

M-estimation in Low-rank Matrix Factorization: a General Framework

Peng Liu^{1,*} Wei Tu^{2,*} Jingyu Zhao³ Yi Liu² Linglong Kong^{2,†} Guodong Li³ Bei Jiang² Hengshuai Yao⁴
Guangjian Tian⁵

¹School of Mathematics, Statistics and Actuarial Science, University of Kent

²Department of Mathematical and Statistical Sciences, University of Alberta

³Department of Statistics and Actuarial Science, University of Hong Kong

⁴ Huawei Hi-Silicon, Canada

⁵Huawei Noah's Ark Lab, Hong Kong, China

Emails: p.liu@kent.ac.uk.

{wei.tu, yliu16, lkong, bei1}@ualberta. {gladys17, gdli}@hku.hk. {hengshuai.yao, tian.guangjian}@huawei.com

Abstract—Many problems in science and engineering can be reduced to the recovery of an unknown large matrix from a small number of random linear measurements. Matrix factorization arguably is the most popular approach for low-rank matrix recovery. Many methods have been proposed using different loss functions, such as the most widely used L_2 loss, more robust choices L_1 and Huber loss, and quantile and expectile loss for skewed data. All of them can be unified into the framework of M-estimation. In this paper, we present a general framework of low-rank matrix factorization based on M-estimation in statistics. The framework mainly involves two steps: we first apply Nesterov's smoothing technique to obtain an optimal smooth approximation for non-smooth loss functions, such as L_1 and quantile loss; secondly, we exploit an alternative updating scheme along with Nesterov's momentum method at each step to minimize the smoothed loss function. Strong theoretical convergence guarantee has been developed for the general framework, and extensive numerical experiments have been conducted to illustrate the performance of the proposed algorithm.

Index Terms—matrix recovery, M-estimation, matrix factorization, robustness, statistical foundation

I. INTRODUCTION

Motivation. In matrix recovery from linear measurements, we are interested in recovering an unknown matrix $X \in \mathbb{R}^{m \times n}$ from $p < mn$ linear measurements $b_i = \text{Tr}(A_i^T X)$, where each $A_i \in \mathbb{R}^{m \times n}$ is a measurement matrix, $i = 1, \dots, p$. Usually it is expensive or even impossible to fully sample the entire matrix X , and we are left with a highly incomplete set of observations. In general it is not always possible to recover X under such settings. However, if we impose a low-rank structure on X , it is possible to exploit this structure and efficiently estimate X . The problem of matrix recovery arises in a wide range of applications, such as collaborate filtering [1], image recovery [2], structure from motion, photometric stereo [3], system identification [4] and computer network tomography [5].

Matrix factorization arguably is the most popular and intuitive approach for low-rank matrix recovery. The basic idea is to decompose the low-rank matrix $X \in \mathbb{R}^{m \times n}$ into the product of two matrices

$$X = U^T V, \quad (1)$$

where $U \in \mathbb{R}^{r \times m}$ and $V \in \mathbb{R}^{r \times n}$. In many practical problems, the rank of a matrix is known or can be estimated in advance; see, for example, the rigid and nonrigid structures from motion as well as image recovery. The matrix U and V can also be interpreted as latent factors that drive the unknown matrix X .

Matrix factorization is usually based on L_2 loss, i.e. square loss, which is optimal for Gaussian errors. However, its performance may be severely deteriorated when the data is contaminated by outliers. For example, in a collaborate filtering system, some popular items have a lot of ratings regardless of whether they are useful, while others have fewer ones. There may even exist shilling attacks, i.e. a user may consistently give positive feedback to their products or negative feedback to their competitors regardless of the items themselves [6]. Recently there are a few attempts to address this problem; see, for example, [2], [7]. However, to the best of our knowledge, they focus on either L_1 or quantile loss, and are only useful in limited scenarios. For low-rank matrix recovery under M-estimation, He *et al.* [8] studied the use of a few smooth loss functions such as Huber and Welsch in this setting. In this paper, we propose a more general framework that is applicable to any M-estimation loss function, smooth or non-smooth, and provide theoretical convergence guarantee on the proposed state-of-the-art algorithm.

This paper introduces a general framework of M-estimation [9]–[11] on matrix factorization. Specifically, we consider loss functions related to M-estimation for matrix factorization. M-estimation is defined in a way similar to the well-known terminology “Maximum Likelihood Estimate (MLE)” in statistics, and has many good properties that are similar to those of MLE. Meanwhile, it still retains intuitive interpretation. The proposed class of loss functions includes well-known L_1 and

*These authors contributed equally to this work.

†Correspondence Author. Part of the work is done during a sabbatical at Huawei.

L_2 loss as special cases. In practice, we choose a suitable M-estimation procedure according to our knowledge of the data and the specific nature of the problem. For example, Huber loss enjoys the property of smoothness as L_2 loss and robustness as L_1 loss.

