} - (\log 7)(m+n)\right)\right), \text{ where } \mathbf{C} \text{ is a deterministic matrix.}$$

The proof of this lemma follows the proof techniques by Litvak et al. [14]. Major steps are shown as follows.

Proof.

(i) We have: $\|\mathbf{B}\| = \max_{\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1} \mathbf{y}^T \mathbf{B} \mathbf{x}$. We discretize the unit balls in \mathbb{R}^n and \mathbb{R}^m by finite ϵ -nets, where $\epsilon \in (0, 1)$. An ϵ -net of a set \mathcal{K} is the subset \mathcal{N} such that for all $\mathbf{x} \in \mathcal{K}$, there exists $\mathbf{y} \in \mathcal{N}$ such that $\|\mathbf{x} - \mathbf{y}\|_2 \leq \epsilon$. Using a construction proof, we can prove that there exists a finite ϵ -net of the unit ball in \mathbb{R}^n with the cardinality of no more than $\left(\frac{2}{\epsilon} + 1\right)^n$. Let \mathcal{N} and \mathcal{M} be the finite ϵ -nets of the unit balls in \mathbb{R}^n and \mathbb{R}^m with minimum cardinality, respectively. Applying the triangle inequality, we have:

$$\|\mathbf{B}\| \leq \frac{1}{(1-\epsilon)^2} \max_{\mathbf{x} \in \mathcal{N}, \mathbf{y} \in \mathcal{M}} \mathbf{y}^T \mathbf{B} \mathbf{x}.$$

We can bound the tail probability $P(\|\mathbf{B}\| \geq u)$ as follows.

$$\mathbb{P}(\|\mathbf{B}\| \geq u) \leq \left(\frac{2}{\epsilon} + 1\right)^{m+n} \max_{\mathbf{x} \in \mathcal{N}, \mathbf{y} \in \mathcal{M}} \mathbb{P}(\mathbf{y}^T \mathbf{B} \mathbf{x} \geq (1-\epsilon)^2 u).$$

We have, b_{ij} are independent b -subgaussian random variables; therefore, $\mathbf{y}^T \mathbf{B} \mathbf{x} = \sum_{i=1}^m \sum_{j=1}^n (x_j y_i) b_{ij}$

is also a b -subgaussian random variable since $\sum_{i=1}^m \sum_{j=1}^n (x_j y_i)^2 = \|\mathbf{x}\|_2^2 \|\mathbf{y}\|_2^2 = 1$. Thus we have:

$$\mathbb{P}(\|\mathbf{B}\| \geq u) \leq \left(\frac{2}{\epsilon} + 1\right)^{m+n} e^{-\frac{(1-\epsilon)^4 u^2}{2b^2}}.$$

Letting $\epsilon = 1/3$, we obtain the inequality

$$\mathbb{P}(\|\mathbf{B}\| \geq u) \leq e^{-\left(\frac{8u^2}{81b^2} - (\log 7)(m+n)\right)}$$

(ii) We have: $\mathbf{y}^T \mathbf{C} \mathbf{B} \mathbf{x} = (\mathbf{C}^T \mathbf{y})^T \mathbf{B} \mathbf{x}$, thus:

$$\mathbb{P}(\mathbf{y}^T \mathbf{C} \mathbf{B} \mathbf{x} \geq u) \leq e^{-\frac{u^2}{2b^2 \|\mathbf{C}^T \mathbf{y}\|_2^2 \|\mathbf{x}\|_2^2}} \leq e^{-\frac{u^2}{2b^2 \|\mathbf{C}\|^2}}.$$

Applying similar arguments, we can then obtain the inequality (ii) of the Lemma. □

We now turn to the proof of the main theorem. We now would like to find conditions on θ and the constants so that Problem (1) has an optimal solution \mathbf{X} of the form

$$\mathbf{X} = \begin{pmatrix} \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix},$$

where $\mathbf{u}_1 > \mathbf{0}$, $\|\mathbf{u}_1\|_2 = 1$, and $\mathbf{v}_1 > \mathbf{0}$, $\|\mathbf{v}_1\|_2 = 1$. If \mathbf{u}_1 and \mathbf{v}_1 are determined, σ_1 can be easily calculated in order to satisfy the condition $\langle \mathbf{A}, \mathbf{X} \rangle = 1$ of the optimal solution. Thus the main task is to find \mathbf{u}_1 and \mathbf{v}_1 if they exist. We will construct them using optimality conditions derived in the previous section for Problem (6) (equivalent to Problem (1)) and its dual, Problem (7). Defining $\mathbf{u} = [\mathbf{u}_1; \mathbf{0}] \in \mathbb{R}^m$ and $\mathbf{v} = [\mathbf{v}_1; \mathbf{0}] \in \mathbb{R}^n$, we can then write the optimality conditions as follows:

There exists \mathbf{Y} and \mathbf{Z} such that $\mathbf{A} = \mathbf{Y} + \mathbf{Z}$ and

$$\mathbf{Y} = \|\mathbf{A}\|_\theta^* (\mathbf{u} \mathbf{v}^T + \mathbf{W}), \quad \mathbf{Z} = \theta \|\mathbf{A}\|_\theta^* \mathbf{V},$$

where $\|\mathbf{W}\| \leq 1$, $\mathbf{W}^T \mathbf{u} = \mathbf{0}$, $\mathbf{W} \mathbf{v} = \mathbf{0}$, and $\|\mathbf{V}\|_\infty \leq 1$, $\mathbf{V}_{11} = \mathbf{e}_M \mathbf{e}_N^T$.

