
Solving the Robust Matrix Completion Problem via a System of Nonlinear Equations

Yunfeng Cai and Ping Li

Cognitive Computing Lab
Baidu Research

No. 10 Xibeiwang East Road, Beijing 100085, China
10900 NE 8th St. Bellevue, WA 98004, USA
{caiyunfeng, liping11}@baidu.com

Abstract

We consider the problem of robust matrix completion, which aims to recover a low rank matrix L_* and a sparse matrix S_* from incomplete observations of their sum $M = L_* + S_* \in \mathbb{R}^{m \times n}$. Algorithmically, the robust matrix completion problem is transformed into a problem of solving a system of nonlinear equations, and the alternative direction method is then used to solve the nonlinear equations. In addition, the algorithm is highly parallelizable and suitable for large scale problems. Theoretically, we characterize the sufficient conditions for when L_* can be approximated by a low rank approximation of the observed M_* . And under proper assumptions, it is shown that the algorithm converges to the true solution linearly. Numerical simulations show that the simple method works as expected and is comparable with state-of-the-art methods.

1 Introduction

Robust matrix completion (RMC) (Chen et al., 2011; Tao and Yuan, 2011; Cherapanamjeri et al., 2017; Klopp et al., 2017; Zeng and So, 2018) aims to recover a low rank matrix L_* and a sparse matrix S_* from a sampling of $M = L_* + S_*$. Mathematically, RMC can be formulated as the following optimization problem (Candès et al., 2011; Chandrasekaran et al., 2011):

$$\begin{aligned} \text{RMC : } & \min_{L, S} \text{rank}(L) + \lambda |\text{supp}(S)|, \\ \text{s.t. } & L_{ij} + S_{ij} = M_{i,j}, \quad (i, j) \in \Omega, \end{aligned}$$

where λ is a tuning parameter, Ω is a subset of $\{1, \dots, m\} \times \{1, \dots, n\}$. When $S = 0$, RMC becomes the matrix completion (MC) problem (Candès and Recht, 2009; Meka et al., 2009; Cai et al., 2010; Candès and Tao, 2010; Jain and Netrapalli, 2015; Liu and Li, 2016); When $\Omega = \{1, \dots, m\} \times \{1, \dots, n\}$, RMC becomes the robust principal component analysis (RPCA) (Jolliffe, 2011). Thus, RMC can be taken as a combination/generalization of MC and RPCA (Candès and Plan, 2010; Jain and Netrapalli, 2015; Jain et al., 2013; Keshavan et al., 2010).

In many scientific and engineering problems, people need to recover a low rank matrix from observed data, e.g., the recommender system (Funk, 2006; Candès and Plan, 2010; Hu and Li, 2017, 2018b,a), social network analysis (Huang et al., 2013), machine learning (Candès and Plan, 2010; Davenport and Romberg, 2016), image inpainting (Bertalmio et al., 2000), computer vision (Candès and Plan, 2010), bioinformatics (Kim et al., 2005), etc.

MC/RPCA/RMC has been studied extensively from an optimization point of view, many algorithms are proposed, and exact recovery is discussed under proper assumptions. Most of well-known algorithms are based on convex optimization, in which the rank of a matrix is relaxed to its nuclear norm (the sum of all singular values), and the number of nonzero entries of a matrix is relaxed to its ℓ_1 -norm (the sum of absolute values of all entries), e.g., (Cai et al., 2010; Candès et al., 2011; Candès and Recht, 2009; Candès and Tao, 2010; Recht et al., 2010; Klopp et al., 2017). However, the computation of the nuclear norm of L , which requires the computation of its singular value decomposition (SVD), is expensive and unsuitable for parallelization, as a result, algorithms based on the nuclear norm relaxation are often not very realistic for large matrices.

To deal with large problems, the low rank matrix can be represented as the product of a tall-skinny matrix

and a short-fat matrix, so that the low rank property is satisfied automatically. But, the optimization problem becomes nonconvex, which makes it difficult to solve. Bi-convex as the problem is, the alternative minimization can be used to solve it more efficiently, e.g., (Jain et al., 2013). In (Yi et al., 2016), gradient descent method is used to solve RMC, which is shown to be fast. In (Cherapanamjeri et al., 2017), projected gradient method with hard-thresholding is used to solve RMC, with nearly-optimal observation and corruption. We refer the readers to (Zeng and So, 2018) and reference therein for more methods. Besides the optimization based method, quite recently, in (Dutta et al., 2019), the RPCA/RMC is solved via an alternating nonconvex projection method. This method does not require any objective function, convex relaxation or surrogate convex constraint.

Contribution. In this paper, we solve the RMC problem via solving a system of nonlinear equations (NLEQ). This method does not require any objective function, convex relaxation or surrogate convex constraint, either. Let $L_* = XY^T$, where $X = [x_1, \dots, x_m]^T \in \mathbb{R}^{m \times r}$, $Y = [y_1, \dots, y_n]^T \in \mathbb{R}^{n \times r}$. The RMC problem can be formulated as the following system of NLEQ:

$$x_i^T y_j = M_{ij}, \text{ for } (i, j) \in \Omega \setminus \text{supp}(S_*), \quad (1)$$

where $\text{supp}(S_*)$ is the support set of S_* . When the rank r and the support set of S_* are both known, solving RMC amounts solving a system of nonlinear equations. Any numerical method for nonlinear equations can be used to solve it (e.g., the steepest descent method, the Newton method, etc.), among which the simplest one is the alternative direction method (ADM): Fixing X (or Y), Y (or X) can be updated via solving an (overdetermined) linear system of equations (usually in least square sense). Thus, solving RMC via solving (1) heavily depends on whether we can solving (1) without knowing r and $\text{supp}(S_*)$. It is worth mentioning here that in (Meka et al., 2009), such a method was proposed to solve the MC problem. And it was stated there that “... its variants outperform most methods in practice. However, analyzing the performance of alternate minimization is a notoriously hard problem.”

The contributions of this paper are three folds. First, for both full and partial observation cases, we characterize some sufficient conditions for when the low rank approximation of the observed M is approximate L_* . Second, we propose to solve RMC via solving a system of NLEQ rather than optimization, and an ADM method, which carefully handles the unknown r and $\text{supp}(S_*)$ issue, is developed. Third, we carefully analyze the convergence of the ADM, and it is shown that under proper assumptions, the ADM converges to the true solution linearly, i.e., exact recovery can be

achieved. So, we give an answer to a problem which is even more difficult than the aforementioned “notoriously hard problem”. In addition, the algorithm is highly parallelizable and naturally suitable for large scale problems. It is also worth mentioning here that the results of this paper are applicable to the MC problem as well as the RPCA problem.

The rest of this paper is organized as follows. In Section 2, we first develop the algorithm, followed by its convergence analysis in Section 3. Numerical experiments are presented in Section 4. Concluding remarks are given in Section 5.

Notation. We shall adopt the MATLAB style convention to access the entries of vectors and matrices. The set of integers from i to j inclusive is $i : j$. For a matrix A , its submatrices $A_{(k:\ell, i:j)}$, $A_{(k:\ell, :)}$, $A_{(:, i:j)}$ consist of intersections of row k to row ℓ and column i to column j , row k to row ℓ and all columns, all rows and column i to column j , respectively. $A_{(j, :)}$ and $A_{(:, k)}$ denote the j th row and k th column of A , respectively. $\|A\|$ stands for the spectral norm of A , $\|A\|_F$ denotes the Frobenius norm, $\|A\|_1 = \sum_{i,j} |a_{ij}|$, $\|A\|_{\max} = \max_{i,j} |a_{ij}|$, $\|A\|_{2, \infty} = \max_i \|A_{(i, :)}\|$, A^\dagger stands for the Moore–Penrose inverse, and $\kappa(A) = \|A\| \|A^\dagger\|$ denote the condition number of A . Denote by $\sigma_j(A)$ for $1 \leq j \leq \min\{m, n\}$ the singular values of A and they are always arranged in a non-increasing order: $\sigma_1(A) \geq \sigma_2(A) \geq \dots \geq \sigma_{\min\{m, n\}}(A)$. $\mathcal{R}(A)$ stands for the range space of A , i.e., $\mathcal{R}(A) = \text{span}\{y \in \mathbb{R}^m \mid y = Ax, x \in \mathbb{R}^n\}$. The vector e_j stands for the j th column of the identity matrix I . Furthermore, for an index set $\Omega \subset \{1, \dots, m\} \times \{1, \dots, n\}$, $|\Omega|$ denotes the cardinality of Ω , $\Pi_\Omega(A) = [\mathcal{I}_{\{(i,j) \in \Omega\}} a_{ij}] \in \mathbb{R}^{m \times n}$, where \mathcal{I} is the indicator function.

2 Algorithm

In this section, we first reformulate the RMC problem as a problem of solving a system of nonlinear equations (NLEQ), show how to solve the NLEQ via an alternative direction method (ADM), then the overall algorithm is summarized.

2.1 Problem Reformulation, Difficulty and Solution

Let $M = XY^T$, where $X = [x_1, \dots, x_m]^T \in \mathbb{R}^{m \times r}$, $Y = [y_1, \dots, y_n]^T \in \mathbb{R}^{n \times r}$. The MC problem can be formulated as

$$x_i^T y_j = M_{ij}, \text{ for } (i, j) \in \Omega. \quad (2)$$

Now recall (1), the task of RMC becomes solving (2) with unknown r and minimum violators. Here by a “violator”, denoted by (i', j') , we mean that $x_{i'}^T y_{j'} \neq M_{i'j'}$, it will be also referred to as an “outlier” hereafter.

The ADM is the simplest method to solve the NLEQ of form (2) – given an initial guess for X , we fix X , (2) becomes a linear system in Y (which is assumed to be overdetermined), we can solve the linear system for Y ; similarly, we fix Y to solve X ; the iteration continues until convergence. As r and $\text{supp}(S_*)$ are unknown, our task is to determine them during the iteration of ADM. Next, we first show how to determine r , then $\text{supp}(S_*)$.

Determine the rank adaptively In the t th iteration of ADM, let the current rank estimation be r_t , Y_t be current estimation for Y and Y_t have orthonormal columns, i.e., $Y_t^T Y_t = I_{r_t}$. By ADM, the estimation of X can be obtained, denote it by \tilde{X}_{t+1} . In order to update r_t , we need to compute the singular values of $\tilde{X}_{t+1} Y_t^T$. Noticing that Y_t is orthonormal, then the top r_t singular values of $\tilde{X}_{t+1} Y_t^T$ will be the singular values of \tilde{X}_{t+1} . Then the singular values can be obtained as follows: first, compute the QR decomposition of \tilde{X}_{t+1} :

$$\tilde{X}_{t+1} = \hat{X}_{t+1} R_{x,t+1},$$

where $\hat{X}_{t+1} \in \mathbb{R}^{m \times r_t}$ is orthonormal, $R_{x,t+1} \in \mathbb{R}^{r_t \times r_t}$; second, compute the SVD of $R_{x,t+1}$:

$$R_{x,t+1} = Q_x \Sigma Q_y^T,$$

where $Q_x \in \mathbb{R}^{r_t \times r_t}$, $Q_y \in \mathbb{R}^{n \times r_t}$ are orthogonal, $\Sigma = \text{diag}(\hat{\sigma}_1, \dots, \hat{\sigma}_{r_t})$ with $\hat{\sigma}_1 \geq \dots \geq \hat{\sigma}_{r_t} \geq 0$. Then the singular values of $\tilde{X}_{t+1} Y_t^T$ are $\hat{\sigma}_1, \dots, \hat{\sigma}_{r_t}$. Similarly, when X_t is the current estimation for X , and $X_t^T X_t = I_{r_t}$, we can compute an estimation for Y , its singular values can be obtained.

Let L_t be the current estimation for L_* . When the rank is underestimated, i.e., $r_t < r$, the residual $\tau_t = \|\Pi_\Omega(L_t + S_t - M)\|_F$ will stagnate, in such case, we increase the estimated rank r_t . When the rank is overestimated, i.e., $r_t > r$, we expect to observe rank deficiency from the singular values $\hat{\sigma}_1, \dots, \hat{\sigma}_{r_t}$. In such a case, we decrease the estimated rank r_t . For the RMC problem, we prefer an overestimated rank over an underestimated rank due to the following reason. The residual τ_t stagnates for two reasons: one is that the estimated rank is smaller than the true rank; the other is that $|\text{supp}(S_*) \setminus \text{supp}(S_t)|$ is large. Then when the residual stagnates, it is difficult for us to make a good choice – to increase the estimated rank or to drop some equalities (of course, those equalities need to be carefully selected) in (2). Increasing the estimated rank when $|\text{supp}(S_*) \setminus \text{supp}(S_t)|$ is large or dropping equalities in (2) when the rank is underestimated will both lead to catastrophic consequences, such as the estimated rank exceeds a prescribed limit, too many “correct” equalities are dropped which will probably result in underdetermined linear systems when updating X (or Y). With an overestimated rank, when the

residual stagnates, we decrease the estimated rank via the singular values of L_t ; if there is no rank deficiency in L_t , we drop some equalities in (2).

When an overestimated rank decreases to the actual rank, it is expected that the estimated rank will remain unchanged in the follow-up iterations. Therefore, we do not need to check the singular values of L_t in each iteration for the sake of efficiency.

Determine $\text{supp}(S_*)$ via outlier detection When a good approximation \hat{L} of L_* is obtained, $S_* = M - L_* \approx M - \hat{L}$. Thus, it is reasonable to detect $(i, j) \in \Omega \cap \text{supp}(S_*)$ from the residual $\{R_{ij} = M_{ij} - \hat{L}_{ij}\}_{(i,j) \in \Omega}$.

Outlier detection has been used for centuries to remove abnormal data. Various outlier detection techniques have been used (Ester et al., 1996; Hodge and Austin, 2004; Xu et al., 2010; Rahmani and Li, 2019; Slawski et al., 2019). In our implementation, we simply determine the outliers as follows: find the top- k values in each row and column of $|R|$ (unavailable entries of $|R|$ are set to zero), and the entries in the intersection are taken as outliers. Alternatively, simply find the top k' values among all entries of $|R|$. Here k, k' are two parameters which can be tuned. In what follows, we denote

$$\mathcal{T}_s(A) = [b_{ij}], \quad (3)$$

where s is the number of the removed outliers, $b_{ij} = A_{(i,j)}$ if $A_{(i,j)}$ is an outlier, $b_{ij} = 0$, otherwise. Of course, one can also try other outlier detection techniques.

2.2 Algorithm details

Now we present Algorithm 1, which summarizes the ADM for RMC described in the previous subsection.

Some implementation details follows.

Initializing Y_0 According to Theorem 2 below, good initial guesses for X and Y can be obtained by computing the SVD of $\Pi_{\Omega_0}(M)$. An iterative procedure (e.g., Krylov subspace method) is usually adopted to accomplish the task, in which matrix vector products $\Pi_{\Omega_0}(M)v$ and $\Pi_{\Omega_0}(M)^T v$ are called several times. A simpler way, which is more efficient and numerically proven to be reliable, is the following: compute $W = \Pi_{\Omega_0}(M)^T \Pi_{\Omega_0}(M)N$, compute an orthonormal basis for W , and set the columns of Y_1 as the basis. Here $N \in \mathbb{R}^{n \times r_0}$ is a random matrix with entries drawn from the standard normal distribution. Such a procedure is essentially one iteration of the subspace method (a generalization of power method to compute several dominant eigenvectors). Since an initial guess for X or Y is sufficient for ADM to run in Algorithm 1, it is indeed unnecessary to compute the estimations for both X and Y .

