

Optimal Exact Least Squares Rank Minimization

Shuo Xiang^{1,2}, Yunzhang Zhu³, Xiaotong Shen³, Jieping Ye^{1,2}

¹Department of Computer Science and Engineering, Arizona State University, AZ 85287

²Center for Evolutionary Medicine and Informatics, The Biodesign Institute, Arizona State University, AZ 85287

³School of Statistics, University of Minnesota, Minneapolis, MN 55455

ABSTRACT

In multivariate analysis, rank minimization emerges when a low-rank structure of matrices is desired as well as a small estimation error. Rank minimization is nonconvex and generally NP-hard, imposing one major challenge. In this paper, we consider a nonconvex least squares formulation, which seeks to minimize the least squares loss function with the rank constraint. Computationally, we develop efficient algorithms to compute a global solution as well as an entire regularization solution path. Theoretically, we show that our method reconstructs the oracle estimator exactly from noisy data. As a result, it recovers the true rank optimally against any method and leads to sharper parameter estimation over its counterpart. Finally, the utility of the proposed method is demonstrated by simulations and image reconstruction from noisy background.

Categories and Subject Descriptors

H.2.8 [Database Management]: Database Applications-Data Mining

General Terms

Algorithms

Keywords

Nonconvex, global optimality, rank minimization

1. INTRODUCTION

In multivariate analysis, estimation of lower-dimensional structures has received attention in statistics, signal processing and machine learning. One type of such structures is the low-rank of matrices [5, 22], where the rank measures the dimension of a multivariate response. Rank minimization approximates multivariate data with the smallest possible rank of matrices. It has many applications in, for instance,

multi-task learning [6, 11], multi-class classification [2], matrix completion [8, 17], collaborative filtering [33, 1], clustering [20, 29], computer vision [35, 21, 18], among others. The central topic this article addresses is least squares rank minimization.

Consider multi-response linear regression in which a k -dimensional response vector \mathbf{z}_i follows

$$\begin{aligned} \mathbf{z}_i &= \Theta^T \mathbf{a}_i + \boldsymbol{\varepsilon}_i; \\ E\boldsymbol{\varepsilon}_i &= \mathbf{0}, \quad \text{Cov}(\boldsymbol{\varepsilon}_i) = \sigma^2 \mathbf{I}_{k \times k}; \quad i = 1, \dots, n, \end{aligned} \quad (1)$$

where \mathbf{a}_i is a p -dimensional design vector, Θ is a $p \times k$ regression parameter matrix, and components of $\boldsymbol{\varepsilon}_i$ are independent. Model (1) reduces to the widely-used linear model in compressed sensing when $k = 1$ and becomes a multivariate autoregressive model with $\mathbf{a}_i = \mathbf{z}_{i-1}$. Denote the rank of Θ as $r(\Theta)$ and rewrite (1) in the matrix form as follows:

$$\mathbf{Z} = \mathbf{A}\Theta + \mathbf{e}, \quad (2)$$

where $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)^T \in \mathbb{R}^{n \times k}$, $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_n)^T \in \mathbb{R}^{n \times p}$ and $\mathbf{e} = (\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_n)^T \in \mathbb{R}^{n \times k}$ are the data, design and error matrices. In (1), we estimate Θ based on n pairs of observation vectors $(\mathbf{a}_i, \mathbf{z}_i)_{i=1}^n$, with *a priori* knowledge that $r(\Theta)$ is relatively small in comparison to $\min(n, k, p)$, where the number of unknown parameters k, p can greatly exceed the sample size n .

Least squares rank minimization, as described, solves

$$\begin{aligned} \min_{\Theta} \quad & \|\mathbf{A}\Theta - \mathbf{Z}\|_F^2 \\ \text{s.t.} \quad & r(\Theta) \leq s, \end{aligned} \quad (3)$$

where $\|\cdot\|_F$ is the Frobenius-norm and s is an integer-valued tuning parameter taking values from $[1, \min(n, k, p)]$. The general rank minimization is nonconvex and NP-hard [23], which is like the L_0 -minimization in univariate analysis. Therefore an exact global solution to (3) is not known as well as its statistical properties, due primarily to discreteness and non-convexity of the rank function.

Estimation under the restriction that $r(\Theta) = r$ has been studied when $n \rightarrow \infty$ with k and p held fixed, see [3, 4, 15, 28, 26]. Two major computational approaches have been proposed for approximating the optimal solution of (3). The first one involves regularization using surrogate function, such as the nuclear-norm, which is a convex envelope of the rank function [13] and can be solved by efficient algorithms [8, 19, 34, 25, 16]. In some cases, the solution of this convex problem coincides with a global minimizer of (3) under certain isometry assumptions [27]. However, these as-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

KDD'12, August 12–16, 2012, Beijing, China.

Copyright 2012 ACM 978-1-4503-1462-6/12/08 ...\$10.00.

sumptions can be strong and difficult to check. Recently, [7] obtained a global minimizer of a regularized version of (3).

The second approach attacks (3) by approximating the rank function iteratively through calculating the largest singular vector using greedy search [30], or by singular value projection (SVP) through a local gradient method [17]. Under weaker isometry assumptions [27, 10, 9] than that of the nuclear-norm approach, these methods guarantee an exact solution of (3) but suffer from the same difficulties as the regularization method [30], although they have achieved promising results on both simulated and real-world data.

Theoretically, some loss error bounds of the first regularization approach are obtained in [24] under Frobenius-norm, and rank selection consistency is established in [7]. Unfortunately, to our best knowledge, whether similar conclusions hold for our formulation (3) remains largely unknown.

In this paper, we have advanced on two fronts. Computationally, we derive a general closed-form for a global minimizer of (3) in Theorem 1, and give a condition under which (3) and its nonconvex regularized counterpart are equivalent with regard to global minimizers, although these two methods are not generally equivalent. Moreover, we develop an efficient algorithm for computing the entire regularization solution path at the cost of computing only one solution for a single regularization parameter. Theoretically, we establish optimality for a global minimizer of (3). More specifically, the proposed method is optimal against any other ones in that it reconstructs the oracle estimator exactly, thus the true rank, under (1). It is important to note that this exact recovery result is a much stronger property than consistency, which is attributed to the discrete nature of the rank function as well as tuning parameter s . Such a result may not be shared by its regularized counterpart with a continuum tuning parameter. In addition, the method enjoys a higher degree of accuracy for parameter estimation than nuclear-norm rank estimation.

After the first draft of this paper was completed, we were aware that [14] and [32] gave an expression of Theorem 1. However, neither paper considered computational and statistical aspects of the solution. Inevitably, some partial overlaps exist between our Theorem 1 and theirs.

The rest of the paper is organized as follows. Section 2 presents a closed-form solution to (3). Section 3 gives an efficient path algorithm for a regularized version of (3). Section 4 is devoted to theoretical investigation, followed by Section 5 discussing methods for tuning. Section 6 presents proofs for all the theorems we develop and Section 7 presents results of empirical evaluations, where several rank minimization methods are compared. Section 8 concludes the paper.

Notation: Table 1 summarizes the notations used in the rest of the paper.

2. PROPOSED METHOD: CLOSED-FORM SOLUTION

This section derives a closed-form solution to (3). The strategy is to simplify (3) through the singular value decomposition (SVD) of matrix \mathbf{A} and utilizing the properties of the rank function. Before proceeding, we present a motivating lemma, known as the Eckart-Young theorem [12].