These conditions come from the properties of the subgradient $\partial \|\mathbf{Y}\|$ and $\partial \|\mathbf{Z}\|_\infty$ and the fact that \mathbf{X} belongs to these sets (up to appropriate scaling factors).

In the following analysis, we will construct (\mathbf{V}, \mathbf{W}) so that the optimality conditions are satisfied. The entries of these matrices will be constructed separately for the four subblocks of \mathbf{A} , starting with the (1,1) block. Breaking the equation $\mathbf{Y} + \mathbf{Z} = \mathbf{A}$ into blocks and scaling by $1/\|\mathbf{A}\|_\theta^*$, we obtain the following more detailed optimality conditions.

$$\mathbf{u}_1 \mathbf{v}_1^T + \mathbf{W}_{11} + \theta \mathbf{e}_M \mathbf{e}_N^T = \mathbf{A}_{11} / \|\mathbf{A}\|_\theta^* = (\sigma_0 \mathbf{u}_0 \mathbf{v}_0^T + \mathbf{R}_{11}) / \|\mathbf{A}\|_\theta^*, \quad (16)$$

$$\mathbf{W}_{ij} + \theta \mathbf{V}_{ij} = \mathbf{A}_{ij} / \|\mathbf{A}\|_\theta^* = \mathbf{R}_{ij} / \|\mathbf{A}\|_\theta^*, \quad (17)$$

where the second line applies to (i, j) equal to (1,2), (2,1) and (2,2). Following this block matrix notation, the remaining optimality conditions to be established are $\mathbf{W}_{11}^T \mathbf{u}_1 = \mathbf{0}$, $\mathbf{W}_{21}^T \mathbf{u}_1 = \mathbf{0}$, $\mathbf{W}_{11} \mathbf{v}_1 =$

$\mathbf{0}$, $\mathbf{W}_{12}\mathbf{v}_1 = \mathbf{0}$, $\|\mathbf{V}\|_\infty \leq 1$, $\|\mathbf{W}\| \leq 1$. We shall establish the latter inequality by proving more specifically that $\|\mathbf{W}_{ij}\| \leq 1/2$ for $(i, j) \in \{1, 2\} \times \{1, 2\}$.

We begin with the (1,1) block of this equation. The conditions $\|\mathbf{W}_{11}\| \leq 1/2$, $\mathbf{W}_{11}^T \mathbf{u}_1 = \mathbf{0}$, $\mathbf{W}_{11} \mathbf{v}_1 = \mathbf{0}$ imply that $\|\mathbf{u}_1 \mathbf{v}_1^T + \mathbf{W}_{11}\| = 1$. This is because the dominant singular triple of $\mathbf{u}_1 \mathbf{v}_1^T + \mathbf{W}_{11}$ must be $(1, \mathbf{u}_1, \mathbf{v}_1)$ by the conditions. Equivalent to $\|\mathbf{u}_1 \mathbf{v}_1^T + \mathbf{W}_{11}\| = 1$ is

$$\left\| \frac{1}{\|\mathbf{A}\|_\theta^*} \mathbf{A}_{11} - \theta \mathbf{V}_{11} \right\| = 1,$$

where, as noted earlier, we are required to take $\mathbf{V}_{11} = \mathbf{e}_M \mathbf{e}_N^T$.

Thus the first necessary condition for \mathbf{X} to be the optimal solution is that there exists $\lambda > 0$ such that $f(\lambda) = \|\lambda \mathbf{A}_{11} - \theta \mathbf{V}_{11}\| = 1$. If such a λ is identified, then \mathbf{u}_1 and \mathbf{v}_1 can be easily found since $\mathbf{u}_1 \mathbf{v}_1^T$ is the rank-one approximation of $\lambda \mathbf{A}_{11} - \theta \mathbf{V}_{11}$. Note that the nonnegativity of \mathbf{u}_1 and \mathbf{v}_1 will require additional conditions which will be discussed later. We have:

$$\begin{aligned} \lambda \mathbf{A}_{11} - \theta \mathbf{V}_{11} &= \lambda [\sigma_0 \mathbf{u}_0 \mathbf{v}_0^T + \mathbf{R}_{11}] - \theta \mathbf{e}_M \mathbf{e}_N^T \\ &= \lambda [\sigma_0 (\mathbf{e}_M + \mathbf{p})(\mathbf{e}_N + \mathbf{q})^T + \mathbf{R}_{11}] - \theta \mathbf{e}_M \mathbf{e}_N^T \\ &= (\lambda \sigma_0 - \theta) \mathbf{e}_M \mathbf{e}_N^T + \lambda [\sigma_0 (\mathbf{e}_M \mathbf{q}^T + \mathbf{p} \mathbf{e}_N^T + \mathbf{p} \mathbf{q}^T) + \mathbf{R}_{11}] \\ &= [\lambda \sigma_0 (1 + c_3) - \theta] \mathbf{e}_M \mathbf{e}_N^T + \lambda [\sigma_0 (\mathbf{e}_M \mathbf{q}^T + \mathbf{p} \mathbf{e}_N^T + \mathbf{p} \mathbf{q}^T) + (\mathbf{R}_{11} - c_3 \sigma_0 \mathbf{e}_M \mathbf{e}_N^T)]. \end{aligned}$$