Algorithm 1 ADM FOR RMC VIA NLEQ

Input: The observed matrix $\Pi_\Omega(M)$, a sparsity level parameter s , an estimated rank r_0 , an upper bound κ for the condition number of L_* , and a tolerance tol .

Output: $X \in \mathbb{R}^{m \times r_t}$, $Y \in \mathbb{R}^{n \times r_t}$ and $S \in \mathbb{R}^{m \times n}$ such that $\|\Pi_\Omega(XY^\top + S - M)\|_F \leq \text{tol}$, $\|S\|_0 \leq s$.

- 1: Set $S_0 = \mathcal{T}_s(M)$, $X_0 = 0$, $Y_0 = 0$, $\Sigma_0 = 0$, $t=1$;
 - 2: Compute $[X_1, \Sigma_1, Y_1] = \text{SVD}_{r_0}((M - S_0)/p')$, where $p' = (|\Omega| - s)/mn$;
 - 3: Compute $R_t = \Pi_\Omega(M - X_t \Sigma_t Y_t^\top)$;
 - 4: Set $S_t = \mathcal{T}_s(R_t)$, $\Omega_t = \Omega \setminus \text{supp}(S_t)$;
 - 5: Compute $\tau_t = \|\Pi_\Omega(M - X_t \Sigma_t Y_t^\top - S_t)\|_F$;
 - 6: **while** $\tau_t > \text{tol}$ **do**
 - 7: Set $t = t + 1$;
 - 8: Solve $\Pi_{\Omega_{t-1}}(\tilde{X}_t Y_{t-1}^\top) = \Pi_{\Omega_{t-1}}(M)$ for \tilde{X}_t ;
 - 9: Compute the QR decomposition $\tilde{X}_t = \tilde{X}_t R_{x,t}$, where \tilde{X}_t has orthonormal columns, $R_{x,t}$ is upper triangular;
 - 10: Compute the SVD $R_{x,t} = Q_x \hat{\Sigma} Q_y^\top$, where $\hat{\Sigma} = \text{diag}(\hat{\sigma}_1, \dots, \hat{\sigma}_{r_{t-1}})$, Q_x, Q_y are orthogonal;
 - 11: Set $r_t = r_{t-1} - |\{j \mid \kappa \hat{\sigma}_j < \hat{\sigma}_1\}|$;
 - 12: Set $\tilde{X}_t = [\tilde{X}_t Q_x]_{(:,1:r_t)}$;
 - 13: Solve $\Pi_{\Omega_{t-1}}(\tilde{X}_t \tilde{Y}_t^\top) = \Pi_{\Omega_{t-1}}(M)$ for \tilde{Y}_t ;
 - 14: Compute the QR decomposition $\tilde{Y}_t = \tilde{Y}_t R_{y,t}$, where \tilde{Y}_t has orthonormal columns, $R_{y,t}$ is upper triangular;
 - 15: Compute the SVD $R_{y,t}^\top = Q_x \hat{\Sigma} Q_y^\top$, where $\hat{\Sigma} = \text{diag}(\hat{\sigma}_1, \dots, \hat{\sigma}_{r_t})$, Q_x, Q_y are orthogonal;
 - 16: Set $r_t = r_t - |\{j \mid \kappa \hat{\sigma}_j < \hat{\sigma}_1\}|$;
 - 17: Set $X_t = [\tilde{X}_t Q_x]_{(:,1:r_t)}$, $Y_t = [\tilde{Y}_t Q_y]_{(:,1:r_t)}$, $\Sigma_t = \hat{\Sigma}_{(1:r_t, 1:r_t)}$;
 - 18: Compute $R_t = \Pi_\Omega(M - X_t \Sigma_t Y_t^\top)$;
 - 19: Set $S_t = \mathcal{T}_s(R_t)$, $\Omega_t = \Omega \setminus \text{supp}(S_t)$;
 - 20: Compute $\tau_t = \|\Pi_\Omega(M - X_t \Sigma_t Y_t^\top - S_t)\|_F$;
 - 21: **end while**
-

Solving X_t and Y_t On Lines 8 and 13, \tilde{X}_t and \tilde{Y}_t can both be solved row by row or simultaneously. And to obtain one row of \tilde{X}_t or \tilde{Y}_t , a small linear system needs to be solved. When the linear system is underdetermined, Algorithm 1 may break down. Therefore, in each row and column, Algorithm 1 requires the number of observed entries (after the removal of the corrupted entries) must be larger than the rank. To be more precise, we need the small linear system to be good conditioned. In general, it is difficult to determine how many rows/columns are needed to ensure the linear system to be good conditioned. Numerically, for a random matrix $A \in \mathbb{R}^{s \times r}$ (generated from a standard normal distribution) with $s = \mathcal{O}(r) > 2r$ is usually good conditioned. So, we may declare that $\mathcal{O}(r) > 2r$ observations in each row and column are sufficient.

In our implementation, the linear systems are solved in the least square sense. One may also choose to minimize ℓ_p -norm ($p \geq 0$) of the residual as in (Zeng and So, 2018).

Computational complexity When the number of observations in each row and column is $\mathcal{O}(r)$, each linear system can be solved in $\mathcal{O}(r^3)$ FLOPS. So, in each iteration, the computational complexity of the linear system solving on Lines 8 and 13 is $\mathcal{O}((m+n)r^3)$. The computational complexity of the QR decomposition on lines 9 and 14 is $\mathcal{O}((m+n)r^2)$. The computational complexity of the SVD is $\mathcal{O}(r^3)$. So, the overall of computational complexity of Algorithm 1 is dominated by the linear system solving. When the number of observations in certain row/column is much larger than r , we may randomly choose $\mathcal{O}(r)$ observations from the row/column, then solve a much smaller linear system of equations. Again, the overall computational complexity in each iteration is $\mathcal{O}((m+n)r^3)$.

Also, note that the linear systems on Line 8 and 13 can be solved in parallel. Therefore, Algorithm 1 are suitable for large scale problems.

Remark 1. When Ω_t is fixed, Algorithm 1 essentially minimizes $\|\Pi_{\Omega_t}(XY^\top - M)\|_F$ via ADM. If gradient method is used to minimize $\|\Pi_{\Omega_t}(XY^\top - M)\|_F$, Algorithm 1 is similar to the GD method in (Yi et al., 2016), except the regularization term $\|U_t^\top U_t - V_t^\top V_t\|_F$ in the loss function.

3 Convergence

This section analyzes the convergence of Algorithm 1. We first study the full observation case, which serves as a motivation for the partial observation case next.

To present the results, we need to define the k canonical angles. Let \mathcal{X}, \mathcal{Y} be two k -dimensional subspaces of \mathbb{R}^n . Let $X, Y \in \mathbb{R}^{n \times k}$ be the orthonormal basis matrices of \mathcal{X} and \mathcal{Y} , respectively, i.e.,

$$\mathcal{R}(X) = \mathcal{X}, \quad X^\top X = I_k, \quad \text{and} \quad \mathcal{R}(Y) = \mathcal{Y}, \quad Y^\top Y = I_k.$$

Denote ω_j for $1 \leq j \leq k$ the singular values of $Y^\top X$ in ascending order, i.e., $\omega_1 \leq \dots \leq \omega_k$. The k canonical angles $\theta_j(\mathcal{X}, \mathcal{Y})$ between \mathcal{X} and \mathcal{Y} are defined by

$$0 \leq \theta_j(\mathcal{X}, \mathcal{Y}) := \arccos \omega_j \leq \frac{\pi}{2} \quad \text{for } 1 \leq j \leq k. \quad (4)$$

They are in descending order, i.e., $\theta_1(\mathcal{X}, \mathcal{Y}) \geq \dots \geq \theta_k(\mathcal{X}, \mathcal{Y})$. Set

$$\Theta(\mathcal{X}, \mathcal{Y}) = \text{diag}(\theta_1(\mathcal{X}, \mathcal{Y}), \dots, \theta_k(\mathcal{X}, \mathcal{Y})). \quad (5)$$

In what follows, we sometimes place a vector or matrix in one or both arguments of $\theta_j(\cdot, \cdot)$ and $\Theta(\cdot, \cdot)$ with the meaning that it is about the subspace spanned by the vector or the columns of the matrix argument.

3.1 Full observation case

Theorem 1. Let $M = L_* + S_* \in \mathbb{R}^{m \times n}$ ($m \geq n$), where L_* is low rank, i.e., $r = \text{rank}(L_*) \ll n$. Let the SVD of M be $M = U\Sigma V^T$, where $U = [U_1 | U_2] = [u_1, \dots, u_r | u_{r+1}, \dots, u_m]$, $V = [V_1 | V_2] = [v_1, \dots, v_r | v_{r+1}, \dots, v_n]$ are orthogonal matrices, $\Sigma = \begin{bmatrix} \text{diag}(\Sigma_1, \Sigma_2) \\ 0 \end{bmatrix}$, $\Sigma_1 = \text{diag}(\sigma_1, \dots, \sigma_r)$, $\Sigma_2 = \text{diag}(\sigma_{r+1}, \dots, \sigma_n)$, and $\sigma_1 \geq \dots \geq \sigma_n$. Let $M_r = U_1 \Sigma_1 V_1^T$ be the best rank r approximation of M . Let the economy sized SVD of L be $L_* = U_* \Sigma_* V_*^T$, where $U_* \in \mathbb{R}^{m \times r}$ and $V_* \in \mathbb{R}^{n \times r}$ both have orthonormal columns, $\Sigma_* = \text{diag}(\sigma_{1*}, \dots, \sigma_{r*})$ with $\sigma_{1*} \geq \dots \geq \sigma_{r*} > 0$. If

$$\|(I - U_* U_*^T) S_* (I - V_* V_*^T)\| < \sigma_{r*}, \quad (6a)$$

$$\max\{\|S_* V_*\|, \|S_*^T U_*\|\} < \sigma_r - \sigma_{r+1}, \quad (6b)$$

then

$$\begin{aligned} \max\{\theta_u, \theta_v\} &\leq \eta, \\ \frac{\|L_* - M_r\|_{\max}}{\|L_*\|} &\leq (\|U_*\|_{2,\infty} \theta_v + \|V_*\|_{2,\infty} \theta_u) \\ &\quad + (1 + 3\|U_*\|_{2,\infty} \|V_*\|_{2,\infty}) \theta_u \theta_v, \end{aligned}$$

where $\theta_u = \|\sin \Theta(U_1, U_*)\|$, $\theta_v = \|\sin \Theta(V_1, V_*)\|$ and $\eta = \frac{\max\{\|S_* V_*\|, \|S_*^T U_*\|\}}{\sigma_r - \sigma_{r+1} - \max\{\|S_* V_*\|, \|S_*^T U_*\|\}}$.

Theorem 1 tells that when (6) holds and η is small, the principal angles between U_1 and U_* , V_1 and V_* will be small, and the best rank- r approximation of M is a good approximation of L_* . Notice that, (6) does not necessarily implies $\|S_*\|$ is small (compared with $\|L_*\|$). In fact, we have the following example, in which $\|S_*\|$ is comparable with $\|L_*\|$ and $\eta = 0$.

Example 1. Let $M = L_* + S_*$, $L_* = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$,

$$S_* = \frac{\rho}{4} \begin{bmatrix} 2 & -1 & 0 & \dots & 0 & -1 \\ -1 & 2 & -1 & 0 & \dots & 0 \\ 0 & -1 & 2 & -1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & -1 & 2 & -1 \\ -1 & 0 & \dots & 0 & -1 & 2 \end{bmatrix}, \text{ where } \mathbf{1}_n \text{ is an } n\text{-by-1}$$

vector of ones, n is even, $\rho \in (-1, 1)$ is a real parameter. One can verify that $\|L_*\| = 1$, $\|S_*\| = \rho$, $\|S_* \mathbf{1}_n\| = 0$, the first two singular values of M are $\sigma_1 = 1$ and $\sigma_2 = |\rho|$, and the economy sized SVD of L can be given by $L = U_* \Sigma_* V_*^T$, where $U_* = V_* = \frac{1}{\sqrt{n}} \mathbf{1}_n$, $\Sigma_* = \sigma_{1*} = 1$. Then $\|(I - U_* U_*^T) S_* (I - V_* V_*^T)\| = \|S_*\| = |\rho| < 1 = \sigma_{1*}$, and $\max\{\|S_* V_*\|, \|S_*^T U_*\|\} = 0 < 1 - |\rho| = \sigma_1 - \sigma_2$. In other words, the assumption (6) holds. Noticing that $\eta = \frac{\max\{\|S_* V_*\|, \|S_*^T U_*\|\}}{\sigma_r - \sigma_{r+1} - \max\{\|S_* V_*\|, \|S_*^T U_*\|\}} = 0$, by Theorem 1, we can conclude that $\|\sin \Theta(U_1, U_*)\| = \|\sin \Theta(V_1, V_*)\| = 0$, $M_1 = L_*$, where U_1, V_1 are the

top left and right singular vectors of M , respectively, and M_1 is the best rank-1 approximation of M .

Now let us assume all entries of M are observed, S_* is sufficiently small, and can be taken as a perturbation to L_* . Let Y_t be guess for Y and have orthonormal columns, let us perform one iteration of ADM. For the ease of illustration, also assume that $r_{t+1} = r_t$. Then the iteration reads: **(S1)** $\tilde{X}_{t+1} = M Y_t$; **(S2)** $\hat{X}_{t+1} = M Y_t G_x$, where G_x is such that \hat{X}_{t+1} is orthonormal; **(S3)** $\tilde{Y}_{t+1} = M^T \hat{X}_{t+1}$; **(S4)** $\hat{Y}_{t+1} = M^T \hat{X}_{t+1} G_y$, where G_y is such that \hat{Y}_{t+1} is orthonormal. So, we get

$$\hat{Y}_{t+1} = (M^T M) Y_t G_x G_y, \quad (7)$$

which is just one iteration of subspace iteration (a generation of power iteration) for computing the dominant eigenspace of $M^T M$ (e.g., (Demmel, 1997; Stewart, 2001; Van Loan and Golub, 2012)). In fact, **(S2)** and **(S4)** are iterations for subspaces spanned by the dominant left and right singular vectors of M , respectively. Classic results tell that the subspaces $\mathcal{R}(X_t)$ and $\mathcal{R}(Y_t)$ converge to the subspaces spanned by the left and right singular vectors of M corresponding with the dominant singular values. When the perturbation is small, $\mathcal{R}(X_t)$ and $\mathcal{R}(Y_t)$ are good approximations for $\mathcal{R}(L_*)$ and $\mathcal{R}(L_*^T)$, respectively. In particular, when $S_* = 0$ and $V_*^T Y_t$ is nonsingular, we have $\|\sin \Theta(X_{t+1}, U_*)\| = 0$, $\|\sin \Theta(Y_{t+1}, V_*)\| = 0$, i.e., one iteration of ADM gives the true solution.