LEMMA 1. *The best s -rank approximation, in terms of the Frobenius-norm, for a t -rank matrix \mathbf{Z} with $t \geq s$, i.e., \mathbf{a}*

Table 1: Notation used throughout the paper

| Notation | Description |
|-----------------------------|----------------------------------------------------------------------------------------|
| \mathbf{A} | The design matrix |
| \mathbf{e} | Error matrix |
| Θ^0 | The ground truth model |
| $\hat{\Theta}^s$ | Estimator obtained by optimizing (3) |
| $\hat{\Theta}_\lambda^*$ | Estimator obtained by optimizing (7) |
| $r(\mathbf{A})$ | The rank of matrix \mathbf{A} |
| $\mathcal{P}_s(\mathbf{Z})$ | The best s -rank approximation of matrix \mathbf{Z} in terms of the Frobenius-norm |

global minimizer Θ^* of

$$\begin{aligned} \min_{\Theta} \quad & \|\Theta - \mathbf{Z}\|_F^2 \\ \text{s.t.} \quad & r(\Theta) \leq s \end{aligned} \quad (4)$$

is given by $\Theta^* = \mathcal{P}_s(\mathbf{Z}) = \mathbf{U}_z \mathbf{D}_s \mathbf{V}_z^T$, where \mathbf{D}_s consists of the largest s singular values of \mathbf{Z} given the SVD of $\mathbf{Z} = \mathbf{U}_z \mathbf{D} \mathbf{V}_z^T$.

Intuitively $\mathcal{P}_s(\mathbf{Z})$ may be viewed as a projection of \mathbf{Z} onto a set of matrices whose ranks are no more than s . Note that (4) is a special case of (3) with matrix \mathbf{A} being the identity matrix. This motivates us to solve (3) through the simpler problem (4).

When \mathbf{A} is nonsingular, clearly (3) has a global minimizer $\mathbf{A}^{-1} \mathcal{P}_s(\mathbf{Z})$ by rank preserversness of any nonsingular matrix in matrix multiplication. When \mathbf{A} is singular, we first assume that $r(\mathbf{A}) \geq s$. However, this assumption is by no means mandatory, which will be discussed later. Given the SVD of $\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^T$, with orthogonal matrices $\mathbf{U} \in \mathbb{R}^{n \times n}$ and $\mathbf{V} \in \mathbb{R}^{p \times p}$ and diagonal matrix $\mathbf{D} \in \mathbb{R}^{n \times p}$, we have

$$\|\mathbf{A}\Theta - \mathbf{Z}\|_F = \|\mathbf{U}^T(\mathbf{A}\Theta - \mathbf{Z})\|_F = \|\mathbf{D}\mathbf{V}^T\Theta - \mathbf{U}^T\mathbf{Z}\|_F.$$

This follows from the fact that the Frobenius-norm is invariant under any orthogonal transformation. Let $\mathbf{Y} = \mathbf{V}^T\Theta$ and $\mathbf{W} = \mathbf{U}^T\mathbf{Z}$, clearly we have $r(\mathbf{Y}) = r(\Theta)$. Now solving (3) amounts to solving an equivalent form:

$$\begin{aligned} \min_{\mathbf{Y}} \quad & \|\mathbf{D}\mathbf{Y} - \mathbf{W}\|_F^2 \\ \text{s.t.} \quad & r(\mathbf{Y}) \leq s. \end{aligned} \quad (5)$$

Consequently, a global minimizer of (3) becomes $\mathbf{V}\mathbf{Y}^*$, where \mathbf{Y}^* is a global minimizer of (5) and is given by the following theorem.

THEOREM 1. *Let \mathbf{D} , \mathbf{Y} , \mathbf{Z} and s be as defined above. If $s \leq r(\mathbf{A})$, then a global minimizer of (5) is given by*

$$\mathbf{Y}^* = \begin{bmatrix} \mathbf{D}_{r(\mathbf{A})}^{-1} \mathcal{P}_s(\mathbf{W}_{r(\mathbf{A})}) \\ \mathbf{a} \end{bmatrix}, \quad (6)$$

where $\mathbf{D}_{r(\mathbf{A})}$ is a diagonal matrix consisting of all the nonzero singular value of \mathbf{A} , \mathbf{a} can be set as the zero matrix, and $\mathbf{W}_{r(\mathbf{A})}$ consists of the first $r(\mathbf{A})$ rows of \mathbf{W} .

Here are some remarks regarding the above theorem:

Remark 1. The solution to problem (5) is generally not unique. Specifically, the matrix \mathbf{a} in (6) need not be fixed at zero, as long as it does not change the rank of \mathbf{Y}^* . However, if \mathbf{A} is of full column rank, i.e., when $r(\mathbf{A}) = p$, then \mathbf{a} vanishes and \mathbf{Y}^* can be uniquely determined. In this case, the optimal solution of (3) is also unique.

Remark 2. The optimal Y^* can also be computed for the general matrix A with an arbitrary rank, i.e., when $r(\mathbf{A}) < s$. See the proof of Theorem 1 in the Section 6.

It is important to note that the value of \mathbf{a} is irrelevant for prediction, but matters for parameter estimation. In other words, when $r(\mathbf{A}) < p$, a global minimizer is not unique, hence that parameter estimation is not identifiable: see Section 4 for a discussion. For simplicity, we set $\mathbf{a} = \mathbf{0}$ for \mathbf{Y}^* subsequently.

In what follows, our estimator is defined as $\hat{\Theta}^s$, as well as an estimated rank \hat{r} . Algorithm 1 below summarizes the main steps for computing $\hat{\Theta}^s$ with regard to $s \leq \min(n, k, p)$, where the LSRM stands for Least Squares Rank Minimization.

Algorithm 1 Exact solution of (3)

Input: $\mathbf{A}, \mathbf{Z}, s \leq r(\mathbf{A})$

Output: A global minimizer Θ of (3)

Function: LSRM($\mathbf{A}, \mathbf{Z}, s$)

```

1: if  $\mathbf{A}$  is nonsingular then
2:    $\Theta = \mathbf{A}^{-1}\mathcal{P}_s(\mathbf{Z})$ 
3: else
4:   Perform SVD on  $\mathbf{A}$ :  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ 
5:   Extract the first  $r$  rows of  $\mathbf{U}^T\mathbf{Z}$  and denote it as
      $\mathbf{W}_{r(\mathbf{A})}$ 
6:    $\Theta = \mathbf{V} \begin{bmatrix} \mathbf{D}_{r(\mathbf{A})}^{-1}\mathcal{P}_s(\mathbf{W}_{r(\mathbf{A})}) \\ \mathbf{0} \end{bmatrix}$ 
7: end if
8: return  $\Theta$ 

```

The complexity of Algorithm 1 is determined mainly by the most expensive operations—matrix inversion and SVD, specifically, at most one matrix inversion and two SVDs. Such operations roughly require a cubic time complexity¹.

3. REGULARIZATION AND SOLUTION PATH

This section studies a regularized counterpart of (3):

$$\min_{\Theta} \|\mathbf{A}\Theta - \mathbf{Z}\|_F^2 + \lambda r(\Theta), \quad (7)$$

where $\lambda > 0$ is a continuous regularization parameter corresponding to s in (3), and Θ_λ^* is a global minimizer of (7). The next theorem establishes an equivalence between (7) and (3) when Θ_λ^* is unique, occurring when $r(\mathbf{A}) = p$. Such a result is not generally anticipated for a nonconvex problem.

THEOREM 2 (EQUIVALENCE). *When $p = r(\mathbf{A})$, (7) has a unique global minimizer. Moreover, (7) and (3) are equivalent with respect to their solutions. For any Θ_λ^* with $\lambda \geq 0$, there exists $1 \leq s^* = r(\Theta_\lambda^*)$ such that $\Theta_\lambda^* = \hat{\Theta}^{s^*}$, and vice versa.*

Next we develop an algorithm for computing an entire solution path for all values of λ with complexity comparable

¹More specifically, for a matrix of dimension $n \times p$, the SVD has a complexity of $O(\min\{n^2p, p^2n\})$, whereas the matrix inversion has a complexity of $O(r(\mathbf{A})^3)$, which can be improved to $O(r(\mathbf{A})^{2.807})$ when the Strassen Algorithm is utilized.

to that of solving (7) at a single λ -value. For motivation, first consider a special case of the identity \mathbf{A} in (7):

$$g(\lambda) = \min_{\Theta} \|\Theta - \mathbf{Z}\|_F^2 + \lambda r(\Theta). \quad (8)$$

3.1 Monotone property

In (8), $r(\Theta)$ decreases as λ increases from 0, where $r(\Theta)$ goes through all integer values from $r(\mathbf{Z})$ down to 0 when λ becomes sufficiently large. In addition, the value of $g(\lambda)$ is nondecreasing as λ increases. The next theorem summarizes these results.