We have: $f(\lambda) \rightarrow +\infty$ when $\lambda \rightarrow +\infty$ since $\mathbf{A}_{11} \neq \mathbf{0}$. Now define $\lambda_0 = \frac{\theta}{\sigma_0(1+c_3)}$ to make the first term vanish, yielding

$$\begin{aligned} \lambda_0 \mathbf{A}_{11} - \theta \mathbf{V}_{11} &= \lambda_0 [\sigma_0 (\mathbf{e}_M \mathbf{q}^T + \mathbf{p} \mathbf{e}_N^T + \mathbf{p} \mathbf{q}^T) + (\mathbf{R}_{11} - c_3 \sigma_0 \mathbf{e}_M \mathbf{e}_N^T)] \\ &= \lambda_0 [\sigma_0 \mathbf{P} + \mathbf{Q}], \end{aligned}$$

where $\mathbf{P} = \mathbf{e}_M \mathbf{q}^T + \mathbf{p} \mathbf{e}_N^T + \mathbf{p} \mathbf{q}^T$ and $\mathbf{Q} = \mathbf{R}_{11} - c_3 \sigma_0 \mathbf{e}_M \mathbf{e}_N^T$. We now bound the spectral norm of \mathbf{P} and \mathbf{Q} as follows.

$$\begin{aligned} \|\mathbf{P}\| &\leq \|\mathbf{e}_M \mathbf{q}^T\| + \|\mathbf{p} \mathbf{e}_N^T\| + \|\mathbf{p} \mathbf{q}^T\| \\ &= \|\mathbf{e}\|_2 \|\mathbf{q}\|_2 + \|\mathbf{p}\|_2 \|\mathbf{e}_N\|_2 + \|\mathbf{p}\|_2 \|\mathbf{q}\|_2 \\ &\leq c_1 \sqrt{MN} + c_2 \sqrt{MN} + c_1 c_2 \sqrt{MN} \\ &= (c_1 + c_2 + c_1 c_2) \sqrt{MN}. \end{aligned}$$

Matrix \mathbf{Q}/σ_0 is random with i.i.d. elements that are b -subgaussian. Thus by Lemma 5(i),

$$\mathbb{P}(\|\mathbf{Q}\| \geq u \sigma_0) \leq \exp\left(-\left(\frac{8u^2}{81b^2} - (\log 7)(M+N)\right)\right),$$

for any $u > 0$. Let us fix $u = (MN)^{3/8}$ to obtain

$$\mathbb{P}\left(\|\mathbf{Q}\| \geq \sigma_0(MN)^{3/8}\right) \leq \exp\left(-\left(\frac{(MN)^{3/4}}{81b^2} - (\log 7)(M+N)\right)\right). \quad (18)$$

For the remainder of this analysis, we will impose the assumption that the event in (18) does not happen. At the end of the proof the right-hand side (18) will be one of the terms in the failure probability of identifying the optimal \mathbf{X} .

Thus, $\|\mathbf{Q}\| \leq o(1)\sigma_0\sqrt{MN}$, so applying the triangle inequality,

$$\begin{aligned} \|\sigma_0\mathbf{P} + \mathbf{Q}\| &\leq \sigma_0(c_1 + c_2 + c_1c_2 + o(1))\sqrt{MN} \\ &\leq \sigma_0c_5\sqrt{MN}. \end{aligned} \quad (19)$$

by (13). (The strict inequality ' $>$ ' in (13) is used in order to absorb the $o(1)$ term.) Therefore,

$$\begin{aligned} f(\lambda_0) &= \lambda_0\|\sigma_0\mathbf{P} + \mathbf{Q}\| \\ &\leq \frac{c_5\theta\sqrt{MN}}{1 + c_3}. \end{aligned}$$

Thus with high probability, $f(\lambda_0) \leq 1$ if

$$\theta < \left(\frac{1 + c_3}{c_5}\right) \frac{1}{\sqrt{MN}}, \quad (20)$$

Inequality (20) is a consequence of (14) stated in the theorem. This inequality implies $f(\lambda_0) \leq 1$, and, due to the continuity of function f , there exists $\lambda^* \geq \lambda_0$ such that $f(\lambda^*) = 1$. We will prove that under some additional conditions, this value λ^* satisfies all other optimality conditions of Problem (1) and indeed $\|\mathbf{A}\|_\theta^* = \frac{1}{\lambda^*}$.