3.2 Partial observation case

For the partial observation case, we study the convergence of Algorithm 1 under the following assumptions:

(A1) For L_* , the column and row incoherence conditions with parameter μ hold, i.e.,

$$\max_{1 \leq i \leq m} \|U_*^T e_i\|^2 \leq \frac{\mu r}{m}, \quad \max_{1 \leq j \leq n} \|V_*^T e_j\|^2 \leq \frac{\mu r}{n}.$$

(A2) S_* has at most ρ -fraction nonzero entries per row and column, i.e.,

$$\|S_{*(i,:)}\|_0 \leq \rho n, \quad \|S_{*(:,j)}\|_0 \leq \rho m, \quad \text{for all } i, j.$$

(A3) Each entry of M is observed independently with probability p .

Besides the notations in Algorithm 1, we also adopt the following notations:

$$\theta_{x,t} = \|\sin \Theta(X_t, U_*)\|, \quad \theta_{y,t} = \|\sin \Theta(Y_t, V_*)\|.$$

The following theorem tells that the SVD of the partial observed matrix (after the removal of outlier) indeed

gives good approximation for U_*, V_* . Furthermore, X_1, Y_1 satisfy a incoherence condition with parameter μ_1 , $\|L_* - X_1 \Sigma_1 Y_1^T\|_{\max}$ is bounded.

Theorem 2. Assume (A1), (A2), (A3) and $m \geq n$. Let $M = L_* + S_* \in \mathbb{R}^{m \times n}$ with $\text{rank}(L_*) = r$, S_0 be obtained as in Algorithm 1. Denote $r'_s = \frac{\|S_0 - S_*\|_F^2}{\|S_0 - S_*\|_F^2}$, $\gamma = \frac{2}{1-\varrho} \sqrt{\frac{2\varrho}{r'_s p}}$. If

$$(\xi + \gamma)\mu r \kappa < \frac{1}{6}, \quad (8)$$

then with probability $\geq 1 - 1/m^{10+\log \alpha}$, it holds that

$$\max\{\theta_{x,1}, \theta_{y,1}\} \leq 3(\xi + \gamma)\mu r \kappa, \quad (9)$$

where $\xi = 6\sqrt{\frac{\alpha}{p'n}}$, $\kappa = \frac{\sigma_{1*}}{\sigma_{r*}}$. Further assume $\mu \ll n$ and that there exists a positive constant $\mu'_1 \ll n$ such that

$$(\xi + \gamma)\mu r \kappa \leq \frac{1}{3}\sqrt{\frac{\mu'_1 r}{m}}, \quad (10)$$

then

$$\begin{aligned} \|X_1\|_{2,\infty} &\leq \sqrt{\frac{\mu_1 r}{m}}, & \|Y_1\|_{2,\infty} &\leq \sqrt{\frac{\mu_1 r}{n}}, \\ \|L_* - X_1 \Sigma_1 Y_1^T\|_{\max} &\leq \|L_*\| \left(\sqrt{\frac{\mu r}{m}} \theta_{y,1} + \sqrt{\frac{\mu r}{n}} \theta_{x,1} \right. \\ &\quad \left. + \theta_{x,1} \theta_{y,1} \right) + \mathcal{O}(n^{-3/2}), \end{aligned}$$

where $\mu_1 = 2(\mu + \mu'_1)$.

Remark 2. When there is no corruption, i.e., $\varrho = 0$, then $\gamma = 0$. Furthermore, when $m = \mathcal{O}(n) \gg 1$, since $p'n = \mathcal{O}(1)r$, we know that $\xi = \sqrt{\frac{\alpha}{p'n}}$ is small, the larger $p'n$ is, the smaller ξ is. By Theorem 1, $\theta_{x,1}$ and $\theta_{y,1}$ will be small. In other words, the SVD $\frac{1}{p'} \Pi_{\Omega_0}(M - S_0) = X_1 \Sigma_1 Y_1^T$, gives good approximation for both U_* and V_* , by X_1 and Y_1 , respectively.

The following theorem, which is motivated by (Drineas and Mahoney, 2018, Lemma 55), establish the bridge between the full observation case and the partial observation case. This gives an upper bound for the distance between the least square solutions between the full observation case and the partial observation case.

Theorem 3. Let $m \geq n$, and denote

$$\begin{aligned} X_{\text{opt}} &= \operatorname{argmin}_X \|XY_t^T - (M - S_t)\|, \\ \tilde{X}_{\text{opt}} &= \operatorname{argmin}_X \|\Pi_{\Omega_t}(XY_t^T - (M - S_t))\|. \end{aligned}$$

Assume that Ω_t can be obtained by sampling each entry of M with probability $p' = p(1 - \varrho)$, $\|Y_t\|_{2,\infty} \leq \sqrt{\frac{\mu'r}{n}}$

for some $\mu' > 0$, and

$$\inf_{X \in \mathbb{R}^{m \times r}} \frac{\|\Pi_{\Omega_t}(XY_t^T)\|}{\|X\|} \geq \sigma \quad (11)$$

for some constant $\sigma > 0$. Then w.p. ≥ 0.99 , it holds

$$\|\tilde{X}_{\text{opt}} - X_{\text{opt}}\| \leq \left(\frac{2}{3} \log(m+n) + 5\right) \frac{\sqrt{\mu'r p'}}{\sigma^2} \|R\|_{\max},$$

where $R = (M - S_t)(I - Y_t Y_t^T)$.

Remark 3. The requirement (11) is critical. The parameter σ reflects the condition number of the least square problem on line 8 of Algorithm 1. What's more, the larger p' is, the larger σ is (in particular, if $p' = 1$, $\sigma = 1$), the smaller the distance between \tilde{X}_{opt} and X_{opt} is, which agrees with our intuition. R is the residual for the full observation case, i.e., $R = X_{\text{opt}} Y_t^T - (M - S_t)$. If the residual is small, the distance between \tilde{X}_{opt} and X_{opt} will be small, too.

Definition 1. Define $\mu' \triangleq \max\{\mu_u, \mu_v\}$, where

$$\begin{aligned} \mu_u &\triangleq \sup_{U \in \mathbb{R}^{m \times r}} \left\{ \frac{m}{r} \|U\|_{2,\infty}^2 \mid \|\sin \Theta(U_*, U)\| \leq \theta_{x,1} \right\}, \\ \mu_v &\triangleq \sup_{V \in \mathbb{R}^{n \times r}} \left\{ \frac{n}{r} \|V\|_{2,\infty}^2 \mid \|\sin \Theta(U_*, V)\| \leq \theta_{y,1} \right\}. \end{aligned}$$

By definition of μ' , we know that if $\theta_{x,t} \leq \theta_{x,1}$, $\theta_{y,t} \leq \theta_{y,1}$ for all t , then X_t, Y_t satisfy the incoherence condition with parameter μ' . Recall Theorem 1, under the assumption of (10), $\theta_{x,1}$ and $\theta_{y,1}$ are quite small (at the order of $\frac{1}{\sqrt{m}}$), then we can show that $\mu' \leq \mu_1$, which implies that μ' is not large.

The next theorem establishes the convergence rate for the ADM, which is the key in our proof of Theorem 5.

Theorem 4. Assume that Ω_t can be obtained by sampling each entry of M with probability p' , $\|Y_t\|_{2,\infty} \leq \sqrt{\frac{\mu'r}{n}}$ for some $\mu' > 0$, (11) and

$$\|L_* - X_t \Sigma_t Y_t^T\|_{\max} \leq c \|L_*\| \theta_{y,t} \sqrt{\frac{\mu r}{m}}, \quad (12)$$

for some constants $c > 0$. Denote

$$\begin{aligned} r_s &= \inf_t \frac{\|S_t - S_*\|_F^2}{\|S_t - S_*\|^2}, & \zeta &= \sqrt{\frac{2s\mu r}{mr_s}}, \\ C_{\text{LS}} &= \left(\frac{2}{3} \log(m+n) + 5\right) \frac{\sqrt{\mu'r p'}}{\sigma^2}, \\ C &= C_{\text{LS}}(1 + 2c\sqrt{2p\varrho n})\sqrt{\mu r}, \\ \epsilon &= c\kappa\zeta, & \phi &= \frac{8\epsilon(\kappa + \sqrt{2\epsilon}) + \sqrt{2}C\kappa/\sqrt{m}}{1 - 2\epsilon - C\kappa/\sqrt{m}}. \end{aligned}$$

Further assume $\theta_{y,t} \leq \frac{1}{\sqrt{2}}$, then w.p. ≥ 0.99 ,

$$\theta_{x,t+1} \leq \phi \theta_{y,t}.$$

Remark 4. The assumption (12) is not a strong requirement as it looks. By Theorem 1, (12) is natural for $t = 1$. For general $t > 1$, it can be shown that there exists a constant $c > 0$ such that (12) holds (see supplementary for details), as long as X_t, Y_t satisfy the incoherence condition. In general, the constant c is at the order of $\mathcal{O}(1)$.

Remark 5. The constant r_s is the infimum of the stable rank of $S_t - S_*$. If we take a random matrix, whose entries are i.i.d. drawn from a normal distribution $\mathcal{N}(0, \sigma^2)$, to approximate $S_t - S_*$. Numerically, we found that $r_s = \mathcal{O}(n)$. Therefore, $\epsilon \approx c\kappa\sqrt{2\rho\mu r}$ is a small number as long as ρ is small, μ and κ are not large. When p' is sufficiently large, $C = \mathcal{O}(1)$. Consequently, when $m = \mathcal{O}(n)$ is large, $\phi \ll 1$, which agrees with our conclusion in the full observation case.

Remark 6. Under similar assumptions as in Theorem 4, it can also be shown that $\theta_{y,t+1} \leq \phi\theta_{x,t+1}$. Then $\theta_{y,t+1} \leq \phi^2\theta_{y,t}$. When $\phi < 1$, $\{\theta_{y,t}\}_t$ is a monotonically decreasing sequence. Combining it with the definition of μ' , we know that Y_t satisfies the incoherence condition with parameter μ' . Similarly, X_t also satisfies the incoherence condition.

Theorem 5. Assume (A1), (A2), (A3) and $m \geq n$. Assume that Ω_t can be obtained by sampling each entry of M with probability $p' = p(1 - \rho)$, $r_t \equiv r$ and

$$\inf_{X \in \mathbb{R}^{m \times r}} \frac{\|\Pi_{\Omega_t}(XY_t^T)\|}{\|X\|} \geq \sigma, \quad \inf_{Y \in \mathbb{R}^{n \times r}} \frac{\|\Pi_{\Omega_t}(\hat{X}_t Y^T)\|}{\|Y\|} \geq \sigma$$

for some $\sigma > 0$. Let r_s, ζ, C, ϵ be the same as in Theorem 4. Then with high probability, it holds that

$$\|M - S_{t+1} - X_{t+1}\Sigma_{t+1}Y_{t+1}^T\| \leq \psi \|M - S_t - X_t\Sigma_t Y_t^T\|,$$

where

$$\psi = \frac{2\sqrt{2}(\kappa + 2\epsilon\sqrt{\frac{mT}{m}} + \frac{C\kappa}{\sqrt{m}})(8\epsilon(\kappa + \sqrt{2}\epsilon) + \sqrt{2}\frac{C\kappa}{\sqrt{m}})}{(1 - 4\sqrt{2}\epsilon\sqrt{\frac{mT}{m}})(1 - 2\epsilon - \frac{C\kappa}{\sqrt{m}})}.$$

Remark 7. If $\psi < 1$, then by Theorem 5, $\{\|M - S_t - X_t\Sigma_t Y_t^T\|\}_t$ is a monotonically decreasing sequence. And in limit, with high probability, it holds

$$\lim_{t \rightarrow \infty} \|M - S_t - X_t\Sigma_t Y_t^T\| = 0.$$

Recall the way we determine S_t , we get $(L_* - X_t\Sigma_t Y_t^T)_{ij} = 0$ for any $(i, j) \notin \text{supp}(S_t)$. Then we can show that

$$X_t\Sigma_t Y_t^T \rightarrow L_*, \quad S_t \rightarrow S_*, \quad \text{as } t \rightarrow \infty,$$

i.e., exact recovery is achieved, with high probability.

4 Experiments

We compare Algorithm 1 (NLEQ) with the gradient descent (GD) method (Yi et al., 2016) and the PG-RMC method (Cherapanamjeri et al., 2017). The codes of GD and PG-RMC are obtained from the lrslibrary (Sobral et al., 2016) in Github.¹

4.1 Synthetic Data

We generate the data matrix $M \in \mathbb{R}^{d \times d}$ as follows. The low rank matrix L_* is given by $L_* = U_*V_*^T$, where the entries of $U_*, V_* \in \mathbb{R}^{d \times r}$ are drawn independently from the Gaussian distribution with mean zero, variance $1/d$. Each entry of the sparse matrix S_* are nonzero with probability ρ , and the nonzero entries of S_* are uniformly drawn from $[-\frac{r}{2d}, \frac{r}{2d}]$. Each entry of $M = L_* + S_*$ is observed independently with probability p .

The results are presented in Figure 1. In Figure 1-(a), we draw the relative residual $\frac{\|R_t\|_F}{\|\Pi_{\Omega_t}(M)\|_F}$ and rank estimation r_t vs. iteration number. We can see that when the initial rank is larger than the true rank, as the iteration continues, the true rank can be revealed from the singular values of $X_t(Y_t)^T$, then the rank estimation drops to the true rank; when the residual stagnates (represented by a big solid dot in the plot), outliers are removed and the residual decreases until convergence. In Figure 1-(b), we plot the relative residual $\frac{\|R_t\|_F}{\|\Pi_{\Omega_t}(M)\|_F}$ vs. total CPU time for different p . We can see that Algorithm 1 works for all three cases, the larger p is, the less iteration number is needed. In Figure 1-(c), we plot total CPU time and the angle $\max\{\|\sin \Theta(X_t, U_*)\|, \|\sin \Theta(Y_t, V_*)\|\}$ vs. matrix size d for different methods. We can see that the CPU time of all three methods grows linearly with respect to the matrix size, and are comparable with each other. The angles of the three methods are all small, which confirm that all methods give the correct results; the angle produced by NLEQ is the smallest. In Figure 1-(d), we plot the relative residual $\frac{\|R_t\|_F}{\|\mathcal{P}_{\Omega_t}(M)\|}$ vs. CPU time for all three methods. The convergence behaviors of three methods are quite different: GD converges almost linearly; PG-RMC at the beginning stage converges linearly with a low converge rate, then converges almost linearly with a larger rate; NLEQ has a zig-zag convergence, which is due to the removal of outliers.