THEOREM 3 (MONOTONE PROPERTY). *Let $r(\mathbf{Z})$ be r . Then the following properties hold:*

(1) *There exists a solution path vector \mathcal{S} of length $r + 2$ satisfying the following:*

$$\mathcal{S}_0 = 0, \quad \mathcal{S}_{r+1} = +\infty, \quad \mathcal{S}_{k+1} > \mathcal{S}_k, \quad k = 0, 1, \dots, r$$

$$\Theta_\lambda^* = \mathcal{P}_{r-k}(\mathbf{Z}), \quad \text{if } \mathcal{S}_k \leq \lambda < \mathcal{S}_{k+1},$$

(2) *Function $g(\lambda)$ is nondecreasing and piecewise linear.*

The monotone property leads to an efficient algorithm for calculating the pathwise solution of (8). Figure 1 displays the solution path by illustrating the function value $g(\lambda)$ and $r(\Theta_\lambda^*)$ as a function of λ .

3.2 Pathwise algorithm

Through the monotone property, we compute the optimal solution of (8) at a particular λ by locating λ in the correct interval in the solution path vector \mathcal{S} , which can be achieved efficiently via a simple binary search. Algorithm 2 describes the main steps.

Algorithm 2 Pathwise solution of (8)

Input: Θ, \mathbf{Z}

Output: Solution path vector \mathcal{S} , pathwise solution Θ

Function: pathwise(Θ, \mathbf{Z})

```

1: Initialize:  $\mathcal{S}_0 = 0, \Theta_0 = \mathbf{Z}, r = r(\mathbf{Z})$ 
2: Perform SVD on  $\mathbf{Z}$ :  $\mathbf{Z} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ 
3: for  $i = r$  down to 1 do
4:    $\mathcal{S}_{r-i+1} = \sigma_i^2$ 
5:    $\Theta_{r-i+1} = \Theta_{r-i} - \sigma_i u_i v_i^T$ 
6: end for
7: return  $\mathcal{S}, \Theta$ 

```

Algorithm 2 requires only one SVD operation, therefore its complexity is of the same order as that of Algorithm 1 at a single s -value. When \mathbf{Z} is a low-rank matrix, existing software for SVD computation such as PROPACK is applicable to further improve computational efficiency.

3.3 Extension to general \mathbf{A}

For a general design matrix \mathbf{A} , note that

$$\|\mathbf{A}\Theta - \mathbf{Z}\|_F^2 = \|\mathbf{D}\mathbf{Y} - \mathbf{W}\|_F^2 = \|\mathbf{W}'\|_F^2 + \|\mathbf{D}_r\mathbf{Y}_r - \mathbf{W}_r\|_F^2,$$

where $\mathbf{W} = \begin{bmatrix} \mathbf{W}_r \\ \mathbf{W}' \end{bmatrix}$. After ignoring the constant term $\|\mathbf{W}'\|_F^2$, we solve

$$\min_{\mathbf{Y}_r} \|\mathbf{D}_r\mathbf{Y}_r - \mathbf{W}_r\|_F^2 + \lambda r(\mathbf{Y}_r).$$

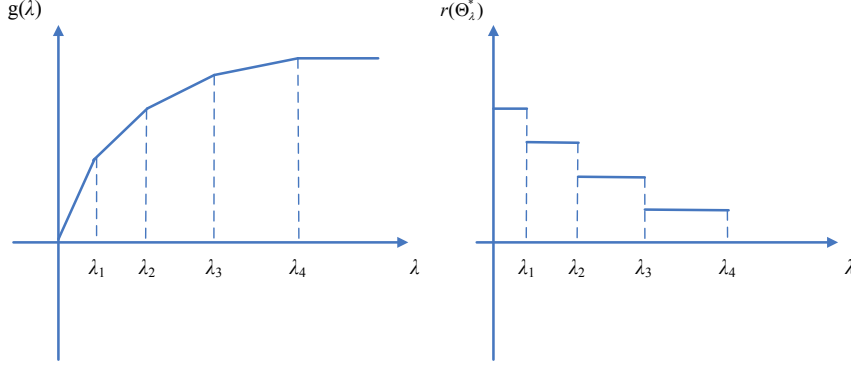


Figure 1: Piecewise linearity of $g(\cdot)$ and the rank of the optimal solution with respect to λ .

Note that \mathbf{D}_r is nonsingular, then the problem reduces to the simple case

$$\min_{\hat{\mathbf{Y}}} \|\hat{\mathbf{Y}} - \mathbf{W}_r\|_F^2 + \lambda r(\hat{\mathbf{Y}}),$$

where $\hat{\mathbf{Y}} = \mathbf{D}_r \mathbf{Y}_r$. The solution path can be obtained directly from Algorithm 2.

4. STATISTICAL PROPERTIES

This section is devoted to the theoretical investigation of least squares rank minimization, which remains largely unexplored, although nuclear-norm regularization has been studied. In particular, we will reveal what is the best performance for prediction as well as the optimal risk for parameter estimation. Moreover, we will establish the optimality of the proposed method. In fact, the proposed method reconstructs the oracle estimator, the optimal one as if the true rank were known in advance. Here the oracle estimator $\hat{\Theta}^0$ is defined as a global minimizer of $\|\mathbf{A}\Theta - \mathbf{Z}\|_F^2$ subject to $r(\Theta) = r_0$, where Θ^0 and $r_0 = r(\Theta^0) \geq 1$ denote the true parameter matrix and the true rank respectively. This leads to the exact rank recovery, in addition to the reconstruction of the optimal performance of the oracle estimator. In other words, the proposed method is optimal against any other methods such as the nuclear-norm rank regularization.

Given the design matrix \mathbf{A} , we study the accuracy of rank recovery as well as prediction and parameter estimation. Let \mathbb{P} and \mathbb{E} be the true probability and expectation under Θ^0 given \mathbf{A} . For rank recovery, we use the metric $\mathbb{P}(\hat{r} \neq r_0)$. For prediction and parameter estimation, we employ the risk $\mathbb{E}K(\hat{\Theta}^s, \Theta^0)$ and $\mathbb{E}\|\hat{\Theta}^s - \Theta^0\|_F^2$ respectively, where

$$\begin{aligned} K(\hat{\Theta}^s, \Theta^0) &= (2\sigma^2 n)^{-1} \sum_{i=1}^n \|\mathbf{a}_i^T (\hat{\Theta}^s - \Theta^0)\|_2^2 \\ &= (2\sigma^2 n)^{-1} \|\mathbf{A}(\hat{\Theta}^s - \Theta^0)\|_F^2 \end{aligned}$$

is the Kullback-Leibler loss and $\|\cdot\|_2$ is the L_2 -norm for a vector. Note that the predictive risk equals to $2\sigma^2 \mathbb{E}K(\hat{\Theta}^s, \Theta^0)$ and parameter estimation is considered when it is identifiable, i.e., $r(\mathbf{A}) = p$.

Now we present the risk bounds under (1) without a Gaussian error assumption.