Let us recall that $\|\lambda^*\mathbf{A}_{11} - \theta\mathbf{V}_{11}\| = 1$, i.e., $\|(\lambda^*\sigma_0(1 + c_3) - \theta)\mathbf{e}_M\mathbf{e}_N^T + \lambda^*(\sigma_0\mathbf{P} + \mathbf{Q})\| = 1$. Applying the fact that $\|\mathbf{e}_M\mathbf{e}_N^T\| = \sqrt{MN}$ and the triangle inequality twice to this equation yields

$$[\lambda^*\sigma_0(1 + c_3) - \theta]\sqrt{MN} - \lambda^*\|\sigma_0\mathbf{P} + \mathbf{Q}\| \leq 1 \leq [\lambda^*\sigma_0(1 + c_3) - \theta]\sqrt{MN} + \lambda^*\|\sigma_0\mathbf{P} + \mathbf{Q}\|.$$

Applying (19) yields

$$[\lambda^*\sigma_0(1 + c_3 - c_5) - \theta]\sqrt{MN} \leq 1 \leq [\lambda^*\sigma_0(1 + c_3 + c_5) - \theta]\sqrt{MN}.$$

Rearranging this chain of inequalities and using the fact that $1 + c_3 - c_5 > 0$, which follows from (12) stated in the theorem, yields

$$\frac{1 + \theta\sqrt{MN}}{\sigma_0(1 + c_3 + c_5)\sqrt{MN}} \leq \lambda^* \leq \frac{1 + \theta\sqrt{MN}}{\sigma_0(1 + c_3 - c_5)\sqrt{MN}}, \quad (21)$$

with high probability.

We wish to establish that $\lambda^* \sigma_0 - \theta \geq 0$. Using the left inequality in (21) yields:

$$\begin{aligned} \lambda^* \sigma_0 - \theta &\geq \frac{1 + \theta \sqrt{MN}}{(1 + c_3 + c_5) \sqrt{MN}} - \theta \\ &= \frac{1 - (c_3 + c_5) \theta \sqrt{MN}}{(1 + c_3 + c_5) \sqrt{MN}}. \end{aligned}$$

Thus, nonnegativity of $\lambda^* \sigma_0 - \theta$ is implied by the inequality $\theta \leq 1/((c_3 + c_5) \sqrt{MN})$, which is a consequence of assumption (14).

Since $\lambda^* \sigma_0 - \theta \geq 0$,

$$\lambda^* \mathbf{A}_{11} - \theta \mathbf{V}_{11} = (\lambda^* \sigma_0 - \theta) \mathbf{e}_M \mathbf{e}_N^T + \lambda^* (\sigma_0 \mathbf{P} + \mathbf{R}_{11}) > \mathbf{0}.$$

Applying the Perron-Frobenius theorem, we obtain the positivity of \mathbf{u}_1 and \mathbf{v}_1 .

We also need $\|\mathbf{W}_{11}\| \leq 1/2$. Recall $\|\mathbf{W}_{11}\| = \sigma_2(\lambda^* \mathbf{A}_{11} - \theta \mathbf{V}_{11})$, the second largest singular value of $\lambda^* \mathbf{A}_{11} - \theta \mathbf{V}_{11}$, since $\lambda^* \mathbf{A}_{11} - \theta \mathbf{V}_{11} = \mathbf{u}_1 \mathbf{v}_1^T + \mathbf{W}_{11}$. Using the well-known fact that

$$\sigma_2(\mathbf{A}) = \min\{\|\mathbf{A} - \mathbf{S}\| : \text{rank}(\mathbf{S}) \leq 1\},$$

we obtain

$$\|\mathbf{W}_{11}\| \leq \|\lambda^* (\sigma_0 \mathbf{P} + \mathbf{Q})\|.$$

Here we selected \mathbf{S} to be $[\lambda^* \sigma_0 (1 + c_3) - \theta] \mathbf{e}_M \mathbf{e}_N^T$. With high probability, we obtain the bound $\|\mathbf{W}_{11}\| \leq \lambda^* \sigma_0 c_5 \sqrt{MN}$ from (19).

Using the upper bound on λ^* from (21), we have:

$$\|\mathbf{W}_{11}\| \leq \frac{(1 + \theta \sqrt{MN}) c_5}{1 + c_3 - c_5}.$$

In order to obtain $\|\mathbf{W}_{11}\| \leq 1/2$, a sufficient condition is

$$\frac{(1 + \theta \sqrt{MN}) c_5}{1 + c_3 - c_5} \leq \frac{1}{2} \tag{22}$$

which is rearranged as

$$\theta \leq \left(\frac{1 + c_3 - 3c_5}{2c_5} \right) \frac{1}{\sqrt{MN}}.$$

The latter inequality follows from (14); the numerator of the right-hand side is positive by (12).

Turning to (17) when $(i, j) = (2, 2)$, we need to find \mathbf{W}_{22} and \mathbf{V}_{22} that satisfy

$$\lambda^* \mathbf{R}_{22} = \mathbf{W}_{22} + \theta \mathbf{V}_{22},$$

$\|\mathbf{W}_{22}\| \leq 1/2$, and $\|\mathbf{V}_{22}\|_\infty \leq 1$. Consider the assignment $\mathbf{V}_{22} = \frac{\lambda^* \sigma_0 c_3}{\theta} \mathbf{e}_{m-M} \mathbf{e}_{n-N}^T$ and $\mathbf{W}_{22} = \lambda^* \mathbf{R}_{22} - \theta \mathbf{V}_{22}$. The coefficient $\lambda^* \sigma_0 c_3 / \theta$ is chosen for the definition of \mathbf{V}_{22} so that the entries of the remainder term \mathbf{W}_{22} have mean zero.