4.2 Foreground-background separation

The next task is foreground-background separation. By stacking up the vectorized video frames, we get a full data matrix. The static background will form a low rank matrix while the foreground can be taken as

¹<https://github.com/andrewssobral/lrslibrary/tree/master/algorithms/mc>

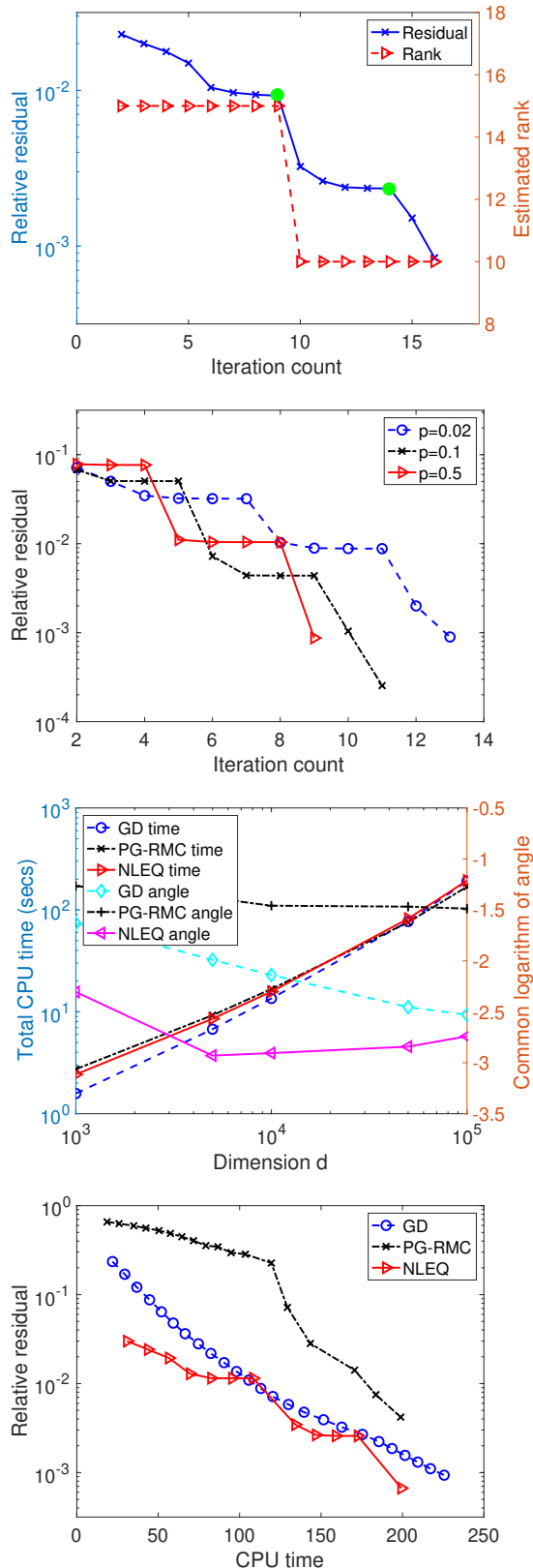


Figure 1: Results on synthetic data. Up to down: (a)-(d). (a) $d = 1e5$, $r = 10$, $p = 0.0015$, $\rho = 0.1$. (b) $d = 1e4$, $r = 10$, $\rho = 0.1$, $p = 0.02, 0.1, 0.5$. (c) $d = 1e3, 5e3, 1e4, 5e4, 1e5$, $r = 10$, $p = 0.15r^2 \log(m)/m$, $\rho = 0.1$. (d) $d = 1e5$, $r = 10$, $p = 0.002$, $\rho = 0.1$.

the sparse component. We apply our method NLEQ, and also GD and PG-RMC to two public benchmarks, the *Bootstrap* and *ShoppingMall*.² Each entry of the data matrix is observed independently w.p. $p = 0.05$. As presented in Figure 2, all three methods are able to separate the foreground from the background, and the backgrounds obtained by three methods are similar.

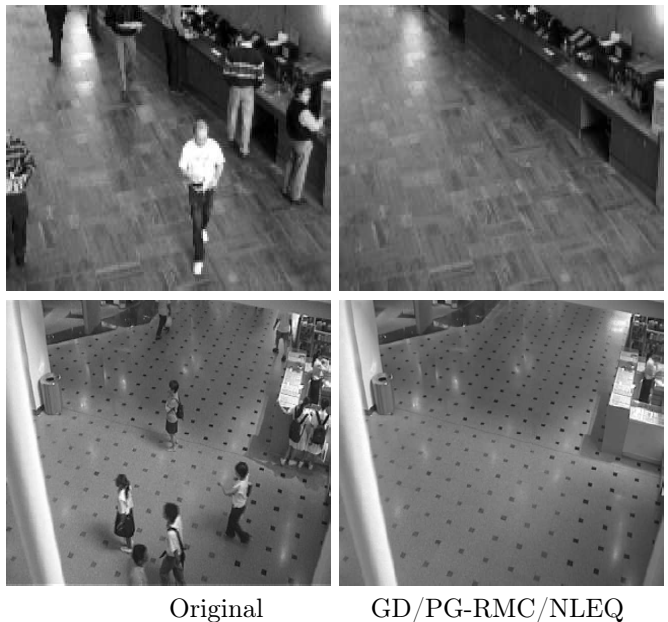


Figure 2: Foreground-background separation.

5 Conclusion

In this paper, we study the RMC problem from an algebraic point of view – transform the RMC problem into a problem of solving an overdetermined nonlinear system of equations (with outliers). This method does not require any objective function, convex relaxation or surrogate convex constraint. Algorithmically, we propose to solve the NLEQ via ADM, in which the true rank and support set of the corruption are determined during the iteration. The algorithm is highly parallelizable and suitable for large scale problems. Theoretically, we characterize the sufficient conditions for when L_* can be approximated by the low rank approximation of M or $\frac{1}{p}\Pi_{\Omega}(M)$. We establish sufficient conditions for $M_r \approx L_*$, where M_r is the best rank r approximation of the observed M . The convergence of the algorithm is guaranteed, and exact recovery is achieved under proper assumptions. Numerical simulations show that the algorithm is comparable with state-of-the-art methods in terms of efficiency and accuracy.

²<http://vis-www.cs.umass.edu/~narayana/castanza/I2Rdataset/>

References

- Marcelo Bertalmío, Guillermo Sapiro, Vicent Caselles, and Coloma Ballester. Image inpainting. In *Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*, pages 417–424, New Orleans, LA, 2000.
- Jian-Feng Cai, Emmanuel J Candès, and Zuowei Shen. A singular value thresholding algorithm for matrix completion. *SIAM J. Optim.*, 20(4):1956–1982, 2010.
- Emmanuel J. Candès and Yaniv Plan. Matrix completion with noise. *Proceedings of the IEEE*, 98(6):925–936, 2010.
- Emmanuel J Candès and Benjamin Recht. Exact matrix completion via convex optimization. *Found. Comput. Math.*, 9(6):717, 2009.
- Emmanuel J. Candès and Terence Tao. The power of convex relaxation: near-optimal matrix completion. *IEEE Trans. Information Theory*, 56(5):2053–2080, 2010.
- Emmanuel J. Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *J. ACM*, 58(3):11:1–11:37, 2011.
- Venkat Chandrasekaran, Sujay Sanghavi, Pablo A Parrilo, and Alan S Willsky. Rank-sparsity incoherence for matrix decomposition. *SIAM J. Optim.*, 21(2):572–596, 2011.
- Yudong Chen, Huan Xu, Constantine Caramanis, and Sujay Sanghavi. Robust matrix completion and corrupted columns. In *Proceedings of the 28th International Conference on Machine Learning (ICML)*, pages 873–880, Bellevue, WA, 2011.
- Yeshwanth Cherapanamjeri, Kartik Gupta, and Prateek Jain. Nearly optimal robust matrix completion. In *Proceedings of the 34th International Conference on Machine Learning, (ICML)*, pages 797–805, Sydney, Australia, 2017.
- Mark A. Davenport and Justin K. Romberg. An overview of low-rank matrix recovery from incomplete observations. *J. Sel. Topics Signal Processing*, 10(4):608–622, 2016.
- Chandler Davis and William Morton Kahan. The rotation of eigenvectors by a perturbation. III. *SIAM J. Numer. Anal.*, 7(1):1–46, 1970.
- James W Demmel. *Applied Numerical Linear Algebra*. SIAM, Philadelphia, PA, 1997.
- Petros Drineas and Michael W Mahoney. Lectures on randomized numerical linear algebra. *The Mathematics of Data*, 25:1, 2018.
- Aritra Dutta, Filip Hanzely, and Peter Richtárik. A nonconvex projection method for robust PCA. In *The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI)*, pages 1468–1476, Honolulu, HI, 2019.
- Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 226–231, Portland, OR, 1996.
- Simon Funk. Netflix update: Try this at home, 2006.
- Victoria Hodge and Jim Austin. A survey of outlier detection methodologies. *Artificial intelligence review*, 22(2):85–126, 2004.
- Jun Hu and Ping Li. Decoupled collaborative ranking. In *Proceedings of the 26th International Conference on World Wide Web (WWW)*, pages 1321–1329, Perth, Australia, 2017.
- Jun Hu and Ping Li. Collaborative multi-objective ranking. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management (CIKM)*, pages 1363–1372, Torino, Italy, 2018a.
- Jun Hu and Ping Li. Collaborative filtering via additive ordinal regression. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining (WSDM)*, pages 243–251, Marina Del Rey, CA, 2018b.
- Jin Huang, Feiping Nie, Heng Huang, Yu Lei, and Chris H. Q. Ding. Social trust prediction using rank-k matrix recovery. In *Proceedings of the 23rd International Joint Conference on Artificial Intelligence (IJCAI)*, pages 2647–2653, Beijing, China, 2013.
- Prateek Jain and Praneeth Netrapalli. Fast exact matrix completion with finite samples. In *Proceedings of The 28th Conference on Learning Theory (COLT)*, pages 1007–1034, Paris, France, 2015.
- Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating minimization. In *Symposium on Theory of Computing Conference (STOC)*, pages 665–674, Palo Alto, CA, 2013.
- Ian Jolliffe. *Principal component analysis*. Springer, 2011.
- Raghunandan H. Keshavan, Andrea Montanari, and Sewoong Oh. Matrix completion from a few entries. *IEEE Trans. Information Theory*, 56(6):2980–2998, 2010.
- Hyunsoo Kim, Gene H. Golub, and Haesun Park. Missing value estimation for DNA microarray gene expression data: local least squares imputation. *Bioinformatics*, 21(2):187–198, 2005.
- Olga Klopp, Karim Lounici, and Alexandre B Tsybakov. Robust matrix completion. *Probability Theory and Related Fields*, 169(1-2):523–564, 2017.

- Guangcan Liu and Ping Li. Low-rank matrix completion in the presence of high coherence. *IEEE Trans. Signal Processing*, 64(21):5623–5633, 2016.
- Raghu Meka, Prateek Jain, and Inderjit S. Dhillon. Matrix completion from power-law distributed samples. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1258–1266, Vancouver, Canada, 2009.
- Mostafa Rahmani and Ping Li. Outlier detection and robust PCA using a convex measure of innovation. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 14200–14210, Vancouver, Canada, 2019.
- Benjamin Recht, Maryam Fazel, and Pablo A Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Rev.*, 52(3):471–501, 2010.
- Martin Slawski, Mostafa Rahmani, and Ping Li. A sparse representation-based approach to linear regression with partially shuffled labels. In *Proceedings of the Thirty-Fifth Conference on Uncertainty in Artificial Intelligence (UAI)*, page 7, Tel Aviv, Israel, 2019.
- Andrews Sobral, Thierry Bouwmans, and El-hadi Zahzah. Lrslibrary: Low-rank and sparse tools for background modeling and subtraction in videos. *Robust Low-Rank and Sparse Matrix Decomposition: Applications in Image and Video Processing*, 2016.
- Gilbert W Stewart. *Matrix algorithms volume 2: eigen-systems*, volume 2. SIAM, 2001.
- Gilbert W Stewart and Ji-Guang Sun. *Matrix Perturbation Theory*. Academic Press, Boston, 1990.
- Min Tao and Xiaoming Yuan. Recovering low-rank and sparse components of matrices from incomplete and noisy observations. *SIAM J. Optim.*, 21(1):57–81, 2011.
- Joel A. Tropp. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015.
- Charles F Van Loan and Gene H Golub. *Matrix Computations*. Johns Hopkins University Press, Baltimore, MD, 4th edition, 2012.
- Huan Xu, Constantine Caramanis, and Sujay Sanghavi. Robust PCA via outlier pursuit. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2496–2504, Vancouver, Canada, 2010.
- Xinyang Yi, Dohyung Park, Yudong Chen, and Constantine Caramanis. Fast algorithms for robust PCA via gradient descent. In *Advances in Neural Information Processing Systems (NIPS)*, pages 4152–4160, Barcelona, Spain, 2016.
- Wen-Jun Zeng and Hing-Cheung So. Outlier-robust matrix completion via p -minimization. *IEEE Trans. Signal Processing*, 66(5):1125–1140, 2018.

Supplementary Materials

6 Preliminary lemmas

The following lemma gives some fundamental results for $\sin \Theta(U, V)$, which can be easily verified via definition.

Lemma 1. *Let $[U, U_c]$ and $[V, V_c]$ be two orthogonal matrices with $U, V \in \mathbb{R}^{n \times k}$. Then*

$$\|\sin \Theta(U, V)\|_{\text{ui}} = \|U_c^T V\|_{\text{ui}} = \|U^T V_c\|_{\text{ui}}.$$

Here $\|\cdot\|_{\text{ui}}$ denotes any unitarily invariant norm, including the spectral norm and Frobenius norm. In particular, for the spectral norm, it holds $\|\sin \Theta(U, V)\| = \|UU^T - VV^T\|$; for the Frobenius norm, it holds $\|\sin \Theta(U, V)\|_F = \frac{1}{\sqrt{2}}\|UU^T - VV^T\|_F$.

The following lemma is the well-known Weyl theorem, which gives the perturbation bound for eigenvalues of Hermitian matrix.

Lemma 2. *(Stewart and Sun, 1990, p.203) For two Hermitian matrices $A, \tilde{A} \in \mathbb{C}^{n \times n}$, let $\lambda_1 \leq \dots \leq \lambda_n$, $\tilde{\lambda}_1 \leq \dots \leq \tilde{\lambda}_n$ be eigenvalues of A, \tilde{A} , respectively. Then*

$$|\lambda_j - \tilde{\lambda}_j| \leq \|A - \tilde{A}\|, \quad \text{for } 1 \leq j \leq n.$$

The following lemma is used to establish the perturbation bound for the invariant subspace of a Hermitian matrix, which is due to Davis and Kahan.

Lemma 3. *(Davis and Kahan, 1970, Theorem 5.1) Let H and M be two Hermitian matrices, and let S be a matrix of a compatible size as determined by the Sylvester equation*

$$HY - YM = S.$$

If either all eigenvalues of H are contained in a closed interval that contains no eigenvalue of M or vice versa, then the Sylvester equation has a unique solution Y , and moreover

$$\|Y\|_{\text{ui}} \leq \frac{1}{\delta} \|S\|_{\text{ui}},$$

where $\delta = \min |\lambda - \omega|$ over all eigenvalues ω of M and all eigenvalues λ of H .