THEOREM 4. *Under (1), the oracle estimator is exactly reconstructed by our method in that $\hat{\Theta}^{r_0} = \Theta^0$ under \mathbb{P} ,*

when $r_0 \leq \min(r(\mathbf{A}), p)$. As a result, exact reconstruction of the optimal performance is achieved by our estimator $\hat{\Theta}^{r_0}$. In particular,

$$\mathbb{E}K(\hat{\Theta}^{r_0}, \Theta^0) \begin{cases} = & \frac{r_0 k}{2n} & \text{if } r_0 = r(\mathbf{A}) \\ \leq & \frac{2\mathbb{E}(\sum_{j=1}^{r_0} \sigma_j^2)}{n} & \text{if } r_0 < r(\mathbf{A}) \end{cases}$$

and

$$\mathbb{E}\|\hat{\Theta}^{r_0} - \Theta^0\|_F^2 \begin{cases} = & \frac{r_0 k}{\sigma_{\min}^2(\mathbf{A}) n} & \text{if } r_0 = r(\mathbf{A}) = p \\ \leq & \frac{\mathbb{E}(\sum_{j=1}^{r_0} \sigma_j^2)}{\sigma_{\min}^2 n} & \text{if } r_0 < r(\mathbf{A}) = p \end{cases},$$

where σ_j and $\sigma_{\min} > 0$ are the j th largest and the smallest nonzero singular value of $\mathbf{e}^* = \mathbf{U}_{r(\mathbf{A})}^T \mathbf{e}$ and $n^{-1/2} \mathbf{A}$, respectively, and $\mathbf{U}_{r(\mathbf{A})}$ denotes the first $r(\mathbf{A})$ rows of \mathbf{U} .

Remark: In general, $\mathbb{E} \sum_{j=1}^{r_0} \sigma_j^2 \leq r_0 \mathbb{E} \sigma_1^2$.

Theorem 4 says that the optimal oracle estimator is exactly reconstructed by our method. Interestingly, the true rank is exactly recovered from noisy data, which is attributed to discreteness of the rank and is analogous to the maximum likelihood estimation over a discrete parameter space. Concerning prediction and parameter estimation, the optimal Kullback-Leibler risk is $\frac{r_0 k}{n}$ but the risk under the Frobenius-norm is $\frac{r_0 k}{\sigma_{\min}^2 n}$. For prediction, only the effective degrees of freedom $\sum_{j=1}^{r_0} \sigma_j^2$ matters as opposed to p , which is in contrast to a rate $\frac{r_0 k p}{n}$ without a rank restriction. This permits p to be much larger than n , or $k, p \gg n$. For estimation, however, p enters the risk through σ_{\min}^2 , where p can not be larger than n , or $\max(k, p) \leq n$.

5. TUNING

As shown in Section 4, theoretically, exact rank reconstruction can be accomplished through tuning. Practically, we employ a predictive measure for rank selection.

The predictive performance of $\hat{\Theta}^s$ is measured by

$$\text{MSE}(\hat{\Theta}^s) = \frac{1}{n} \mathbb{E} \|\mathbf{A} \hat{\Theta}^s - \mathbf{Z}\|_F^2,$$

which is proportional to the risk $R(\hat{\Theta}^s)$, where the expectation is taken with respect to (\mathbf{Z}, \mathbf{A}) .

To estimate s over integer values in $[0, \dots, \min(n, p, k)]$, one may cross-validate through a tuning data set, which estimates the MSE. Alternatively, one may use the generalized degrees of freedom [31] through data perturbation without a tuning set.

$$\widehat{\text{GDF}}(\hat{\Theta}^s) = n^{-1} \|\mathbf{Z} - \mathbf{A}\hat{\Theta}^s\|^2 + 2n^{-1} \sum_{i=1}^n \sum_{l=1}^k \widehat{\text{Cov}}(\mathbf{z}_{il}, (\mathbf{a}_i^T \hat{\Theta}^s)_l), \quad (9)$$

where $\widehat{\text{Cov}}(\mathbf{z}_{il}, (\mathbf{a}_i^T \hat{\Theta}^s)_l)$ is the estimated covariance between the l th component of \mathbf{z}_i and the l th component of $\mathbf{a}_i^T \hat{\Theta}^s$. In the case that \mathbf{e}_i in (1) follows $N(\mathbf{0}, \sigma^2 \mathbf{I}_{k \times k})$, the method of data perturbation of [31] is applicable. Specifically, sample \mathbf{e}_i^* from $N(\mathbf{0}, \sigma^2 \mathbf{I}_{k \times k})$ and let

$$\begin{aligned} \mathbf{Z}^* &= (1 - \tau)\mathbf{Z} + \tau \tilde{\mathbf{Z}} \\ \widehat{\text{Cov}}(\mathbf{z}_{il}, (\mathbf{a}_i^T \hat{\Theta}^s)_l) &= \tau^{-1} \text{Cov}^*(\mathbf{z}_{il}^*, (\mathbf{a}_i^T \hat{\Theta}^{*s})_l) \end{aligned}$$

where $\text{Cov}^*(\mathbf{z}_{il}^*, (\mathbf{a}_i^T \hat{\Theta}^{*s})_l)$ is the sample covariance with the Monte Carlo size T . For the types of problems we consider, we fixed T to be n .

6. PROOF OF THEOREMS

In this section we present detailed proofs to the theorems developed in the previous sections. We first present a technical lemma to be used in the proof of Theorem 4.

LEMMA 2. Suppose \mathbf{A} and \mathbf{B} are two $n_1 \times n_2$ matrices. Then,

$$\langle \mathbf{A}, \mathbf{B} \rangle \leq \|\mathbf{A}\|_F \|\mathcal{P}_{r(\mathbf{A})}(\mathbf{B})\|_F, \quad (10)$$

where $\langle \mathbf{A}, \mathbf{B} \rangle = \text{Tr}(\mathbf{A}^T \mathbf{B}) = \text{Tr}(\mathbf{B}^T \mathbf{A})$, Tr denotes the trace, and $r(\mathbf{A})$ is the rank of \mathbf{A} .

PROOF. Let the singular value decomposition of \mathbf{A} and \mathbf{B} be $\mathbf{A} = \mathbf{U}_1 \Sigma_1 \mathbf{V}_1^T$ and $\mathbf{B} = \mathbf{U}_2 \Sigma_2 \mathbf{V}_2^T$ where \mathbf{U}_i and $\mathbf{V}_i, i = 1, 2$, are orthogonal matrices, and Σ_1 and Σ_2 are diagonal matrices with their diagonal elements being the singular values of \mathbf{A} and \mathbf{B} , respectively. Then

$$\begin{aligned} \langle \mathbf{A}, \mathbf{B} \rangle &= \text{Tr}(\mathbf{V}_1 \Sigma_1^T \mathbf{U}_1^T \mathbf{U}_2 \Sigma_2 \mathbf{V}_2^T) \\ &= \text{Tr}(\Sigma_1^T \mathbf{U}_1^T \mathbf{U}_2 \Sigma_2 \mathbf{V}_2^T \mathbf{V}_1) \\ &\equiv \text{Tr}(\Sigma_1^T \mathbf{U} \Sigma_2 \mathbf{V}^T), \end{aligned}$$

where $\mathbf{U} = \mathbf{U}_1^T \mathbf{U}_2$ and $\mathbf{V} = \mathbf{V}_1^T \mathbf{V}_2$ continue to be orthogonal. Let the ordered singular values of \mathbf{A} be $\sigma_1 \geq \dots \geq \sigma_{r(\mathbf{A})}$ and $\tilde{\mathbf{B}} = (\tilde{b}_{ij}) = \mathbf{U} \Sigma_2 \mathbf{V}^T$. By Cauchy-Schwarz's inequality,

$$\begin{aligned} \text{Tr}(\Sigma_1^T \mathbf{U} \Sigma_2 \mathbf{V}^T) &= \text{Tr}(\Sigma_1^T \tilde{\mathbf{B}}) \\ &= \sum_{i=1}^{r(\mathbf{A})} \sigma_i \tilde{b}_{ii} \\ &\leq \sqrt{\sum_{i=1}^{r(\mathbf{A})} \sigma_i^2} \sqrt{\sum_{i=1}^{r(\mathbf{A})} \tilde{b}_{ii}^2} \\ &= \|\mathbf{A}\|_F \sqrt{\sum_{i=1}^{r(\mathbf{A})} \tilde{b}_{ii}^2}. \end{aligned} \quad (11)$$

Similarly, let the ordered singular values of \mathbf{B} be $\eta_1 \geq \dots \geq \eta_{r(\mathbf{B})}$. Then it suffices to show that $\sum_{i=1}^{r(\mathbf{A})} \tilde{b}_{ii}^2 \leq \sum_{i=1}^{r(\mathbf{A})} \eta_i^2$.