The requirement $\|\mathbf{V}_{22}\|_\infty \leq 1$ is satisfied if and only if $\lambda^* \sigma_0 c_3 / \theta \leq 1$. Because of the upper bound on λ^* established by (21), this requirement is satisfied if

$$\frac{(1 + \theta \sqrt{MN}) c_3}{\theta (1 + c_3 - c_5) \sqrt{MN}} \leq 1 \Leftrightarrow \theta \geq \left(\frac{c_3}{1 - c_5} \right) \frac{1}{\sqrt{MN}}, \quad c_5 < 1.$$

This inequality is assured by (12) and (15). (In particular, (12) implies $c_5 < 1$.)

To bound $\|\mathbf{W}_{22}\|$, consider $\mathbf{W}_{22}/(\lambda^* \sigma_0)$, which is a random matrix with i.i.d. elements that are b -subgaussian. Applying Lemma 5(i) to $\mathbf{W}_{22}/(\lambda^* \sigma_0)$ and taking $u = 1/(2\lambda^* \sigma_0)$ yields

$$\mathbb{P}(\|\mathbf{W}_{22}\| \geq 1/2) \leq \exp\left(-\left(\frac{2}{81b^2(\lambda^* \sigma_0)^2} - (\log 7)(m - M + n - N)\right)\right).$$

Use the upper bound on λ^* from (21) to obtain the following tail bound:

$$\mathbb{P}\left(\|\mathbf{W}_{22}\| \geq \frac{1}{2}\right) \leq \exp\left(-\left(\frac{2(1 + c_3 - c_5)^2}{81b^2(1 + \theta \sqrt{MN})^2} MN - (\log 7)(m - M + n - N)\right)\right).$$

From (22) we obtain

$$\frac{(1 + c_3 - c_5)^2}{(1 + \theta \sqrt{MN})^2} \geq 4c_5^2 \tag{23}$$

hence

$$\mathbb{P}\left(\|\mathbf{W}_{22}\| \geq \frac{1}{2}\right) \leq \exp\left(-\left(\frac{8c_5^2 MN}{81b^2} - (\log 7)(m - M + n - N)\right)\right). \tag{24}$$

Now consider (17) when $(i, j) = (1, 2)$. Again we need to find \mathbf{W}_{12} and \mathbf{V}_{12} such that

$$\lambda^* \mathbf{R}_{12} = \mathbf{W}_{12} + \theta \mathbf{V}_{12},$$

$\|\mathbf{W}_{12}\| \leq 1/2$, $\|\mathbf{V}_{12}\|_\infty \leq 1$, and $\mathbf{W}_{12}^T \mathbf{u}_1 = \mathbf{0}$. We construct \mathbf{W}_{12} and \mathbf{V}_{12} column by column as follows:

$$\mathbf{V}_{12}(:, i) = \frac{\lambda^* \mathbf{R}_{12}(:, i)^T \mathbf{u}_1}{\theta \|\mathbf{u}_1\|_1} \mathbf{e}_M, \quad \mathbf{W}_{12}(:, i) = \lambda^* \mathbf{R}_{12}(:, i) - \theta \mathbf{V}_{12}(:, i)$$

all $i = 1, \dots, n - N$. By construction we have $\mathbf{W}_{12}(:, i)^T \mathbf{u}_1 = 0$ for all $i = 1, \dots, n - N$. Now consider the requirement that $\|\mathbf{V}_{12}(:, i)\|_\infty \leq 1$ for all $i = 1, \dots, n - N$. The requirement is equivalent to

$$\frac{\lambda^* \mathbf{R}_{12}(:, i)^T \mathbf{u}_1}{\theta \|\mathbf{u}_1\|_1} \leq 1$$

for all $i = 1, \dots, n - N$. Subtract $\lambda^* c_3 \sigma_0 / \theta$ from both sides and apply the identity $\mathbf{e}_M^T \mathbf{u}_1 = \|\mathbf{u}_1\|_1$ to obtain

$$\frac{\lambda^* (\mathbf{R}_{12}(:, i)^T - c_3 \sigma_0 \mathbf{e}_M^T) \mathbf{u}_1}{\theta \|\mathbf{u}_1\|_1} \leq 1 - \frac{\lambda^* c_3 \sigma_0}{\theta}.$$

We will establish this inequality in two steps. First, we establish that $\lambda^*c_3\sigma_0/\theta \leq 1/2$. Because of (21), it suffices to establish that

$$\frac{c_3(1 + \theta\sqrt{MN})}{\theta(1 + c_3 - c_5)\sqrt{MN}} \leq \frac{1}{2}. \quad (25)$$

This can be rearranged into

$$\theta \geq \frac{2c_3}{(1 - c_3 - c_5)\sqrt{MN}},$$

which follows from (15).