For a rectangular matrix $A \in \mathbb{R}^{m \times n}$ (without loss of generality, assume $m \geq n$), let the SVD of A be $A = U\Sigma V^T$, where $U = [U_1 | U_2 | U_3] = [u_1, \dots, u_k | u_{k+1}, \dots, u_r | u_{r+1}, \dots, u_m] \in \mathbb{R}^{m \times m}$, $V = [V_1 | V_2 | V_3] = [v_1, \dots, v_k | v_{k+1}, \dots, v_r | v_{r+1}, \dots, v_n] \in \mathbb{R}^{n \times n}$ are orthogonal matrices, and $\Sigma = \begin{bmatrix} \text{diag}(\Sigma_1, \Sigma_2) & 0_{r \times (n-r)} \\ 0_{(m-r) \times r} & 0_{(m-r) \times (n-r)} \end{bmatrix}$, $\Sigma_1 = \text{diag}(\sigma_1, \dots, \sigma_k)$, $\Sigma_2 = \text{diag}(\sigma_{k+1}, \dots, \sigma_r)$ with $\sigma_1 \geq \dots \geq \sigma_r > 0$, $k \leq r = \text{rank}(A)$. Then the spectral decomposition of $\begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix}$ can be given by

$$\begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix} = X \text{diag}(\Sigma_1, \Sigma_2, -\Sigma_1, -\Sigma_2, 0_{n-r}, 0_{m-r}) X^T, \quad (13)$$

where $X = \frac{1}{\sqrt{2}} \begin{bmatrix} U & -U & 0 & \sqrt{2}U_3 \\ V & V & \sqrt{2}V_3 & 0 \end{bmatrix}$ is an orthogonal matrix.

With the help of (13) and Lemmas 2 and 3, we are able to prove Lemma 4, which established an error bound for singular vectors.

Lemma 4. *Given $A \in \mathbb{R}^{m \times n}$ ($m \geq n$), let the SVD of A be given as above. Let $\hat{\sigma}_j, \hat{u}_j, \hat{v}_j$ be respectively the approximate singular values, right and left singular vectors of A satisfying that $\hat{U} = [\hat{u}_1, \dots, \hat{u}_k] \in \mathbb{R}^{m \times k}$ and $\hat{V} = [\hat{v}_1, \dots, \hat{v}_k] \in \mathbb{R}^{n \times k}$ are both orthonormal, $\hat{\Sigma} = \hat{U}^T A \hat{V} = \text{diag}(\hat{\sigma}_1, \dots, \hat{\sigma}_k)$ with $\hat{\sigma}_1 \geq \dots \geq \hat{\sigma}_k > 0$. Let*

$$E = A\hat{V} - \hat{U}\hat{\Sigma}, \quad F = A^T\hat{U} - \hat{V}\hat{\Sigma}. \quad (14)$$

If

$$\|(I_m - \widehat{U}\widehat{U}^T)A(I_n - \widehat{V}\widehat{V}^T)\| < \hat{\sigma}_k, \quad \max\{\|E\|, \|F\|\} < \sigma_k - \sigma_{k+1},$$

then

$$\max\{\Theta_u, \Theta_v\} \leq \eta, \quad \frac{\|U_1 \Sigma_1 V_1^T - \widehat{U}\widehat{\Sigma}\widehat{V}^T\|_{\max}}{\|A\|} \leq (\|U_1\|_{2,\infty}\Theta_v + \|V_1\|_{2,\infty}\Theta_u) + (1 + 3\|U_1\|_{2,\infty}\|V_1\|_{2,\infty})\Theta_u\Theta_v,$$

where $\Theta_u = \|\sin \Theta(U_1, \widehat{U})\|$, $\Theta_v = \|\sin \Theta(V_1, \widehat{V})\|$, $\eta = \frac{\max\{\|E\|, \|F\|\}}{\sigma_k - \sigma_{k+1} - \max\{\|E\|, \|F\|\}}$.

Proof. Let

$$\begin{aligned} H &= \begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix}, & X_1 &= \frac{1}{\sqrt{2}} \begin{bmatrix} U_1 & -U_1 \\ V_1 & V_1 \end{bmatrix}, \\ \widehat{X}_1 &= \frac{1}{\sqrt{2}} \begin{bmatrix} \widehat{U} & -\widehat{U} \\ \widehat{V} & \widehat{V} \end{bmatrix}, & X_2 &= \frac{1}{\sqrt{2}} \begin{bmatrix} U_2 & -U_2 & 0 & \sqrt{2}U_3 \\ V_2 & V_2 & \sqrt{2}V_3 & 0 \end{bmatrix}. \end{aligned}$$

By calculations, we have

$$\|X_1 X_1^T - \widehat{X}_1 \widehat{X}_1^T\|_{\text{ui}} = \|\text{diag}(U_1 U_1^T - \widehat{U}\widehat{U}^T, V_1 V_1^T - \widehat{V}\widehat{V}^T)\|_{\text{ui}} \quad (15)$$

By simple calculations, we have

$$H\widehat{X}_1 - \widehat{X}_1 \text{diag}(\widehat{\Sigma}, -\widehat{\Sigma}) = \frac{1}{\sqrt{2}} \begin{bmatrix} A\widehat{V} - \widehat{U}\widehat{\Sigma} & A\widehat{V} - \widehat{U}\widehat{\Sigma} \\ A^T\widehat{U} - \widehat{V}\widehat{\Sigma} & -A^T\widehat{U} + \widehat{V}\widehat{\Sigma} \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} E & E \\ F & -F \end{bmatrix} \triangleq R, \quad (16a)$$

$$HX_2 - X_2 \text{diag}(\Sigma_2, -\Sigma_2, 0, 0) = 0, \quad (16b)$$

where (16a) uses (14), (16b) uses the SVD of A . Then it follows from (16a) that

$$\|R\| = \left\| \text{diag}(E, F) \frac{1}{\sqrt{2}} \begin{bmatrix} I_k & I_k \\ I_k & -I_k \end{bmatrix} \right\| = \|\text{diag}(E, F)\| = \max\{\|E\|, \|F\|\}. \quad (17)$$

Pre-multiplying (16a) by X_2^T and using (16b), we have

$$X_2^T R = X_2^T H\widehat{X}_1 - X_2^T \widehat{X}_1 \text{diag}(\widehat{\Sigma}, -\widehat{\Sigma}) = \text{diag}(\Sigma_2, -\Sigma_2, 0, 0) X_2^T \widehat{X}_1 - X_2^T \widehat{X}_1 \text{diag}(\widehat{\Sigma}, -\widehat{\Sigma}). \quad (18)$$

To apply Lemma 3 to (18), we need to estimate the gap between the eigenvalues of $\text{diag}(\widehat{\Sigma}, -\widehat{\Sigma})$ and those of $\text{diag}(\Sigma_2, -\Sigma_2, 0, 0)$. Using (16a) and $\widehat{U}^T A \widehat{V} = \widehat{\Sigma}$, we have

$$(H - R\widehat{X}_1^T - \widehat{X}_1 R^T)\widehat{X}_1 = H\widehat{X}_1 - R = \widehat{X}_1 \text{diag}(\widehat{\Sigma}, -\widehat{\Sigma}), \quad (19)$$

which implies that $\pm\hat{\sigma}_j$ are eigenvalues of $H - R\widehat{X}_1^T - \widehat{X}_1 R^T$, and the corresponding eigenvectors are $\frac{1}{\sqrt{2}} \begin{bmatrix} \pm\hat{u}_j \\ \hat{v}_j \end{bmatrix}$, for $j = 1, \dots, k$. Next, we declare that $\hat{\sigma}_1, \dots, \hat{\sigma}_k$ are the k largest eigenvalues of $H - R\widehat{X}_1^T - \widehat{X}_1 R^T$. This is because

$$\begin{aligned} & \max_{\widehat{X}_1^T x=0} \frac{x^T (H - R\widehat{X}_1^T - \widehat{X}_1 R^T) x}{x^T x} \\ & \leq \|(I - \widehat{X}_1 \widehat{X}_1^T)(H - R\widehat{X}_1^T - \widehat{X}_1 R^T)(I - \widehat{X}_1 \widehat{X}_1^T)\| \\ & = \|(I - \widehat{X}_1 \widehat{X}_1^T)H(I - \widehat{X}_1 \widehat{X}_1^T)\| \\ & = \left\| \begin{bmatrix} I_m - \widehat{U}\widehat{U}^T & 0 \\ 0 & I_n - \widehat{V}\widehat{V}^T \end{bmatrix} \begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix} \begin{bmatrix} I_m - \widehat{U}\widehat{U}^T & 0 \\ 0 & I_n - \widehat{V}\widehat{V}^T \end{bmatrix} \right\| \\ & = \|(I_m - \widehat{U}\widehat{U}^T)A(I_n - \widehat{V}\widehat{V}^T)\| < \hat{\sigma}_k. \end{aligned}$$

Therefore, by Lemma 2, we have

$$|\sigma_j - \hat{\sigma}_j| \leq \|R\hat{X}_1^T + \hat{X}_1R^T\|, \quad \text{for } j = 1, \dots, k. \quad (20)$$

Together with (17), we get

$$\begin{aligned} |\sigma_j - \hat{\sigma}_j| &\leq \|R\hat{X}_1^T + \hat{X}_1R^T\| = \max_j |\lambda_j([R, \hat{X}_1] \begin{bmatrix} \hat{X}_1^T \\ R^T \end{bmatrix})| = \max_j |\lambda_j(\begin{bmatrix} \hat{X}_1^T \\ R^T \end{bmatrix} [R, \hat{X}_1])| \\ &= \max_j |\lambda_j(\begin{bmatrix} 0 & I_k \\ R^T R & 0 \end{bmatrix})| = \|R\| = \max\{\|E\|, \|F\|\}. \end{aligned} \quad (21)$$

Here we uses the property that for any two matrix $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times m}$, the nonzero eigenvalues of AB and BA are the same.

Now by the assumption that $\max\{\|E\|, \|F\|\} < \sigma_k - \sigma_{k+1}$, we have

$$\hat{\sigma}_k - \sigma_{k+1} = \sigma_k - \sigma_{k+1} + \hat{\sigma}_k - \sigma_k \geq \sigma_k - \sigma_{k+1} - \max\{\|E\|, \|F\|\} > 0, \quad (22)$$

therefore, the eigenvalues of $\text{diag}(\Sigma_2, -\Sigma_2, 0, 0)$ lie in $[-\sigma_{k+1}, \sigma_{k+1}]$, which has no eigenvalues of $\text{diag}(\hat{\Sigma}, -\hat{\Sigma})$. We are able to apply Lemma 3 to (18), which yields

$$\|X_2^T \hat{X}_1\|_{\text{ui}} \leq \frac{\|X_2^T R\|_{\text{ui}}}{\sigma_k - \sigma_{k+1} - \max\{\|E\|, \|F\|\}}. \quad (23)$$

Using (15), Lemma 1, (22) and (23), we get

$$\max\{\Theta_u, \Theta_v\} = \|\sin \Theta(X_1, \hat{X}_1)\| = \|X_2^T \hat{X}_1\| \leq \frac{\|X_2^T R\|}{\sigma_k - \sigma_{k+1} - \max\{\|E\|, \|F\|\}} \leq \eta. \quad (24)$$

Let

$$\hat{U} = U\Gamma_u = [U_1, U_2, U_3] \begin{bmatrix} \Gamma_{u1} \\ \Gamma_{u2} \\ \Gamma_{u3} \end{bmatrix}, \quad \hat{V} = V\Gamma_v = [V_1, V_2, V_3] \begin{bmatrix} \Gamma_{v1} \\ \Gamma_{v2} \\ \Gamma_{v3} \end{bmatrix}, \quad (25)$$

where $\Gamma_{u1} \in \mathbb{R}^{k \times k}$, $\Gamma_{u2} \in \mathbb{R}^{(r-k) \times k}$, $\Gamma_{u3} \in \mathbb{R}^{(m-r) \times k}$, $\Gamma_{v1} \in \mathbb{R}^{k \times k}$, $\Gamma_{v2} \in \mathbb{R}^{(r-k) \times k}$, $\Gamma_{v3} \in \mathbb{R}^{(n-r) \times k}$, and $\begin{bmatrix} \Gamma_{u1} \\ \Gamma_{u2} \\ \Gamma_{u3} \end{bmatrix}$, $\begin{bmatrix} \Gamma_{v1} \\ \Gamma_{v2} \\ \Gamma_{v3} \end{bmatrix}$ are both orthonormal. By (24), we have

$$\|\begin{bmatrix} \Gamma_{u2} \\ \Gamma_{u3} \end{bmatrix}\| = \Theta_u, \quad \sigma_{\min}(\Gamma_{u1}) = \sqrt{1 - \Theta_u^2}, \quad \|\begin{bmatrix} \Gamma_{v2} \\ \Gamma_{v3} \end{bmatrix}\| = \Theta_v, \quad \sigma_{\min}(\Gamma_{v1}) = \sqrt{1 - \Theta_v^2}. \quad (26)$$

Substituting (25) into $\hat{U}^T A \hat{V} = \hat{\Sigma}$ and using the SVD of A , we have

$$\hat{\Sigma} = [\Gamma_{u1}^T, \Gamma_{u2}^T, \Gamma_{u3}^T] \text{diag}(\Sigma_1, \Sigma_2, 0_{(m-r) \times (n-r)}) \begin{bmatrix} \Gamma_{v1} \\ \Gamma_{v2} \\ \Gamma_{v3} \end{bmatrix} = \Gamma_{u1}^T \Sigma_1 \Gamma_{v1} + \Gamma_{u2}^T \Sigma_2 \Gamma_{v2}. \quad (27)$$

Then it follows that

$$\begin{aligned} \|\Sigma_1 - \Gamma_{u1} \hat{\Sigma} \Gamma_{v1}^T\| &= \|(\Sigma_1 - \Gamma_{u1} \Gamma_{u1}^T \Sigma_1) + (\Gamma_{u1} \Gamma_{u1}^T \Sigma_1 - \Gamma_{u1} \Gamma_{u1}^T \Sigma_1 \Gamma_{v1} \Gamma_{v1}^T) - \Gamma_{u1} \Gamma_{u2}^T \Sigma_2 \Gamma_{v2} \Gamma_{v1}^T\| \\ &\leq \|I - \Gamma_{u1} \Gamma_{u1}^T\| \|\Sigma_1\| + \|\Gamma_{u1} \Gamma_{u1}^T\| \|I - \Gamma_{v1} \Gamma_{v1}^T\| \|\Sigma_1\| + \|\Gamma_{u2}\| \|\Gamma_{v2}\| \|\Sigma_2\| \\ &\leq (\Theta_u^2 + \Theta_v^2 + \Theta_u \Theta_v) \|\Sigma_1\|. \end{aligned} \quad (28)$$