Assume, without loss of generality, that $\eta_i = 0$ if $i > r(\mathbf{B})$. Let $n = \min(n_1, n_2)$. By Cauchy-Schwarz's inequality,

$$\begin{aligned} \sum_{i=1}^{r(\mathbf{A})} \tilde{b}_{ii}^2 &= \sum_{i=1}^{r(\mathbf{A})} \left(\sum_{k=1}^n u_{ik} \eta_k v_{ik} \right)^2 \\ &\leq \sum_{i=1}^{r(\mathbf{A})} \left(\sum_{k=1}^n u_{ik}^2 \eta_k^2 \right) \left(\sum_{k=1}^n v_{ik}^2 \right) \\ &\leq \sum_{i=1}^{r(\mathbf{A})} \left(\sum_{k=1}^n u_{ik}^2 \eta_k^2 \right) \\ &= \sum_{k=1}^n \eta_k^2 \left(\sum_{i=1}^{r(\mathbf{A})} u_{ik}^2 \right) \\ &\leq \sum_{k=1}^{r(\mathbf{A})} \eta_k^2, \end{aligned}$$

where the last step uses the fact that $\sum_{i=1}^{r(\mathbf{A})} u_{ik}^2 \leq 1$ and $\sum_{k=1}^n \sum_{i=1}^{r(\mathbf{A})} u_{ik}^2 = \sum_{i=1}^{r(\mathbf{A})} \sum_{k=1}^n u_{ik}^2 \leq r(\mathbf{A})$. A combination of the above bounds lead to the desired results. This completes the proof. \square

6.1 Proof of Theorem 1

First partition \mathbf{D} and \mathbf{W} as follows:

$$\mathbf{D} = \begin{bmatrix} \mathbf{D}_{r(\mathbf{A})} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{W} = \begin{bmatrix} \mathbf{W}_{r(\mathbf{A})} \\ \mathbf{W}' \end{bmatrix},$$

then

$$\begin{aligned} \mathbf{D}\mathbf{Y} - \mathbf{W} &= \begin{bmatrix} \mathbf{D}_{r(\mathbf{A})} \mathbf{Y}_{r(\mathbf{A})} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{W}_{r(\mathbf{A})} \\ \mathbf{W}' \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{D}_{r(\mathbf{A})} \mathbf{Y}_{r(\mathbf{A})} - \mathbf{W}_{r(\mathbf{A})} \\ -\mathbf{W}' \end{bmatrix}. \end{aligned}$$

Evidently, only the first $r(\mathbf{A})$ rows of \mathbf{Y} are involved in minimizing $\|\mathbf{D}\mathbf{Y} - \mathbf{W}\|_2^2$, which amounts to computing the global minimizer $\mathbf{Y}_{r(\mathbf{A})}^*$ of

$$\arg \min_{\mathbf{Y}_{r(\mathbf{A})}} \|\mathbf{D}_{r(\mathbf{A})} \mathbf{Y}_{r(\mathbf{A})} - \mathbf{W}_{r(\mathbf{A})}\|_F^2.$$

Then $\mathbf{Y}_{r(\mathbf{A})}^* = \mathbf{D}_{r(\mathbf{A})}^{-1} \mathcal{P}_s(\mathbf{W}_{r(\mathbf{A})})$ by non-singularity of $\mathbf{D}_{r(\mathbf{A})}$ and Lemma 1 with $s \leq r(\mathbf{A})$. For $s > r(\mathbf{A})$, recall that, only the upper part of \mathbf{Y}^* is relevant in minimizing (5). The result then follows. This completes the proof.

6.2 Proof of Theorem 2

For any Θ_λ^* with $\lambda > 0$, let $s^* = r(\Theta_\lambda^*)$. Next we prove by contradiction that $\Theta_\lambda^* = \hat{\Theta}^{s^*}$. Suppose $\Theta_\lambda^* \neq \hat{\Theta}^{s^*}$. By uniqueness of $\hat{\Theta}^{s^*}$ given in Theorem 1 and the definition of minimization, $\|\mathbf{A}\hat{\Theta}^{s^*} - \mathbf{Z}\|_F^2 < \|\mathbf{A}\Theta_\lambda^* - \mathbf{Z}\|_F^2$. This, together with $r(\hat{\Theta}^{s^*}) = r(\Theta_\lambda^*)$, implies that

$$\|\mathbf{A}\hat{\Theta}^{s^*} - \mathbf{Z}\|_F^2 + \lambda r(\hat{\Theta}^{s^*}) < \|\mathbf{A}\Theta_\lambda^* - \mathbf{Z}\|_F^2 + \lambda r(\Theta_\lambda^*).$$

This contradicts to the fact that Θ_λ^* is a minimizer of (7). The converse is revealed in the proof of Theorem 3.

6.3 Proof of Theorem 3

We prove the first conclusion by constructing such a solution path vector \mathcal{S} . Let $\mathcal{S}_0 = 0$, $\mathcal{S}_{r+1} = +\infty$. Define \mathcal{S}_k for $1 \leq k \leq r$ as the solution of equation:

$$\|\mathcal{P}_{r-k+1}(\mathbf{Z}) - \mathbf{Z}\|_F^2 + \mathcal{S}_k(r-k+1) = \|\mathcal{P}_{r-k}(\mathbf{Z}) - \mathbf{Z}\|_F^2 + \mathcal{S}_k(r-k).$$

It follows that

$$\begin{aligned} \mathcal{S}_k &= \|\mathcal{P}_{r-k}(\mathbf{Z}) - \mathbf{Z}\|_F^2 - \|\mathcal{P}_{r-k+1}(\mathbf{Z}) - \mathbf{Z}\|_F^2 \\ &= \sum_{j=r-k+1}^r \sigma_j^2 - \sum_{j=r-k+2}^r \sigma_j^2 \\ &= \sigma_{r-k+1}^2, \end{aligned} \quad (12)$$

where σ_j is the j th largest non-zero singular value of \mathbf{Z} . By (12), \mathcal{S}_k is increasing. In addition, by definition of \mathcal{S}_k and \mathcal{S}_{k+1} , whenever λ falls into the interval $[\mathcal{S}_k, \mathcal{S}_{k+1})$, the rank of a global minimizer $\hat{\Theta}^*$ of (8) would be no more than $r-k$ and larger than $r-k-1$. In other words, $\hat{\Theta}_\lambda^*$ is always of rank $r-k$ and is given by $\mathcal{P}_{r-k}(\mathbf{Z})$. Therefore, the constructed solution path vector \mathcal{S} satisfies all the requirements in the theorem.

Moreover, when $\mathcal{S}_k \leq \lambda < \mathcal{S}_{k+1}$,

$$\begin{aligned} g(\lambda) &= \|\mathcal{P}_{r-k}(\mathbf{Z}) - \mathbf{Z}\|_F^2 + \lambda r(\mathcal{P}_{r-k}(\mathbf{Z})) \\ &= \|\mathcal{P}_{r-k}(\mathbf{Z}) - \mathbf{Z}\|_F^2 + (r-k)\lambda. \end{aligned} \quad (13)$$

Since $\mathcal{P}_{r-k}(\mathbf{Z})$ is independent of λ , $g(\lambda)$ is a nondecreasing linear function of λ in each interval $[\mathcal{S}_k, \mathcal{S}_{k+1})$. Combined with the definition of the solution path vector \mathcal{S} , we conclude that $g(\lambda)$ is nondecreasing and piecewise linear with each element of \mathcal{S} as a kink point, as shown in Figure 1. The proof is completed.