Second, we establish that with probability exponentially close to 1,

$$\frac{\lambda^*(\mathbf{R}_{12}(:, i)^T - c_3\sigma_0\mathbf{e}_M^T)\mathbf{u}_1}{\theta\|\mathbf{u}_1\|_1} \leq \frac{1}{2}.$$

Notice that $r_{ji}/\sigma_0 - c_3$ is b -subgaussian; thus,

$$\frac{1}{\sigma_0}(\mathbf{R}_{12}(:, i) - c_3\sigma_0\mathbf{e}_M)^T\mathbf{u}_1 = \frac{1}{\sigma_0}\sum_{j=1}^M u_1(j)(R_{12}(j, i) - c_3\sigma_0)$$

is also b -subgaussian since $\|\mathbf{u}_1\|_2 = 1$. Thus, by (10), taking $x = (\mathbf{R}_{12}(:, i) - c_3\sigma_0\mathbf{e}_M)^T\mathbf{u}_1/\sigma_0$ and taking $t = \theta\|\mathbf{u}_1\|_1/(2\lambda^*\sigma_0)$,

$$\begin{aligned} \mathbb{P}\left(\frac{\lambda^*(\mathbf{R}_{12}(:, i)^T - c_3\sigma_0\mathbf{e}_M^T)\mathbf{u}_1}{\theta\|\mathbf{u}_1\|_1} > \frac{1}{2}\right) &\leq \exp(-\theta^2\|\mathbf{u}_1\|_1^2/(8b^2(\lambda^*\sigma_0)^2)). \\ &\leq \exp(-c_3^2\|\mathbf{u}_1\|_1^2/(2b^2)). \end{aligned} \quad (26)$$

since, as noted above $\theta/(\lambda^*\sigma_0) \geq 2c_3$.

To proceed, we now need a lower bound for $\|\mathbf{u}_1\|_1$. Let \mathbf{F} denote $\lambda^*\mathbf{A}_{11} - \theta\mathbf{V}_{11}$, which is equal to $[\lambda^*\sigma_0(1 + c_3) - \theta]\mathbf{e}_M\mathbf{e}_N^T + \lambda^*(\sigma_0\mathbf{P} + \mathbf{Q})$. We know that \mathbf{u}_1 is the first (left) singular vector of \mathbf{F} . Letting $\mathbf{X}_0 = [\lambda^*\sigma_0(1 + c_3) - \theta]\mathbf{e}_M\mathbf{e}_N^T$ and $\mathbf{E} = \lambda^*(\sigma_0\mathbf{P} + \mathbf{Q})$, we then have $\mathbf{F} = \mathbf{X}_0 + \mathbf{E}$, and \mathbf{X}_0 is a rank-one matrix with a single nonzero singular value equal to $(\lambda^*\sigma_0(1 + c_3) - \theta)\sqrt{MN}$ and with left singular vector \mathbf{e}_M/\sqrt{M} and right singular vector \mathbf{e}_N/\sqrt{N} . Furthermore, since $\|\mathbf{F}\| = 1$, we know that the singular value of \mathbf{X}_0 is at least $1 - \|\mathbf{E}\|$ by Corollary 8.6.2 of [11]. Thus, by Theorem 8.6.5 of [11],

$$\left\|\mathbf{u}_1 - \frac{\mathbf{e}_M}{\sqrt{M}}\right\| \leq \frac{4\|\mathbf{E}\|}{1 - \|\mathbf{E}\|} \leq 4/5, \quad (27)$$

provided that $\|\mathbf{E}\| \leq 1/6$. Thus, the next step in the analysis is to show that $\|\mathbf{E}\| \leq 1/6$. This follows from the following sequence of inequalities:

$$\begin{aligned} \|\mathbf{E}\| &= \lambda^*\|\sigma_0\mathbf{P} + \mathbf{Q}\| \\ &\leq \frac{c_5(1 + \theta\sqrt{MN})}{1 + c_3 - c_5} \\ &\leq 1/6, \end{aligned}$$

where the second line holds with high probability according to (19) and the third line follows from (22) and (12).

Thus, we have established that $\|\mathbf{E}\| \leq 1/6$ with high probability, which in turn implies that

$$\begin{aligned}
\|\mathbf{u}_1\|_1 &= \mathbf{e}_M^T \mathbf{u}_1 \\
&= (\mathbf{e}_M - M^{1/2} \mathbf{u}_1 + M^{1/2} \mathbf{u}_1)^T \mathbf{u}_1 \\
&\geq M^{1/2} \mathbf{u}_1^T \mathbf{u}_1 - |(\mathbf{e}_M - M^{1/2} \mathbf{u}_1)^T \mathbf{u}_1| \\
&\geq M^{1/2} \mathbf{u}_1^T \mathbf{u}_1 - \|\mathbf{e}_M - M^{1/2} \mathbf{u}_1\| \cdot \|\mathbf{u}_1\| \\
&= M^{1/2} - M^{1/2} \|M^{-1/2} \mathbf{e}_M - \mathbf{u}_1\| \\
&\geq M^{1/2} - (4/5)M^{1/2},
\end{aligned}$$

where the last line is obtained from (27). This gives a lower bound of $M^{1/2}/5$ on $\|\mathbf{u}\|_1$. Thus, substituting this into (26) yields

$$\mathbb{P}\left(\frac{\lambda^*(\mathbf{R}_{12}(:,i)^T - c_3 \sigma_0 \mathbf{e}_M) \mathbf{u}_1}{\theta \|\mathbf{u}_1\|_1} > \frac{1}{2}\right) \leq \exp(-\sigma_0^2 c_3^2 M / (50b^2)).$$