Finally, using (26), (27), (28) and $\|\Gamma_{u1}\| \leq 1$, $\|\Gamma_{v1}\| \leq 1$, $\|\widehat{\Sigma}\| \leq \|A\|$, we have

$$\begin{aligned}
 \|U_1 \Sigma_1 V_1^T - \widehat{U} \widehat{\Sigma} \widehat{V}^T\|_{\max} &= \max_{i,j} |e_i^T (U_1 \Sigma_1 V_1^T - \widehat{U} \widehat{\Sigma} \widehat{V}^T) e_j| \\
 &= \max_{i,j} |e_i^T (U_1 \Sigma_1 V_1^T - U \Gamma_u \widehat{\Sigma} \Gamma_v^T V^T) e_j| \\
 &\leq \max_{i,j} |e_i^T (U_1 \Sigma_1 V_1^T - U_1 \Gamma_{u1} \widehat{\Sigma} \Gamma_{v1}^T V_1^T) e_j| + \|[U_2, U_3] \begin{bmatrix} \Gamma_{u3}^{u2} \\ \Gamma_{u3} \end{bmatrix} \widehat{\Sigma} \begin{bmatrix} \Gamma_{v3}^{v2} \\ \Gamma_{v3} \end{bmatrix}^T [V_2, V_3]^T\| \\
 &\quad + \max_{i,j} \left(|e_i^T [U_2, U_3] \begin{bmatrix} \Gamma_{u3}^{u2} \\ \Gamma_{u3} \end{bmatrix} \widehat{\Sigma} \Gamma_{v1}^T V_1^T e_j| + |e_i^T U_1 \Gamma_{u1} \widehat{\Sigma} \begin{bmatrix} \Gamma_{v3}^{v2} \\ \Gamma_{v3} \end{bmatrix}^T [V_2, V_3]^T e_j| \right) \\
 &\leq \max_{i,j} (3 \|e_i^T U_1\| \|e_j^T V_1\| \|A\| \Theta_u \Theta_v + \|A\| \Theta_u \Theta_v + \|e_j^T V_1\| \|A\| \Theta_u + \|e_i^T U_1\| \|A\| \Theta_v) \\
 &\leq \|A\| ((\|U_1\|_{2,\infty} \Theta_v + \|V_1\|_{2,\infty} \Theta_u) + (1 + 3\|U_1\|_{2,\infty} \|V_1\|_{2,\infty}) \Theta_u \Theta_v),
 \end{aligned}$$

completing the proof. \square

Lemma 5. (Tropp, 2015, Corollary 6.1.2) Let $\mathbf{S}_1, \dots, \mathbf{S}_n$ be independent random matrices with common dimension $d_1 \times d_2$, and assume that each matrix has uniformly bounded deviation from its mean:

$$\|\mathbf{S}_k - \mathbb{E}(\mathbf{S}_k)\| \leq L, \quad \text{for each } k = 1, \dots, n.$$

Let $\mathbf{Z} = \sum_{k=1}^n \mathbf{S}_k$, $v(\mathbf{Z})$ denote the matrix covariance statistic of the sum:

$$\begin{aligned}
 v(\mathbf{Z}) &= \max\{\|\mathbb{E}[(\mathbf{Z} - \mathbb{E}(\mathbf{Z}))(\mathbf{Z} - \mathbb{E}(\mathbf{Z}))^H]\|, \|\mathbb{E}[(\mathbf{Z} - \mathbb{E}(\mathbf{Z}))^H(\mathbf{Z} - \mathbb{E}(\mathbf{Z}))]\|\} \\
 &= \max\{\|\mathbb{E}\left[\sum_{k=1}^n (\mathbf{S}_k - \mathbb{E}(\mathbf{S}_k))(\mathbf{S}_k - \mathbb{E}(\mathbf{S}_k))^H\right]\|, \|\mathbb{E}\left[\sum_{k=1}^n (\mathbf{S}_k - \mathbb{E}(\mathbf{S}_k))^H(\mathbf{S}_k - \mathbb{E}(\mathbf{S}_k))\right]\|\}.
 \end{aligned}$$

Then for all $t \geq 0$,

$$\mathbb{P}\{\|\mathbf{Z} - \mathbb{E}(\mathbf{Z})\| \geq t\} \leq (d_1 + d_2) \cdot \exp\left(\frac{-t^2/2}{v(\mathbf{Z}) + Lt/3}\right).$$

Lemma 6. For any linear homogeneous function $F: \mathbb{R}^k \rightarrow \mathbb{R}^{m \times n}$, assume that the linear system of equations $F(x) = C$ either has a unique solution or has no solution at all. Then it holds

$$\operatorname{argmin}_x \|F(x) - C\| = \operatorname{argmin}_x \|F(x) - C\|_F.$$

Proof. For any $A, B \in \mathbb{R}^{m \times n}$, define $\langle A, B \rangle = \operatorname{trace}(A^T B)$. It is easy to see that $\langle \cdot, \cdot \rangle$ is an inner product over $\mathbb{R}^{m \times n}$. Denote the range space of $F(\cdot)$ by \mathcal{F} , and its orthogonal complement space by \mathcal{F}^\perp . Write $C = C_{\text{LS}} + C$ such that $C_{\text{LS}} \in \mathcal{F}$, and $C \in \mathcal{F}^\perp$. Then the solutions to $\min \|F(x) - C\|$ and $\min \|F(x) - C\|_F$ are nothing but the solutions to $F(x) = C_{\text{LS}}$. Since $C_{\text{LS}} \in \mathcal{F}$, $F(x) = C_{\text{LS}}$ has at least a solution. By the assumption, the solution should be unique. The proof is completed. \square

Lemma 7. Let $L_* \in \mathbb{R}^{m \times n}$ with $m \geq n$, let the SVD of L_* be $L_* = U_* \Sigma_* V_*^T$, where $U_* \in \mathbb{R}^{m \times r}$, $V_* \in \mathbb{R}^{n \times r}$ are orthonormal, $\Sigma_* = \operatorname{diag}(\sigma_{1*}, \dots, \sigma_{r*})$ with $\sigma_{1*} \geq \dots \geq \sigma_{r*} > 0$. Let $G \in \mathbb{R}^{m \times n}$ be a perturbation to L_* , $X \in \mathbb{R}^{m \times r}$, $Y \in \mathbb{R}^{n \times r}$ have full column rank. Denote $\theta_x = \|\sin \Theta(U_*, X)\|$, $\theta_y = \|\sin \Theta(V_*, Y)\|$. Then

$$\min_{X,Y} \|L_* - G - XY^T\| \geq \sigma_{r*} \max\{\sqrt{1 - \theta_x^2} \theta_y, \sqrt{1 - \theta_y^2} \theta_x\} \sqrt{1 - \theta_x^2} \sqrt{1 - \theta_y^2} - \|G\|.$$

Proof. Let $U_{*,c}, V_{*,c}$ be such that $U = [U_*, U_{*,c}]$, $V = [V_*, V_{*,c}]$ are orthogonal. Let $\widehat{X} = U_* C_x + U_{*,c} S_x$, $\widehat{Y} = V_* C_y + V_{*,c} S_y$, where the columns of \widehat{X} , \widehat{Y} form the orthonormal basis for $\mathcal{R}(X)$ and $\mathcal{R}(Y)$, respectively, $C_x^T C_x + S_x^T S_x = I_r$, $C_y^T C_y + S_y^T S_y = I_r$. By Lemma 1, we know that $\|S_x\| = \theta_x$, $\|S_y\| = \theta_y$.

Noticing that

$$\begin{aligned}
 \min_{X,Y} \|L_* - XY^T\|^2 &= \min_D \|U^T L_* V - U^T \widehat{X} \widehat{Y}^T V\|^2 = \min_D \left\| \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix} - \begin{bmatrix} C_x \\ S_x \end{bmatrix} D [C_y^T, S_y^T] \right\|^2 \\
 &= \left\| \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix} - \begin{bmatrix} C_x \\ S_x \end{bmatrix} [C_x, S_x]^T \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} C_y \\ S_y \end{bmatrix} [C_y^T, S_y^T] \right\|^2,
 \end{aligned}$$

we have

$$\begin{aligned}
 \min_{X,Y} \|L_* - XY^T\|^2 &\geq \max \left\{ \left\| C_x [C_x, S_x]^T \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} C_y \\ S_y \end{bmatrix} S_y^T \right\|^2, \left\| S_x [C_x, S_x]^T \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} C_y \\ S_y \end{bmatrix} C_y^T \right\|^2 \right\} \\
 &\geq \max \{ \sigma_{\min}^2(C_x) \|S_y\|^2, \sigma_{\min}^2(C_y) \|S_x\|^2 \} \sigma_{\min}^2 \left([C_x, S_x]^T \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} C_y \\ S_y \end{bmatrix} \right) \\
 &\geq \max \{ (1 - \theta_x^2) \theta_y^2, (1 - \theta_y^2) \theta_x^2 \} (\sigma_{r*} \sqrt{1 - \theta_x^2} \sqrt{1 - \theta_y^2})^2
 \end{aligned}$$

Combining it with the fact that $\|L_* - G - XY^T\| \geq \|L_* - XY^T\| - \|G\|$ for any X, Y , we get the conclusion. \square

Lemma 8. *Let L_* , G be the same as in Lemma 7. Let $X = (L_* - G)Y$, where $Y \in \mathbb{R}^{n \times r}$ is orthonormal. Denote $\theta_x = \|\sin \Theta(U_*, X)\|$, $\theta_y = \|\sin \Theta(V_*, Y)\|$. If $\|G\| < \sigma_{r*} \sqrt{1 - \theta_y^2}$, then*

$$\sigma_r(X) \geq \sigma_{r*} \sqrt{1 - \theta_y^2} - \|G\|, \quad \theta_x \leq \frac{\|G\|}{\sigma_r \sqrt{1 - \theta_y^2} - \|G\|}.$$

Proof. By Lemma 2 and Lemma 1, we have

$$\begin{aligned}
 \sigma_r(X) &= \sigma_r((L_* - G)Y) \geq \sigma_r(L_* Y) - \|GY\| \geq \sigma_r(\Sigma_* V_*^T Y) - \|G\| \geq \sigma_{r*} \sigma_{\min}(V_*^T Y) - \|G\| \\
 &= \sigma_{r*} \sigma_{\min}^{\frac{1}{2}}(Y^T V_* V_*^T Y) - \|G\| \geq \sigma_{r*} \sigma_{\min}^{\frac{1}{2}}(I_r - Y^T(I - V_* V_*^T)Y) - \|G\| \\
 &= \sigma_{r*} \sqrt{1 - \|(I - V_* V_*^T)Y\|^2} - \|G\| = \sigma_{r*} \sqrt{1 - \theta_y^2} - \|G\| > 0.
 \end{aligned} \tag{29}$$

Therefore, X has full column rank. Denote $G_x = (X^T X)^{-\frac{1}{2}}$, $\hat{X} = X G_x$. Then \hat{X} and $X = AY$ can be rewritten as $\hat{X} = AY G_x$. Using Lemma 1 and (29), we have

$$\|\theta_x\| = \|U_{*,c}^T \hat{X}\| = \|U_{*,c}^T (L_* - G) Y G_x\| \leq \|G Y G_x\| \leq \|G\| \|G_x\| \leq \frac{\|G\|}{\sigma_r(X)} \leq \frac{\|G\|}{\sigma_{r*} \sqrt{1 - \theta_y^2} - \|G\|}.$$

The proof is completed. \square

Lemma 9. *Let $U, X \in \mathbb{R}^{m \times r}$ both have orthonormal columns. It holds $\|X\|_{2,\infty} \leq \|U\|_{2,\infty} + \|\sin \Theta(U, X)\|$.*

Proof. Let U_c be such that $[U, U_c]$ is an orthogonal matrix. We can write $X = U C_x + U_c S_x$, where $C_x^T C_x + S_x^T S_x = I_r$. By Lemma 1, we have $\|\sin \Theta(U, X)\| = \|U_c^T X\| = \|S_x\|$. Then for any $1 \leq i \leq m$, we have

$$\|e_i^T X\| = \|e_i^T U C_x + e_i^T U_c S_x\| \leq \|e_i^T U\| + \|S_x\|,$$

the conclusion follows. \square

Lemma 10. *(Jain and Netrapalli, 2015, Lemmas 8,10) Let $A \in \mathbb{R}^{m \times n}$ with $m \geq n$. Suppose Ω is obtained by sampling each entry of A with probability $p \in [\frac{1}{4m}, 0.5]$. Then w.p. $\geq 1 - 1/m^{10+\log \alpha}$,*

$$\left\| \frac{1}{p} \Pi_{\Omega}(A) - A \right\| \leq \frac{6\sqrt{\alpha m}}{\sqrt{p}} \|A\|_{\max}.$$

7 Proof for Main Theorems

7.1 Proof of Theorem 1

Proof of Theorem 1. First, it holds $\|(I - U_* U_*^T) M (I - V_* V_*^T)\| = \|(I - U_* U_*^T) S_* (I - V_* V_*^T)\|$. Then by assumption, we have $\|(I - U_* U_*^T) M (I - V_* V_*^T)\| < \sigma_{r*}$.

Second, we have

$$\begin{aligned}\|E\| &= \|MV_* - U_*\Sigma_*\| = \|L_*V_* - U_*\Sigma_* + S_*V_*\| = \|S_*V_*\|, \\ \|F\| &= \|M^T U_* - V_*\Sigma_*\| = \|L_*^T U_* - V_*\Sigma_* + S_*^T U_*\| = \|S_*^T U_*\|.\end{aligned}$$

It follows

$$\max\{\|E\|, \|F\|\} = \max\{\|S_*V_*\|, \|S_*^T U_*\|\} < \sigma_r - \sigma_{r+1}.$$

Then applying Lemma 4 gives the conclusion. \square

7.2 Proof of Theorem 2

Throughout the rest of this section, we follow the notations in Algorithm 1. Besides that, we will also adopt the following notations. Denote

$$r = \text{rank}(L_*), \quad \kappa_* = \kappa_2(L_*), \quad p' = p(1 - \varrho), \quad \Omega_t = \Omega / \text{supp}(S_t), \quad G_t = S_t - S_*. \quad (30)$$

The SVDs of L_* is given by

$$L_* = [U_*, U_{*,c}] \text{diag}(\Sigma_*, 0) [V_*, V_{*,c}]^T, \quad (31)$$

where $[U_*, U_{*,c}]$ and $[V_*, V_{*,c}]$ are orthogonal matrices $U_* \in \mathbb{R}^{m \times r}$ and $V_* \in \mathbb{R}^{n \times r}$, $\Sigma_* = \text{diag}(\sigma_{1*}, \dots, \sigma_{r*})$ with $\sigma_{1*} \geq \dots \geq \sigma_{r*} > 0$. Further denote

$$\theta_{x,t} = \|\sin \Theta(U_*, X_t)\|, \quad \theta_{y,t} = \|\sin \Theta(V_*, Y_t)\|. \quad (32)$$

Lemma 11. $\|S_t - S_*\|_{\max} \leq 2\|\Pi_{\Omega}(X_t \Sigma_t Y_t^T - L_*)\|_{\max}$ for $t = 0, 1, \dots$.

Proof. Denote $\Phi_* = \text{supp}(S_*)$, $\Phi_t = \text{supp}(S_t)$, it is obvious that $S_t - S_*$ is supported on $\Phi_t \cup \Phi_*$ and $\Phi_t \cup \Phi_* \subset \Omega$. Now we claim that

$$\|\Pi_{\Omega}(S_t - S_*)\|_{\max} \leq 2\|\Pi_{\Omega}(X_t \Sigma_t Y_t^T - L_*)\|_{\max}.$$

To show the claim, it suffices to consider the following two cases.