6.4 Proof of Theorem 4

The proof uses direct calculations. First we bound the Kullback-Leibler loss. By Theorem 1,

$$\mathbf{A}\hat{\Theta}^{r_0} = \mathbf{U}_{r(\mathbf{A})}\mathcal{P}_{r_0}(\mathbf{W}_{r(\mathbf{A})}),$$

where

$$\mathbf{W}_{r(\mathbf{A})} = \mathbf{D}_{r(\mathbf{A})}(\mathbf{V}^T\Theta^0)_{r(\mathbf{A})} + (\mathbf{U}^T\mathbf{e})_{r(\mathbf{A})}.$$

Denote $\mathbf{B} = \mathbf{D}_{r(\mathbf{A})}(\mathbf{V}^T\Theta^0)_{r(\mathbf{A})}$ with rank $r(\mathbf{B}) \leq r_0$. Since the Frobenius-norm is invariant under orthogonal transformation, it follows that

$$\begin{aligned} &\|\mathbf{A}\hat{\Theta}^{r_0} - \mathbf{A}\Theta^0\|_F^2 \\ &= \|\mathbf{U}_{r(\mathbf{A})}\mathcal{P}_{r_0}(\mathbf{W}_{r(\mathbf{A})}) - \mathbf{U}\mathbf{D}\mathbf{V}^T\Theta^0\|_F^2 \\ &= \|\mathbf{U}_{r(\mathbf{A})}\mathcal{P}_{r_0}(\mathbf{W}_{r(\mathbf{A})}) - \mathbf{U}_{r(\mathbf{A})}\mathbf{D}_{r(\mathbf{A})}(\mathbf{V}^T\Theta^0)_{r(\mathbf{A})}\|_F^2 \\ &= \|\mathcal{P}_{r_0}(\mathbf{W}_{r(\mathbf{A})}) - \mathbf{D}_{r(\mathbf{A})}(\mathbf{V}^T\Theta^0)_{r(\mathbf{A})}\|_F^2 \\ &= \|\mathcal{P}_{r_0}(\mathbf{B} + (\mathbf{U}^T\mathbf{e})_{r(\mathbf{A})}) - \mathbf{B}\|_F^2, \\ &= \|\mathcal{P}_{r_0}(\mathbf{B} + \mathbf{e}^*) - \mathbf{B}\|_F^2, \end{aligned}$$

where $\mathbf{e}^* = (\mathbf{U}^T\mathbf{e})_{r(\mathbf{A})}$. From the definition of $\mathcal{P}_{r_0}(\mathbf{B} + \mathbf{e}^*)$ we can conclude that:

$$\|\mathcal{P}_{r_0}(\mathbf{B} + \mathbf{e}^*) - \mathbf{B} - \mathbf{e}^*\|_F^2 \leq \|\mathbf{B} - \mathbf{B} - \mathbf{e}^*\|_F^2 = \|\mathbf{e}^*\|_F^2,$$

which implies that,

$$\begin{aligned} \|\mathcal{P}_{r_0}(\mathbf{B} + \mathbf{e}^*) - \mathbf{B}\|_F^2 &\leq 2\langle \mathcal{P}_{r_0}(\mathbf{B} + \mathbf{e}^*) - \mathbf{B}, \mathbf{e}^* \rangle \\ &\leq 2\|\mathcal{P}_{r_0}(\mathbf{B} + \mathbf{e}^*) - \mathbf{B}\|_F \|\mathcal{P}_{r_0}(\mathbf{e}^*)\|_F, \end{aligned}$$

where the last inequality follows from Lemma 2. Thus,

$$\|\mathcal{P}_{r_0}(\mathbf{B} + \mathbf{e}^*) - \mathbf{B}\|_F^2 \leq 4\|\mathcal{P}_{r_0}(\mathbf{e}^*)\|_F^2 = 4\sum_{j=1}^{r_0} \sigma_j^2. \quad (14)$$

The risk bounds then follow.

Second we bound $\|\hat{\Theta}^{r_0} - \Theta^0\|_F^2$, which is equal to

$$\begin{aligned} &\|\mathbf{V} \begin{bmatrix} \mathbf{D}_{r(\mathbf{A})}^{-1}\mathcal{P}_{r_0}(\mathbf{W}_{r(\mathbf{A})}) \\ \mathbf{0} \end{bmatrix} - \Theta^0\|_F^2 \\ &= \left\| \begin{bmatrix} \mathbf{D}_{r(\mathbf{A})}^{-1}\mathcal{P}_{r_0}(\mathbf{W}_{r(\mathbf{A})}) \\ \mathbf{0} \end{bmatrix} - \mathbf{V}^T\Theta^0 \right\|_F^2 \\ &= \|\mathbf{D}_{r(\mathbf{A})}^{-1}\mathcal{P}_{r_0}(\mathbf{W}_{r(\mathbf{A})}) - (\mathbf{V}^T\Theta^0)_{r(\mathbf{A})}\|_F^2 + \|(\mathbf{V}^T\Theta^0)_{r(\mathbf{A})^c}\|_F^2 \\ &\leq \frac{1}{\sigma_{\min}^2 n} \|\mathcal{P}_{r_0}(\mathbf{W}_{r(\mathbf{A})}) - \mathbf{D}_{r(\mathbf{A})}(\mathbf{V}^T\Theta^0)_{r(\mathbf{A})}\|_F^2 \\ &\quad + \|(\mathbf{V}^T\Theta^0)_{r(\mathbf{A})^c}\|_F^2, \end{aligned}$$

where $\sigma_{r(\mathbf{A})}(n^{-1/2}\mathbf{A}) = n^{1/2}\sigma_{\min}$. If $p = r(\mathbf{A})$, then the last term vanishes. Thus

$$\begin{aligned} \|\hat{\Theta}^{r_0} - \Theta^0\|_F^2 &\leq \frac{1}{\sigma_{r(\mathbf{A})}^2 n} \|\mathcal{P}_{r_0}(\mathbf{W}_{r(\mathbf{A})}) - \mathbf{D}_{r(\mathbf{A})}(\mathbf{V}^T\Theta^0)_{r(\mathbf{A})}\|_F^2 \\ &\leq \frac{4\sum_{j=1}^{r_0} \sigma_j^2}{\sigma_{\min}^2 n}. \end{aligned}$$

Finally, if $r(\mathbf{A}) \geq r_0$, $\mathbb{E}\|\hat{\Theta}^{r_0} - \Theta^0\|_F^2 \leq \frac{4}{\sigma_{r(\mathbf{A})}^2} \mathbb{E}(\sum_{j=1}^{r_0} \sigma_j^2)$ and $\mathbb{E}\|\mathbf{A}\hat{\Theta}^{r_0} - \mathbf{A}\Theta^0\|_F^2 \leq 4\mathbb{E}(\sum_{j=1}^{r_0} \sigma_j^2)$. In particular, if $r(\mathbf{A}) = r_0$, then

$$\begin{aligned} \mathbb{E}\|\hat{\Theta}^{r_0} - \Theta^0\|_F^2 &= \frac{1}{\sigma_{\min}^2} \frac{r_0 k}{n} \\ \mathbb{E}\|\mathbf{A}\hat{\Theta}^{r_0} - \mathbf{A}\Theta^0\|_F^2 &= \frac{r_0 k}{n}. \end{aligned}$$

7. EMPIRICAL EVALUATIONS

This section examines effectiveness of the proposed method and compares it with the nuclear-norm regularization as well as the SVP method [17]. One benefit of our method is that it can evaluate the approximation quality of the inexact ones. For the SVP, we choose a default initial value $\mathbf{0}$ for this local method since no other choices are guaranteed to deliver better performance. Our numerical experiment, which is not reported in here, suggests that the SVP method is indeed sensitive to the choice of the initial value. For nuclear-norm regularization, we select a suitable regularization parameter value giving the solution satisfying the rank constraint in (3).

7.1 Synthetic Data

Simulations are performed under (2). First, the $n \times p$ design matrix \mathbf{A} is sampled, with each entry being iid $N(10, 1)$. Second, the $p \times k$ truth Θ^0 is generated by multiplying a $p \times r$ matrix and an $r \times k$ matrix, each entry of which has a normal distribution over $N(10, 1)$. The data matrix \mathbf{Z} is then sampled according to (2) with \mathbf{e} following iid $N(0, \sigma^2)$ with $\sigma = 0.5$.