This shows that one column of \mathbf{V}_{12} exceeds norm $1/2$ with exponentially small probability. Applying the union bound over all the columns, we find

$$\mathbb{P}(\|\mathbf{V}_{12}\|_\infty \geq 1) \leq (n - N) \exp(-\sigma_0^2 c_3^2 M / (50b^2)). \tag{28}$$

Thus, we have established that $\|\mathbf{V}_{12}\|_\infty \leq 1$ with probability exponentially close to 1.

We now consider the matrix \mathbf{W}_{12} , which can be written as $\mathbf{W}_{12} = \lambda^* \mathbf{D} \mathbf{R}_{12}$, where $\mathbf{D} \in \mathbb{R}^{M \times M}$ is given by

$$\mathbf{D} = \mathbf{I} - \frac{1}{\|\mathbf{u}_1\|_1} \mathbf{e}_M \mathbf{u}_1^T.$$

Notice that we can equivalently write

$$\mathbf{W}_{12} = \lambda^* \mathbf{D} (\mathbf{R}_{12} - c_3 \sigma_0 \mathbf{e}_M \mathbf{e}_{n-N}^T),$$

since $\mathbf{D} \mathbf{e}_M = \mathbf{0}$. The matrix $\mathbf{R}_{12} - c_3 \sigma_0 \mathbf{e}_M \mathbf{e}_{n-N}^T$ is a subgaussian matrix scaled by σ_0 . Furthermore,

$$\|\mathbf{D}_{12}\| \leq 1 + \frac{\sqrt{M}}{\|\mathbf{u}\|_1} \leq 6,$$

with high probability, since $\|\mathbf{u}_1\|_1 \geq M^{1/2}/5$. Thus, Lemma 5(ii) applied to $\mathbf{W}_{12}/(\lambda^* \sigma_0)$, taking $u = 1/(2\lambda^* \sigma_0)$, yields

$$\mathbb{P}(\|\mathbf{W}_{12}\| \geq 1/2) \leq \exp\left(-\left(\frac{2}{36 \cdot 81b^2 \sigma_0^2 (\lambda^*)^2} - (\log 7)(M + n - N)\right)\right).$$

Apply the upper bound on λ^* from (21) to obtain

$$\mathbb{P}(\|\mathbf{W}_{12}\| \geq 1/2) \leq \exp\left(-\left(\frac{2(1+c_3-c_5)^2 MN}{36 \cdot 81b^2(1+\theta\sqrt{MN})^2} - (\log 7)(M+n-N)\right)\right).$$

Now finally we apply (23) to obtain

$$\mathbb{P}(\|\mathbf{W}_{12}\| \geq 1/2) \leq \exp\left(-\left(\frac{8c_5^2 MN}{36 \cdot 81b^2} - (\log 7)(M+n-N)\right)\right). \quad (29)$$

The same construction and analysis applies to \mathbf{V}_{21} and \mathbf{W}_{21} , and the same results are obtained except with the roles of (M, m) and (N, n) interchanged. Thus,

$$\mathbb{P}(\|\mathbf{V}_{21}\|_\infty \geq 1) \leq (m-M) \exp(-\sigma_0^2 c_3^2 N / (50b^2)), \quad (30)$$

and

$$\mathbb{P}(\|\mathbf{W}_{21}\| \geq 1/2) \leq \exp\left(-\left(\frac{8c_5^2 MN}{36 \cdot 81b^2} - (\log 7)(N+m-M)\right)\right). \quad (31)$$

From the analysis of all four blocks of \mathbf{V} and \mathbf{W} , we have:

$$\|\mathbf{V}\|_\infty = \max\{\|\mathbf{V}_{11}\|_\infty, \|\mathbf{V}_{12}\|_\infty, \|\mathbf{V}_{21}\|_\infty, \|\mathbf{V}_{22}\|_\infty\} \leq 1,$$

where $\mathbf{V}_{11} = \mathbf{e}_M \mathbf{e}_N^T$. With a high probability, we also have $\|\mathbf{W}\| \leq 1$ since

$$\|\mathbf{W}\|^2 \leq \|\mathbf{W}_{11}\|^2 + \|\mathbf{W}_{12}\|^2 + \|\mathbf{W}_{22}\|^2 + \|\mathbf{W}_{21}\|^2 \leq 1.$$

By the union bound, the probability of failure of the main result is at most the sum of the probabilities of the failure at each step. Therefore, the failure of the convex relaxation to find the claimed optimal \mathbf{X} is at most the sum of the right-hand sides of (18), (24), (30), (29), (28), and (31). We require these probabilities to be exponentially small. We assure that (18) is exponentially small by requiring

$$MN \geq k_1(M+N)^{4/3} \quad (32)$$

where

$$k_1 > ((\log 7)81b^2)^{4/3}. \quad (33)$$

Next, all of (24), (29), (31) are exponentially small provided that

$$MN \geq k_2(m+n) \quad (34)$$

where

$$k_2 > \frac{(\log 7)36 \cdot 81b^2}{8c_5^2}. \quad (35)$$

Finally, to ensure that (28) and (30) tend to 0 exponentially fast requires that M grow as fast as $\Omega(\log(n-N))$ and similarly N grows as fast as $\Omega(\log(m-M))$, but this is already a consequence of (32) and (34).