Case (1) For any $(i, j) \in \Phi_t$, it holds $(S_t)_{(i,j)} = (L_* + S_* - X_t \Sigma_t Y_t^T)_{(i,j)}$. Then it follows that

$$|(S_t - S_*)_{(i,j)}| = |(L_* - X_t \Sigma_t Y_t^T)_{(i,j)}| \leq \|\Pi_{\Omega}(X_t \Sigma_t Y_t^T - L_*)\|_{\max}.$$

Case (2) For any $(i, j) \in \Phi_* \setminus \Phi_t$, it holds $(S_t)_{(i,j)} = 0$. If $|(S_t - S_*)_{(i,j)}| = |(S_*)_{(i,j)}| > 2\|\Pi_{\Omega}(X_t \Sigma_t Y_t^T - L_*)\|_{\max}$, then

$$|(L_* + S_* - X_t \Sigma_t Y_t^T)_{(i,j)}| > \|\Pi_{\Omega}(L_* - X_t \Sigma_t Y_t^T)\|_{\max}.$$

Noticing that S_* only changes s entries of $\Pi_{\Omega}(L_* - X_t \Sigma_t Y_t^T)$, we know that the (i, j) entry of $|\Pi_{\Omega}(L_* + S_* - X_t \Sigma_t Y_t^T)|$ is larger than the $(s+1)$ st largest entry of $|\Pi_{\Omega}(L_* + S_* - X_t \Sigma_t Y_t^T)|$. This contradicts with $(i, j) \notin \Phi_t$. \square

Lemma 12. Assume (A1). Denote $r'_s = \frac{\|S_0 - S_*\|_F^2}{\|S_0 - S_*\|^2}$. Let S_0 be obtained as in Algorithm 1. It holds

$$\|S_0 - S_*\| \leq 2\sqrt{\frac{2\varrho p}{r'_s} \mu r} \|L_*\|.$$

Proof. First, for any i, j , we have $L_{ij} = e_i^T U_* \Sigma_* V_*^T e_j$. Using (A1), we have

$$|L_{ij}| \leq \|e_i^T U_*\| \|\Sigma_*\| \|e_j^T V_*\| \leq \frac{\mu r}{\sqrt{mn}} \|L_*\|,$$

and hence

$$\|L_*\|_{\max} \leq \frac{\mu r}{\sqrt{mn}} \|L_*\|. \quad (33)$$

By Lemma 11, we have

$$\|S_0 - S_*\|_{\max} \leq 2\|\Pi_{\Omega}(L_*)\|_{\max} \leq \frac{2\mu r}{\sqrt{mn}} \|L_*\|. \quad (34)$$

Therefore, using (33), (34) and (A2), we have

$$\|S_0 - S_*\|_F \leq \sqrt{2s} \|S_0 - S_*\|_{\max} \leq 2\sqrt{2s} \|\Pi_{\Omega}(L_*)\|_{\max} \leq 2\sqrt{2s} \|L_*\|_{\max} \leq 2\sqrt{2\varrho p} \mu r \|L_*\|. \quad (35)$$

By the definition of r'_s , it follows that

$$\|S_0 - S_*\| \leq \frac{\|S_0 - S_*\|_F}{\sqrt{r'_s}} \leq 2\sqrt{\frac{2\varrho p}{r'_s}} \mu r \|L_*\|.$$

The proof is completed. \square

Proof of Theorem 2. By (A3), Lemma 10 and (33), w.p. $\geq 1 - 1/m^{10+\log \alpha}$, it holds

$$\|\frac{1}{p'} \Pi_{\Omega_0}(L_*) - L_*\| \leq \frac{6\sqrt{\alpha m}}{\sqrt{p'}} \|L_*\|_{\max} \leq \xi \mu r \|L_*\|. \quad (36)$$

Using Lemma 12 and (36), we have w.p. $\geq 1 - 1/m^{10+\log \alpha}$,

$$\|\frac{1}{p'} \Pi_{\Omega_0}(M - S_0) - L_*\| \leq \|\frac{1}{p'} \Pi_{\Omega_0}(L_*) - L_*\| + \frac{1}{p'} \|S_0 - S_*\| \leq (\xi + \gamma) \mu r \|L_*\|. \quad (37)$$

Let the SVD of $X_1^T L_* Y_1$ be $X_1^T L_* Y_1 = P \tilde{\Sigma} Q^T$, where P, Q are orthogonal matrices, $\tilde{\Sigma} = \text{diag}(\tilde{\sigma}_1, \dots, \tilde{\sigma}_r)$. Denote $\tilde{U} = X_1 P$, $\tilde{V} = Y_1 Q$, and let

$$E = L_* \tilde{V} - \tilde{U} \tilde{\Sigma}, \quad F = L_*^T \tilde{U} - \tilde{V} \tilde{\Sigma}. \quad (38)$$

Then it follows that

$$\|\Sigma_1 - X_1^T L_* Y_1\| = \|X_1^T [\frac{1}{p'} \Pi_{\Omega_0}(M - S_0) - L_*] Y_1\| \leq \|\frac{1}{p'} \Pi_{\Omega_0}(M - S_0) - L_*\|. \quad (39)$$

Using (37) and (39), by calculations, we get

$$\begin{aligned} \|E\| &= \|L_* \tilde{V} - \tilde{U} \tilde{\Sigma}\| = \|L_* Y_1 - X_1 P \tilde{\Sigma} Q^T\| = \|L_* Y_1 - X_1 X_1^T L_* Y_1\| \\ &\leq \|L_* Y_1 - X_1 \Sigma_1\| + \|X_1 (\Sigma_1 - X_1^T L_* Y_1)\| \\ &= \|L_* Y_1 - \frac{1}{p'} \Pi_{\Omega_0}(M - S_0) Y_1\| + \|\Sigma_1 - X_1^T L_* Y_1\| \\ &\leq 2\|\frac{1}{p'} \Pi_{\Omega_0}(M - S_0) - L_*\| \leq 2(\xi + \gamma) \mu r \|L_*\|, \quad \text{w.p. } \geq 1 - 1/m^{10+\log \alpha}. \end{aligned} \quad (40)$$

Similarly, we get

$$\|F\| \leq 2(\xi + \gamma) \mu r \|L_*\|, \quad \text{w.p. } \geq 1 - 1/m^{10+\log \alpha}. \quad (41)$$

Next, we only need to show $\max\{\|E\|, \|F\|\} \leq \sigma_{r^*}$ and $\|(I_m - \tilde{U} \tilde{U}^T) L_* (I_n - \tilde{V} \tilde{V}^T)\| < \tilde{\sigma}_r$. Once these two inequalities hold, we may apply Lemma 4.

For the first inequality, using (40), (41) and the assumption $(\xi + \gamma) \mu \kappa r < \frac{1}{6}$, we get

$$\max\{\|E\|, \|F\|\} \leq 2(\xi + \gamma) \mu r \|L_*\| < \sigma_{r^*}, \quad \text{w.p. } \geq 1 - 1/m^{10+\log \alpha}. \quad (42)$$

For the second inequality, using (35) and (36), we have

$$\left\| \frac{1}{p'} \Pi_{\Omega}(M - S_0) - L_* \right\| \leq \left\| \frac{1}{p'} \Pi_{\Omega}(L_*) - L_* \right\| + \frac{1}{p'} \|S_* - S_0\| \leq (\xi + \gamma) \mu r \|L_*\|. \quad (43)$$

Then using Lemma 2, (37), (39) and (43), we have

$$|\tilde{\sigma}_r - \sigma_{r*}| \leq |\tilde{\sigma}_r - \hat{\gamma}_{r,0}| + |\hat{\gamma}_{r,0} - \sigma_{r*}| \leq \|X_1^T L_* Y_1 - \Sigma_1\| + \left\| \frac{1}{p'} \Pi_{\Omega_0}(M - S_0) - L_* \right\| \leq 2(\xi + \gamma) \mu r \|L_*\|.$$

It follows that

$$\tilde{\sigma}_r \geq \sigma_{r*} - 2(\xi + \gamma) \mu r \|L_*\|. \quad (44)$$

Then using the assumption $(\xi + \gamma) \mu \kappa r < \frac{1}{6}$, (37) and (44), we have

$$\begin{aligned} \|(I_m - \tilde{U} \tilde{U}^T) L_* (I_n - \tilde{V} \tilde{V}^T)\| &= \|(I_m - X_1 X_1^T) [L_* - \frac{1}{p'} \Pi_{\Omega_0}(M - S_0)] (I_n - Y_1 Y_1^T)\| \\ &\leq \|L_* - \frac{1}{p'} \Pi_{\Omega_0}(M - S_0)\| \leq (\xi + \gamma) \mu r \|L_*\| < \sigma_{r*} - 2(\xi + \gamma) \mu r \|L_*\| \leq \tilde{\sigma}_r. \end{aligned}$$

Now using (40), (41), the assumption $(\xi + \gamma) \mu \kappa r < \frac{1}{6}$ and Lemma 4, we have

$$\max\{\theta_{x,1}, \theta_{y,1}\} = \max\{\|\sin \Theta(U_*, \tilde{U})\|, \|\sin \Theta(V_*, \tilde{V})\|\} \leq \frac{2(\xi + \gamma) \mu r \kappa}{1 - 1/3} = 3(\xi + \gamma) \mu r \kappa, \quad (45a)$$

$$\|L_* - \tilde{U} \tilde{\Sigma} \tilde{V}^T\|_{\max} / \|L_*\| \leq (\|U_*\|_{2,\infty} \theta_{y,1} + \|V_*\|_{2,\infty} \theta_{x,1}) + (1 + 3\|U_*\|_{2,\infty} \|V_*\|_{2,\infty}) \theta_{x,1} \theta_{y,1}. \quad (45b)$$

Using the assumption $(\xi + \gamma) \mu \kappa r < \frac{1}{3} \sqrt{\frac{\mu_1' r}{m}}$, by (45a), we have $\max\{\theta_{x,1}, \theta_{y,1}\} \leq \sqrt{\frac{\mu_1' r}{m}}$. On the other hand, assumption **(A1)** implies that

$$\|U_*\|_{2,\infty} \leq \sqrt{\frac{\mu r}{m}}, \quad \|V_*\|_{2,\infty} \leq \sqrt{\frac{\mu r}{n}}. \quad (46)$$

Then it follows from Lemma 9 that

$$\|X_1\|_{2,\infty} \leq \|U_*\|_{2,\infty} + \|\sin \Theta(X_1, U_*)\| \leq \sqrt{\frac{\mu r}{m}} + \sqrt{\frac{\mu_1' r}{m}} \leq \sqrt{\frac{\mu_1 r}{m}}, \quad (47a)$$

$$\|Y_1\|_{2,\infty} \leq \|V_*\|_{2,\infty} + \|\sin \Theta(Y_1, V_*)\| \leq \sqrt{\frac{\mu r}{n}} + \sqrt{\frac{\mu_1' r}{m}} \leq \sqrt{\frac{\mu_1 r}{n}}. \quad (47b)$$

Using the assumption $(\xi + \gamma) \mu \kappa r < \frac{1}{3} \sqrt{\frac{\mu_1' r}{m}}$, (37), (39), (45b), (46) and (47), by calculations, we have

$$\begin{aligned} \|L_* - X_1 \Sigma_1 Y_1^T\|_{\max} / \|L_*\| &\leq \|L_* - \tilde{U} \tilde{\Sigma} \tilde{V}^T\|_{\max} / \|L_*\| + \|\tilde{U} \tilde{\Sigma} \tilde{V}^T - X_1 \Sigma_1 Y_1^T\|_{\max} / \|L_*\| \\ &= \|L_* - \tilde{U} \tilde{\Sigma} \tilde{V}^T\|_{\max} / \|L_*\| + \|X_1 (X_1^T L_* Y_1 - \Sigma_1) Y_1^T\|_{\max} / \|L_*\| \\ &\leq \|L_* - \tilde{U} \tilde{\Sigma} \tilde{V}^T\|_{\max} / \|L_*\| + \|X_1\|_{2,\infty} \|X_1^T L_* Y_1 - \Sigma_1\| \|Y_1\|_{2,\infty} / \|L_*\| \\ &\leq \|L_* - \tilde{U} \tilde{\Sigma} \tilde{V}^T\|_{\max} / \|L_*\| + \|X_1\|_{2,\infty} \|Y_1\|_{2,\infty} (\xi + \gamma) \mu r \kappa, \\ &\leq (\|U_*\|_{2,\infty} \theta_{y,1} + \|V_*\|_{2,\infty} \theta_{x,1}) + (1 + 3\|U_*\|_{2,\infty} \|V_*\|_{2,\infty}) \theta_{x,1} \theta_{y,1} \\ &\quad + \|X_1\|_{2,\infty} \|Y_1\|_{2,\infty} \frac{1}{3} \sqrt{\frac{\mu_1' r}{m}} \\ &\leq \left(\sqrt{\frac{\mu r}{m}} \theta_{y,1} + \sqrt{\frac{\mu r}{n}} \theta_{x,1} + \theta_{x,1} \theta_{y,1} \right) + \mathcal{O}(n^{-3/2}), \end{aligned}$$

which completes the proof. \square

7.3 Proof of Theorem 3

Proof of Theorem 3. First, we give an upper bound for $\sup_{X \in \mathbb{R}^{m \times r}} \|\Pi_{\Omega_t}(R)\Pi_{\Omega_t}(XY_t^T)^T\|/\|X\|$. Let $\{\delta_{ij}\}$ be an independent family of BERNOLLI(p') random variables, $X^T = [x_1, \dots, x_m] \in \mathbb{R}^{r \times m}$ be arbitrary nonzero matrix with $\|X\| = 1$, and $Y_t^T = [y_1, \dots, y_n]$. Denote $E_{ij} = e_i e_j^T$, $R = [r_{ij}]$, $\mathbf{W}_{il} = \sum_{j,k} \delta_{ij} r_{ij} E_{ij} \delta_{lk} x_k^T y_l E_{kl}^T$, $\mathbf{Z} = \sum_{i,l} \mathbf{Z}_{i,l}$. By calculations, we have

$$\begin{aligned} \mathbb{E}(\mathbf{W}_{il}) &= p'^2 \sum_{j,k} r_{ij} E_{ij} y_k^T x_l E_{kl} = p'^2 \sum_j r_{ij} E_{ij} y_j^T x_l E_{jl} = p'^2 R_{(i,:)} Y_t x_l E_{il} = 0, \\ \|\mathbf{W}_{il}\| &\leq \sqrt{p'n} \max |r_{ij} x_j^T y_l| \leq \sqrt{\mu' r p'} \|R\|_{\max}, \\ \|\mathbb{E}[\sum_{i,l} \mathbf{W}_{il} \mathbf{W}_{il}^T]\| &= \|\mathbb{E}[\sum_{i,l} (\sum_{j,k} \delta_{ij} r_{ij} E_{ij} \delta_{lk} x_k^T y_l E_{kl}^T) (\sum_{j',k'} \delta_{ij'} r_{ij'} E_{ij'} \delta_{lk'} x_{k'}^T y_l E_{k'l}^T)^T]\| = 0, \\ \|\mathbb{E}[\sum_{i,l} \mathbf{W}_{il}^T \mathbf{W}_{il}]\| &= 0. \end{aligned}$$