For predictive performance and rank recovery, we compute the MSE $\|\mathbf{A}(\Theta^0 - \hat{\Theta}^{\hat{r}})\|_F^2$, the absolute difference $|\hat{r} - r_0|$ and record the training time for each method, averaged over 100 simulation replications on a test set of size $10n$, where \hat{r} is tuned over integers in $[0, \min(n, p)]$ by an independent tuning set of size $2n$. We consider three possible situations, i.e., when $k = r_0 < p$, $k > p > r_0$ and $p > k > r_0$. To illustrate our theoretical conclusion in Theorem 4 that the prediction error bound does not depend on the value of p , we add one more case of $p > k > r_0$ with a larger value of p . The simulation results are summarized in Tables 2-4.

Table 2: Prediction results for the synthetic data: averaged MSEs as well as their standard deviations, for three competing methods based on the selected tuning parameters over 100 simulation replications. Our, SVP, and Nuclear-Norm refer to our method, the SVP method and nuclear-norm regularization method.

| Set-up | Our | SVP | Nuclear-Norm |
|----------------------------|--------------------|------------------------|------------------------|
| $k = r_0 < p, 5 = 5 < 99$ | 0.0110 (0.0065) | 0.4591 (0.0386) | 0.4130 (0.0360) |
| $k > p > r_0, 50 > 40 > 5$ | 12.0010 (0.6125) | 10166.5052 (1236.3237) | 7709.0779 (960.6982) |
| $p > k > r_0, 50 > 40 > 5$ | 13.1945 (0.5682) | 10185.7699 (1132.3515) | 8116.2653 (943.6058) |
| $p > k > r_0, 99 > 40 > 5$ | 121.4019 (70.3054) | 14597.0728 (1185.7200) | 13153.3217 (1056.5155) |

Table 3: Rank recovery for the synthetic data: averaged values of $|\hat{r} - r_0|$ as well as their standard deviations, for three competing methods based on the selected tuning parameters over 100 simulation replications. Our, SVP, and Nuclear-Norm refer to our method, the SVP method and nuclear-norm regularization method.

| Set-up | Our | SVP | Nuclear-Norm |
|----------------------------|-----------------|-----------------|-----------------|
| $k = r_0 < p, 5 = 5 < 99$ | 0.8400 (0.8005) | 0.0000 (0.0000) | 0.0000 (0.0000) |
| $k > p > r_0, 50 > 40 > 5$ | 0.0000 (0.0000) | 3.5300 (1.6358) | 1.6400 (0.6439) |
| $p > k > r_0, 50 > 40 > 5$ | 0.0000 (0.0000) | 4.0800 (1.3830) | 1.6500 (0.5925) |
| $p > k > r_0, 99 > 40 > 5$ | 0.0000 (0.0000) | 2.9800 (1.9121) | 0.2200 (0.4399) |

As suggested in Tables 2 and 3, our exact method is much more precise in prediction in all cases and rank recovery except one case of $k = r_0 = 5 < p = 99$, than the other two methods. This is in agreement with the theoretical results in Theorem 4 that exact reconstruction of the oracle estimator is achieved through tuning. Note that the best MSE value does not necessarily yield best rank recovery, as in the case of $k = r_0 = 5 < p = 99$, which is due to the bias/variance trade-off phenomenon. As indicated in Table 4, our method is, on average, 10-20 times faster than the other two.

7.2 MIT logo Recovery

Next we examine the three competing methods for reconstructing the MIT logo image, which was studied in [27, 17]. The original logo is displayed in Figure 2, where we use the gray image of size 44×85 and the image matrix has rank 7.



Figure 2: Original MIT logo image.

Our objective is to reconstruct this image from its noisy version and examine the quality of reconstruction. Towards this end, we sample the design matrix \mathbf{A} with each entry being iid $N(0, 1)$, where the sample size n ranges from 20 to 80. To generate a noisy version, we add random error sampled from $N(0, 0.5^2)$ to each element of the sample data. The reconstruction results are displayed in Figure 3 for the three methods with the default initial value $\mathbf{0}$ for the SVP.

Visually, our method delivers highly competitive performance as compared to the other two methods, as displayed in Figure 3, and yields nearly perfect reconstruction when the size of the design matrix \mathbf{A} becomes larger, say $n = 60$ and $n = 80$. For all the methods, better reconstruction can be reached as n increases, and comparable results are

achieved when n is small, say $n = 20$ and $n = 40$. This conclusion is consistent with our theoretical analysis.

In practice, the exact rank of the matrix to be estimated is unknown but may reasonably be assumed to be small. In this sense, s needs to be tuned or estimated. Next we investigate the effect of estimation of s with regard to the reconstruction quality, measured by the MSE and the recovery rate. The latter is defined as the ratio $\|\hat{\Theta}^s - \Theta^0\|_F / \|\Theta^0\|_F$, commonly used to measure the quality of parameter estimation. In particular, we display both of them as a function of s in Figures 4 and 5, where s is defined in (3).

For the relative recovery error, a clear transition of our solution occurs around $s = 44$, after which a perfect recovery is achieved, whereas no improvement occurs for the other two methods when s increases. In the case of an under-determined \mathbf{A} , i.e., \mathbf{A} is not of full column rank, all three methods produce similar recovery results. For the MSEs, our method yields the perfect result zero when $s \geq 7$ and a reasonably small value when $s = 3, 5$, whereas the other two methods lead to elevated MSEs as s increases. This is in accordance with our theory which suggests that our method constructs the oracle estimator giving perfect image reconstruction when $s \leq 7$.

8. CONCLUSION

This paper considers a nonconvex least squares formulation based on the rank constraint/regularization. We establish the optimality of the global solution, against any other ones. Experimental results on synthetic and real data demonstrate the efficiency and effectiveness of the proposed algorithm. In future work, we plan to expand the present work to a general loss function.

9. ACKNOWLEDGEMENT

This work was supported by NSF grants IIS-0953662 and DMS-0906616, NIH grants R01LM010730, 2R01GM081535-01, and R01HL105397.

10. REFERENCES

- [1] J. Abernethy, F. Bach, T. Evgeniou, and J. Vert.

Table 4: Run time for the synthetic data: average training time in seconds for three competing methods based on the selected tuning parameters over 100 simulation replications. Our, SVP, and Nuclear-Norm refer to our method, the SVP method and nuclear-norm regularization.

| Set-up | Our | SVP | Nuclear-Norm |
|----------------------------|-----------------|------------------|-----------------|
| $k = r_0 < p, 5 = 5 < 99$ | 0.0224 (0.0045) | 14.1523 (0.3933) | 0.9766 (0.0074) |
| $k > p > r_0, 50 > 40 > 5$ | 0.0349 (0.0031) | 41.2122 (1.8511) | 5.8937 (0.2945) |
| $p > k > r_0, 50 > 40 > 5$ | 0.0368 (0.0025) | 29.5856 (1.1786) | 5.7447 (0.1173) |
| $p > k > r_0, 99 > 40 > 5$ | 0.0542 (0.0039) | 79.9850 (3.6724) | 7.6046 (0.1011) |



Figure 3: Reconstruction of a noisy version of the MIT logo with varying sampling size n . From left to right (for each case of n): SVP; nuclear-norm regularization; our method.