5 Conclusions

We have shown that a convex relaxation can find a large, approximately rank-one submatrix of a much larger noisy matrix provided that the dimensions of the larger matrix are no larger than the square of the dimensions of the smaller matrix, and provided certain upper bounds are satisfied on the level of the noise.

It is interesting to note that our result also applies to the maximum biclique problem, which was introduced in Section 1 as a special case of LAROS. In particular, if G is a bipartite graph (U, V, E) containing a biclique given by $U^* \times V^*$, where $|U| = m$, $|V| = n$, $|U^*| = M$, $|V^*| = N$, and if the remaining edges of E (i.e., those not in $U^* \times V^*$) are inserted at random with probability $1/2$, then the U -to- V adjacency matrix has the form (11) in which $\sigma = 1$, $c_1 = c_2 = 0$, $c_3 = 1/2$, $b = 1/(8 \log 2)^{1/2}$. (This is not quite correct since in this case $\mathbf{R}_{11} = \mathbf{0}$. However, our analysis covers this case as well.) Thus, our algorithm with parameter $\theta = O(1/(MN)^{1/2})$ finds the planted biclique when $M \sim N$, $m \sim n$, and $M \geq \Omega(m^{1/2})$. The same result was obtained earlier by Ames and Vavasis [1] using a different convex relaxation. Theirs has the advantage that M, N do not need to be known or estimated in advance, but ours solves a more general class of problems.

References

- [1] B. Ames and S. Vavasis. Nuclear norm minimization for the planted clique and biclique problems. Under review, Math. Prog., URL: <http://arxiv.org/abs/0901.3348>, 2009.
- [2] N. Asgarian and R. Greiner. Using rank-1 biclusters to classify microarray data. Technical report, Department of Computing Science and the Alberta Ingenuity Center for Machine Learning, University of Alberta, Edmonton, AB, Canada, 2008.
- [3] S. Bergmann, J. Ihmels, and N. Barkai. Iterative signature algorithm for the analysis of large-scale gene expression data. *Physical Review E*, 67:031902, 2003.
- [4] Michael Biggs, Ali Ghodsi, and Stephen A. Vavasis. Nonnegative matrix factorization via rank-one downdating. In *Proceedings of the 2008 International Conference on Machine Learning*, 2008. Proceedings published online at <http://icml2008.cs.helsinki.fi/abstracts.shtml>, Preliminary version of the full paper in arxiv.org, 0808.0120.

- [5] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, UK, 2004.
- [6] E. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772, 2009.
- [7] V. Chandrasekaran, S. Sanghavi, P.A. Parrilo, and A. Willsky. Rank-sparsity incoherence for matrix decomposition. See <http://arxiv.org/abs/0906.2220>, 2009.
- [8] J. Cohen and U. Rothblum. Nonnegative ranks, decompositions and factorizations of nonnegative matrices. *Linear Algebra and its Applications*, 190:149–168, 1993.
- [9] M. Fazel. *Matrix Rank Minimization with Applications*. PhD thesis, Stanford University, Stanford, CA, 2001.
- [10] N. Gillis and F. Glineur. Using underapproximations for sparse nonnegative matrix factorization. *Pattern Recognition*, 43:1676–1687, 2010.
- [11] G. H. Golub and C. F. Van Loan. *Matrix Computations, 3rd Edition*. Johns Hopkins University Press, Baltimore, 1996.
- [12] T. Hofmann. Probabilistic latent semantic analysis. In Kathryn B. Laskey and Henri Prade, editors, *UAI '99: Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence, Stockholm, Sweden, July 30-August 1, 1999*, pages 289–296. Morgan Kaufmann, 1999.
- [13] D. Lee and H. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401:788–791, 1999.
- [14] A. Litvak, A. Pajor, M. Rudelson, and N. Tomczak-Jaegermann. Smallest singular values of random matrices and geometry of random polytopes. *Advances in Mathematics*, 195:491–523, 2005.
- [15] R. Peeters. The maximum edge biclique problem is NP-complete. *Discrete Applied Mathematics*, 131:651–654, 2003.
- [16] B. Recht, M. Fazel, and P. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010.
- [17] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, NJ, 1970.

- [18] L. B. Thomas. Problem 73-14, rank factorization of nonnegative matrices by Berman and Plemmons. *SIAM Review*, 16:393–394, 1974.
- [19] Stephen A. Vavasis. On the complexity of nonnegative matrix factorization. *SIAM J. Optim.*, 20(3):1364–1377, 2009.
- [20] K. Ziętak. Properties of linear approximations of matrices in the spectral norm. *Linear Algebra Applications*, 183:41–60, 1993.