By Lemma 5, we have $\mathbb{P}\{\|\mathbf{W}\| > t\} \leq (m+n) \exp\left(-\frac{3t/2}{\sqrt{\mu' r p'} \|R\|_{\max}}\right)$. Let $t = \frac{2}{3}(\log(m+n) + 5)\sqrt{\mu' r p'} \|R\|_{\max}$, then w.p. ≥ 0.99 , it holds

$$\|\mathbf{W}\| \leq \frac{2}{3}(\log(m+n) + 5)\sqrt{\mu' r p'} \|R\|_{\max}. \quad (48)$$

Second, It is easy to see that $X_{\text{opt}} = (M - S_t)Y_t$. By calculations, we have

$$\begin{aligned} &\min_X \|\Pi_{\Omega_t}(XY_t^T) - \Pi_{\Omega_t}(M - S_t)\|^2 \\ &= \min_{\Delta X} \|\Pi_{\Omega_t}((X_{\text{opt}} + \Delta X)Y_t^T) - \Pi_{\Omega_t}((M - S_t)Y_t Y_t^T + (M - S_t)(I - Y_t Y_t^T))\|^2 \\ &= \min_{\Delta X} \|\Pi_{\Omega_t}(\Delta X Y_t^T) - \Pi_{\Omega_t}(R)\|^2. \end{aligned} \quad (49)$$

$$\quad (50)$$

Then we declare that (50) is minimized when $\Delta X = \tilde{X}_{\text{opt}} - X_{\text{opt}}$. This is because (49) is minimized when $X = \tilde{X}_{\text{opt}}$ and $X = X_{\text{opt}} + \Delta X$. Thus, we have

$$\|\tilde{X}_{\text{opt}} - X_{\text{opt}}\| = \|\Delta X\| \leq \frac{\sup_{X \in \mathbb{R}^{m \times r}} \|\Pi_{\Omega_t}(R)\Pi_{\Omega_t}(XY_t^T)^T\|}{\sigma^2}. \quad (51)$$

Substituting (48) into (51), we get the conclusion. \square

7.4 Proof of Theorem 4

Lemma 13. Denote $r_s = \inf_t \frac{\|S_t - S_*\|_F^2}{\|S_t - S_*\|_F^2}$, $\zeta = \sqrt{\frac{2s\mu r}{mr_s}}$. If $\|L_* - X_t \Sigma_t Y_t^T\|_{\max} \leq c_t \|L_*\| \sqrt{\frac{\mu r}{m}}$ for some positive parameter c_t , then

$$\|S_t - S_*\| \leq 2c_t \|L_*\| \zeta, \quad |\gamma_{j,t} - \sigma_{j*}| \leq 2c_t \|L_*\| \zeta.$$

Proof. Using Lemma 11, by simple calculations, we have

$$\|S_t - S_*\| \leq \frac{\|S_t - S_*\|_F}{\sqrt{r_s}} \leq \sqrt{\frac{2s}{r_s}} \|S_t - S_*\|_{\max} \leq 2\sqrt{\frac{2s}{r_s}} \|\Pi_{\Omega}(L_* - X_t \Sigma_t Y_t^T)\|_{\max} \leq 2c_t \|L_*\| \sqrt{\frac{2s\mu r}{mr_s}} = 2c_t \|L_*\| \zeta.$$

Then by Lemma 2, we know that

$$|\gamma_{j,t} - \sigma_{j*}| \leq \|(M - S_t) - L_*\| = \|S_t - S_*\| \leq 2c_t \|L_*\| \zeta.$$

The proof is completed. \square

Proof of Theorem 4. First, denote $\bar{X}_{t+1} = (M - S_t)Y_t$, then we know that \bar{X}_{t+1} is the solution to $\min_X \|M - S_t - XY_t^T\|$. Also note that \tilde{X}_{t+1} on line 8 of Algorithm 1 is the solution to $\min_X \|\Pi_{\Omega_t}(M - S_t - XY_t^T)\|$. Then by Theorem 3, we have

$$\|\bar{X}_{t+1} - \tilde{X}_{t+1}\| \leq C_{\text{LS}}\|(M - S_t)(I - Y_t Y_t^T)\|_{\max}, \quad \text{w.p.} \geq 0.99.$$

Then it follows that from Lemma 1, Lemma 11 and Lemma 13 that

$$\begin{aligned} \|\bar{X}_{t+1} - \tilde{X}_{t+1}\| &\leq C_{\text{LS}}(\|L_*(I - Y_t Y_t^T)\|_{\max} + \|(S_t - S_*)(I - Y_t Y_t^T)\|_{\max}) \\ &\leq C_{\text{LS}}(\|L_*\|\sqrt{\frac{\mu r}{m}}\theta_{y,t} + \|S_t - S_*\|_{2,\infty}) \leq C_{\text{LS}}(\|L_*\|\sqrt{\frac{\mu r}{m}}\theta_{y,t} + \sqrt{2p\varrho n}\|S_t - S_*\|_{\max}) \\ &\leq C_{\text{LS}}(\|L_*\|\sqrt{\frac{\mu r}{m}}\theta_{y,t} + 2\sqrt{2p\varrho n}\|L_* - X_t \Sigma_t Y_t^T\|_{\max}) \leq \frac{C}{\sqrt{m}}\|L_*\|\theta_{y,t}. \end{aligned} \quad (52)$$

Second, using Lemma 13 and $4c\kappa\zeta < 1$, we have

$$\|S_t - S_*\| < 2c\theta_{y,t}\|L_*\|\zeta \leq \sqrt{2}c\|L_*\|\zeta < \frac{\sigma_{r^*}}{\sqrt{2}} \leq \sigma_{r^*}\sqrt{1 - \theta_{y,t}^2}, \quad (53)$$

Then by Lemma 8, we know that

$$\|\sin \Theta(\bar{X}_{t+1}, U_*)\| \leq \frac{\|S_t - S_*\|}{\sigma_{r^*}\sqrt{1 - \theta_{y,t}^2} - \|S_t - S_*\|}. \quad (54)$$

Using (54), the assumption $\|L_* - X_t \Sigma_t Y_t^T\|_{\max} \leq c\|L_*\|\theta_{y,t}\sqrt{\frac{\mu r}{m}}$, Lemma 13 and $\theta_{y,t} \leq \frac{1}{\sqrt{2}}$, we get

$$\|\sin \Theta(\bar{X}_{t+1}, U_*)\| \leq \frac{2c\|L_*\|\zeta\theta_{y,t}}{\frac{\sigma_{r^*}}{\sqrt{2}} - 2c\|L_*\|\zeta\theta_{y,t}} \leq \frac{2\sqrt{2}c\kappa\zeta\theta_{y,t}}{1 - 2c\kappa\zeta} < 4\sqrt{2}c\kappa\zeta\theta_{y,t}. \quad (55)$$

Therefore, using Lemma 1, (52) and (55), we have

$$\begin{aligned} \|\theta_{x,t+1}\| &= \|U_{*,c}^T X_{t+1}\| = \|U_{*,c}^T \tilde{X}_{t+1} R_{x,t+1}^{-1}\| \leq \|U_{*,c}^T \bar{X}_{t+1} R_{x,t+1}^{-1}\| + \|U_{*,c}^T (\tilde{X}_{t+1} - \bar{X}_{t+1}) R_{x,t+1}^{-1}\| \\ &\leq \|R_{x,t+1}^{-1}\|(\|\sin \Theta(\bar{X}_{t+1}, U_*)\| \|\bar{X}_{t+1}\| + \|\tilde{X}_{t+1} - \bar{X}_{t+1}\|) \\ &\leq \frac{1}{\sigma_r(\tilde{X}_{t+1})} \left(4\sqrt{2}c\kappa\zeta \|\bar{X}_{t+1}\| + \frac{C}{\sqrt{m}} \|L_*\| \right) \theta_{y,t}. \end{aligned} \quad (56)$$

Now using Lemma 2, (52), $\theta_{y,t} \leq \frac{1}{\sqrt{2}}$ and (53), we have

$$\begin{aligned} \|\bar{X}_{t+1}\| &= \|(M - S_t)Y_t\| \leq \|L_* Y\| + \|S_t - S_*\| \leq \|L_*\| + \sqrt{2}c\|L_*\|\zeta, \\ \sigma_r(\tilde{X}_{t+1}) &\geq \sigma_r(\bar{X}_{t+1}) - \frac{C}{\sqrt{m}}\|L_*\|\theta_{y,t} \geq \sigma_r((M - S_t)Y_t) - \frac{C}{\sqrt{2m}}\|L_*\| \geq \sigma_r(L_* Y_t) - \|S_t - S_*\| - \frac{C}{\sqrt{2m}}\|L_*\| \\ &\geq \sigma_{r^*}\sqrt{1 - \theta_{y,t}^2} - \sqrt{2}c\|L_*\|\zeta - \frac{C}{\sqrt{2m}}\|L_*\| \geq \frac{\sigma_{r^*}}{\sqrt{2}} - \sqrt{2}c\|L_*\|\zeta - \frac{C}{\sqrt{2m}}\|L_*\|. \end{aligned}$$

Substituting them into (56), we get the conclusion. \square

7.5 Proof of Theorem 5

Lemma 14. Follow the notations and assumptions in Lemma 1. Then

$$\|L_* - \hat{X}_{t+1} R_{x,t+1} Y_t^T\|_{\max} \leq \left((1 + C_{\text{LS}} \sqrt{\frac{\mu' r}{n}}) \sqrt{\frac{\mu r}{m}} + (1 + C_{\text{LS}} \sqrt{2p\varrho n}) \sqrt{\frac{\mu' r}{n}} 2c\zeta \right) \|L_*\| \theta_{y,t}.$$

Proof. Direct calculations give rise to

$$\begin{aligned} \|L_* - \tilde{X}_{t+1}Y_t^T\|_{\max} &\leq \|L_* - (M - S_t)Y_tY_t^T\|_{\max} + \|(M - S_t)Y_tY_t^T - \tilde{X}_{t+1}Y_t^T\|_{\max} \\ &\leq \|L_* - (M - S_t)Y_tY_t^T\|_{\max} + \|(M - S_t)Y_t - \tilde{X}_{t+1}\| \sqrt{\frac{\mu'r}{n}} \end{aligned} \quad (57a)$$

$$\leq \|L_* - (M - S_t)Y_tY_t^T\|_{\max} + C_{\text{LS}} \sqrt{\frac{\mu'r}{n}} \|(M - S_t)(I - Y_tY_t^T)\|_{\max} \quad (57b)$$

$$\begin{aligned} &\leq (1 + C_{\text{LS}} \sqrt{\frac{\mu'r}{n}}) \|L_*(I - Y_tY_t^T)\|_{\max} + (1 + C_{\text{LS}} \sqrt{2\varrho pn}) \sqrt{\frac{\mu'r}{n}} \|S_t - S_*\|_{\max} \\ &\leq \left((1 + C_{\text{LS}} \sqrt{\frac{\mu'r}{n}}) \sqrt{\frac{\mu r}{m}} + (1 + C_{\text{LS}} \sqrt{2\varrho pn}) \sqrt{\frac{\mu'r}{n}} 2c\zeta \right) \|L_*\| \theta_{y,t} \end{aligned} \quad (57c)$$

where (57a) uses $\|Y_t\|_{2,\infty} \leq \sqrt{\frac{\mu'r}{m}}$, (57b) uses Theorem 3, (57c) uses the SVD of L_* , $\|U_*\|_{2,\infty} \leq \sqrt{\frac{\mu r}{m}}$ and Lemme 1. \square

Proof of Theorem 5. First, by Lemma 7, we have

$$\|M - S_t - X_t \Sigma_t Y_t^T\| \geq \sigma_{r_*} \max\{\sqrt{1 - \theta_{x,t}^2} \theta_{y,t}, \sqrt{1 - \theta_{y,t}^2} \theta_{x,t}\} \sqrt{1 - \theta_{x,t}^2} \sqrt{1 - \theta_{y,t}^2} - \|S_t - S_*\|$$

Then using (53), $\theta_{x,t} \leq \frac{1}{\sqrt{2}}$ and $\theta_{y,t} \leq \frac{1}{\sqrt{2}}$, we get

$$\|M - S_t - X_t \Sigma_t Y_t^T\| \geq \frac{\sigma_{r_*}}{2\sqrt{2}} \theta_{y,t} - 2c \|L_*\| \sqrt{\frac{\mu r}{m}} \theta_{y,t}. \quad (58)$$

Second, by calculations, we have

$$\begin{aligned} \|M - S_t - \hat{X}_{t+1} \tilde{Y}_{t+1}^T\| &\leq \|(I - \hat{X}_{t+1} \hat{X}_{t+1}^T)(M - S_t)\| + \|\hat{X}_{t+1} \hat{X}_{t+1}^T (M - S_t) - \hat{X}_{t+1} \tilde{Y}_{t+1}^T\| \\ &\leq \|(I - \hat{X}_{t+1} \hat{X}_{t+1}^T) L_*\| + \|S_t - S_*\| + \|\hat{X}_{t+1}^T (M - S_t) - \tilde{Y}_{t+1}^T\| \\ &\leq \|L_*\| \theta_{x,t+1} + 2c \|L_*\| \zeta \sqrt{\frac{\mu r}{m}} \theta_{x,t+1} + \frac{C}{\sqrt{m}} \|L_*\| \theta_{x,t+1} \end{aligned} \quad (59)$$

$$\leq (1 + 2c\zeta \sqrt{\frac{\mu r}{m}} + \frac{C}{\sqrt{m}}) \phi \|L_*\| \theta_{y,t}, \quad (60)$$

where the first two terms of (59) use Lemma 1 and (53), respectively, and the last term can be obtained similar to (52), with the help of Lemma 14.

Then it follows that

$$\|M - S_{t+1} - X_{t+1} \Sigma_{t+1} Y_{t+1}^T\| \leq \|M - S_t - X_{t+1} \Sigma_{t+1} Y_{t+1}^T\| \quad (61a)$$

$$\leq (1 + 2c \sqrt{\frac{\mu r}{m}} + \frac{C}{\sqrt{m}}) \phi \|L_*\| \theta_{y,t} \quad (61b)$$

$$\begin{aligned} &\leq \frac{(1 + 2c\zeta \sqrt{\frac{\mu r}{m}} + \frac{C}{\sqrt{m}}) \phi \|L_*\|}{\frac{\sigma_{r_*}}{2\sqrt{2}} - 2c\zeta \|L_*\| \sqrt{\frac{\mu r}{m}}} \|M - S_t - X_t \Sigma_t Y_t^T\| \\ &= \psi \|M - S_t - X_t \Sigma_t Y_t^T\|, \end{aligned} \quad (61c)$$

where (61a) uses Lemma 6, (61b) uses (60), (61c) uses (58). The proof is completed. \square