- Low-rank matrix factorization with attributes. *Arxiv preprint cs/0611124*, 2006.
- [2] Y. Amit, M. Fink, N. Srebro, and S. Ullman. Uncovering shared structures in multiclass classification. In *Proceedings of the 24th Annual International Conference on Machine Learning*, pages 17–24. ACM, 2007.
- [3] T. Anderson. Estimating linear restrictions on regression coefficients for multivariate normal distributions. *The Annals of Mathematical Statistics*, 22(3):327–351, 1951.
- [4] T. Anderson. Asymptotic distribution of the reduced rank regression estimator under general conditions. *The Annals of Statistics*, 27(4):1141–1154, 1999.
- [5] T. André, R. Nowak, and B. Van Veen. Low-rank estimation of higher order statistics. *Signal Processing, IEEE Transactions on*, 45(3):673–685, 1997.
- [6] A. Argyriou, T. Evgeniou, and M. Pontil. Multi-task feature learning. *Advances in Neural Information Processing Systems*, 19:41, 2007.
- [7] F. Bunea, Y. She, and M. Wegkamp. Optimal selection of reduced rank estimators of high-dimensional matrices. *The Annals of Statistics*, 39(2):1282–1309, 2011.
- [8] J. Cai, E. Candès, and Z. Shen. A singular value thresholding algorithm for matrix completion. *SIAM Journal of Optimization*, 20(4):1956–1982, 2010.
- [9] E. Candès and Y. Plan. Matrix completion with noise. *Arxiv preprint arXiv:0903.3131*, 2009.
- [10] E. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772, 2009.
- [11] J. Chen, J. Liu, and J. Ye. Learning incoherent sparse and low-rank patterns from multiple tasks. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1179–1188. ACM, 2010.
- [12] C. Eckart and G. Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936.
- [13] M. Fazel, H. Hindi, and S. Boyd. A rank minimization heuristic with application to minimum order system approximation. In *American Control Conference, 2001. Proceedings of the 2001*, volume 6, pages 4734–4739. IEEE, 2001.
- [14] S. Friedland and A. Torokhti. Generalized rank-constrained matrix approximations. *SIAM Journal on Matrix Analysis and Applications*, 29(2):656–659, 2007.
- [15] A. Izenman. Reduced-rank regression for the multivariate linear model. *Journal of multivariate analysis*, 5(2):248–264, 1975.
- [16] M. Jaggi and M. Suvolsky. A Simple Algorithm for Nuclear Norm Regularized Problems. In *Proceedings of the 27th Annual International Conference on Machine Learning*, 2010.
- [17] P. Jain, R. Meka, and I. Dhillon. Guaranteed rank minimization via singular value projection. *Advances in Neural Information Processing Systems*, 23:937–945, 2010.
- [18] H. Ji, C. Liu, Z. Shen, and Y. Xu. Robust video denoising using low rank matrix completion. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 1791–1798. IEEE, 2010.
- [19] S. Ji and J. Ye. An accelerated gradient method for trace norm minimization. In *Proceedings of the 26th*

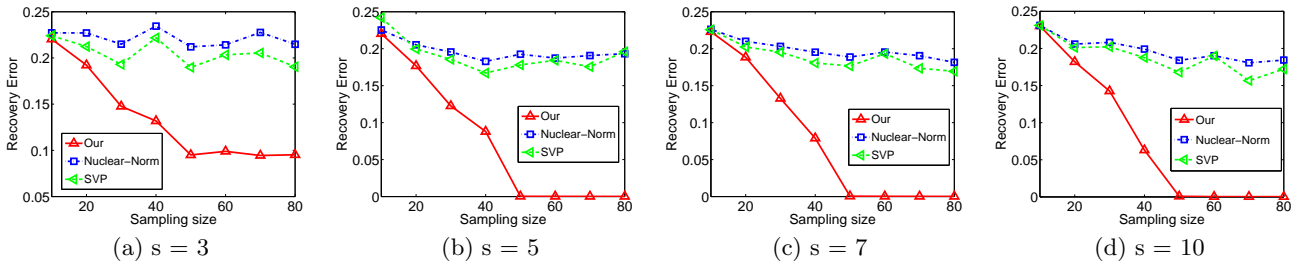


Figure 4: Relative recovery error as a function of the sampling size for the MIT logo image under different rank constraints.

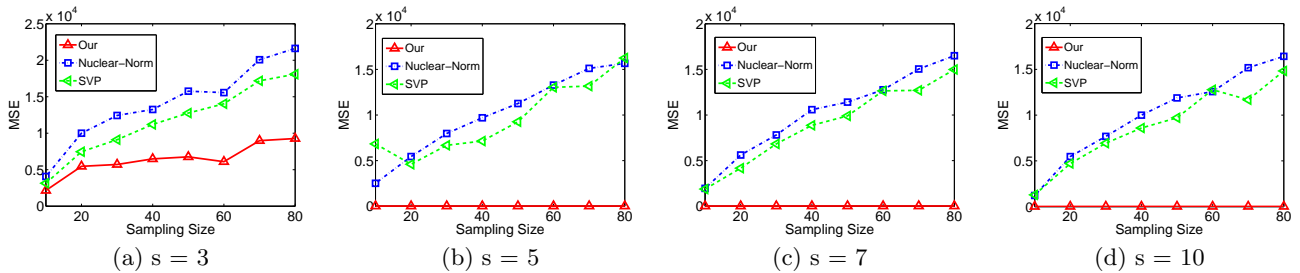


Figure 5: MSE as a function of the sampling size for the MIT logo image under different rank constraints. Note that the MSE of our method in (b) is not zero but indistinguishable from zero, which is unlike (c) and (d) in which the MSEs are identical to zero.

Annual International Conference on Machine Learning, pages 457–464, 2009.

- [20] B. Kulis, A. Surendran, and J. Platt. Fast low-rank semidefinite programming for embedding and clustering. In *Eleventh International Conference on Artificial Intelligence and Statistics, AISTATS 2007*, 2007.
- [21] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma. Robust recovery of subspace structures by low-rank representation. *Arxiv preprint arXiv:1010.2955*, 2010.
- [22] X. Luo. High dimensional low rank and sparse covariance matrix estimation via convex minimization. *Arxiv preprint arXiv:1111.1133*, 2011.
- [23] B. K. Natarajan. Sparse approximation solutions to linear systems. *SIAM J. Comput.*, 24(2):227–234, 1995.
- [24] S. Negahban and M. J. Wainwright. Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *Annals of Statistics*, 39(2):1069–1097, 2011.
- [25] T. Pong, P. Tseng, S. Ji, and J. Ye. Trace norm regularization: Reformulations, algorithms, and multi-task learning. *SIAM Journal on Optimization*, 20:3465–3489, 2010.
- [26] C. Rao. Matrix approximations and reduction of dimensionality in multivariate statistical analysis. *Multivariate analysis*, 5:3–22, 1980.
- [27] B. Recht, M. Fazel, and P. Parrilo. Guaranteed Minimum-Rank Solutions of Linear Matrix Equations via Nuclear Norm Minimization. *SIAM Review*, 52(471), 2010.
- [28] G. Reinsel and R. Velu. *Multivariate reduced-rank regression: theory and applications*. Springer New York, 1998.
- [29] B. Savas and I. Dhillon. Clustered low rank approximation of graphs in information science applications. In *Proceedings of the SIAM Data Mining Conference.*, 2011.
- [30] S. Shalev-Shwartz, A. Gonen, and O. Shamir. Large-Scale Convex Minimization with a Low-Rank Constraint. In *Proceedings of the 28th Annual International Conference on Machine Learning*, 2011.
- [31] X. Shen and H. Huang. Optimal model assessment, selection, and combination. *Journal of the American Statistical Association*, 101(474):554–568, 2006.
- [32] D. Sondermann. Best approximate solutions to matrix equations under rank restrictions. *Statistical Papers*, 27(1):57–66, 1986.
- [33] N. Srebro, J. Rennie, and T. Jaakkola. Maximum-margin matrix factorization. *Advances in Neural Information Processing Systems*, 17:1329–1336, 2005.
- [34] K. Toh and S. Yun. An accelerated proximal gradient algorithm for nuclear norm regularized linear least squares problems. *Pacific J. Optim*, 6:615–640, 2010.
- [35] J. Wright, A. Ganesh, S. Rao, and Y. Ma. Robust principal component analysis: Exact recovery of corrupted low-rank matrices via convex optimization. *Advances in Neural Information Processing Systems*, 2009.