The loss functions of some M-estimation procedures are smooth, such as the L_2 loss, while some others are non-smooth such as L_1 and quantile loss. Note that the resulting objective functions all have a bilinear structure due to the decomposition at Eq. (1). For non-smooth cases, we first consider Nesterov's smoothing method to obtain an optimal smooth approximation [12], and the bilinear structure is preserved. The alternating minimization method is hence used to search for the solutions. At each step, we employ Nesterov's momentum method to accelerate the convergence, and it turns out to find the global optima easily. Figure 1 gives the flowchart of the proposed algorithm. Theoretical convergence analysis is conducted for both smooth and non-smooth loss functions.

Contributions. We summarize our contributions below.

1. We propose to do matrix factorization based on the loss functions of M-estimation procedures. The proposed framework is very general and applicable to any M-estimate loss function, which gives us flexibility in selecting a suitable loss for specific problems.
2. We propose to use Nesterov's smoothing technique to obtain an optimal smooth approximation when the loss function is non-smooth.
3. We consider the Nesterov's momentum method, rather than gradient descent methods, to perform the optimization at each step of the alternating minimization, which greatly accelerates the convergence.
4. We provide theoretical convergence guarantees for the proposed algorithm.
5. We illustrate the usefulness of our method by conducting extensive numerical experiments on both synthetic data and real data.

II. METHODOLOGY FRAMEWORK

Let $X^* \in \mathbb{R}^{m \times n}$ be the target low-rank matrix, and $A_i \in \mathbb{R}^{m \times n}$ with $1 \leq i \leq p$ be given measurement matrices. Here A_i 's can be the same or different with each other. We assume the observed signals $b = (b_1, \dots, b_p)^\top$ have the structure

$$b_i = \langle A_i, X^* \rangle + \epsilon_i, \quad i = 1, \dots, p, \quad (2)$$

where $\langle A_i, X \rangle := \text{Tr}(A_i^\top X)$ and ϵ_i is the error term. Suppose that the rank of matrix X^* is no more than r with $r \ll \min(m, n, p)$. We then have the decomposition $X^* = U^{*\top} V^*$, and the matrix can be recovered by solving a non-convex optimization problem

$$\min_{U \in \mathbb{R}^{r \times m}, V \in \mathbb{R}^{r \times n}} \frac{1}{p} \sum_{i=1}^p \mathcal{L}(b_i - \langle A_i, U^\top V \rangle), \quad (3)$$

where $\mathcal{L}(\cdot)$ is the loss function used an M-estimation procedure. For example, $\mathcal{L}(y) = y^2$ for the L_2 loss and $\mathcal{L}(y) = |y|$ for the L_1 loss. Here $\mathcal{L}(\cdot)$ usually is convex. Let $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow$

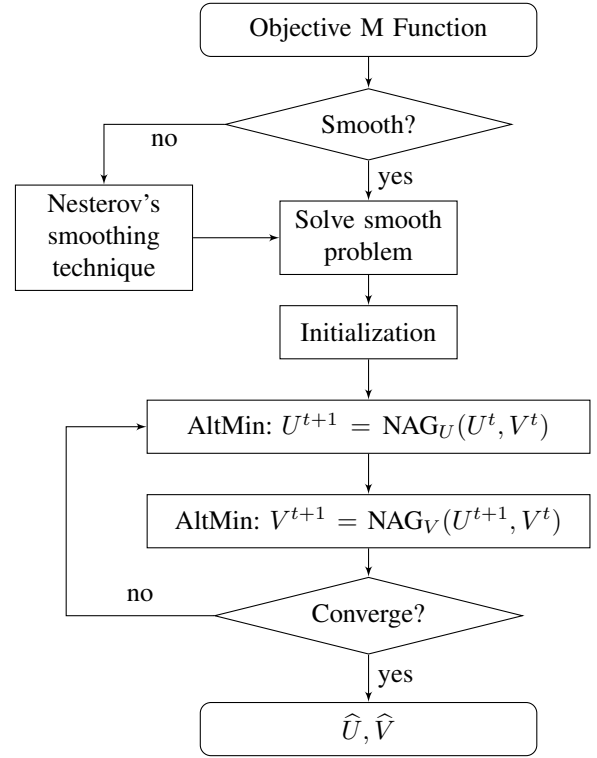


Figure 1: Flowchart of the whole algorithm

\mathbb{R}^p be an affine transformation with the i^{th} entry of $\mathcal{A}(X)$ being $\langle A_i, X \rangle$, and $\mathcal{M}(x) = p^{-1} \sum_{i=1}^p \mathcal{L}(x_i)$ for a vector $x = (x_1, \dots, x_p)^\top$. We can then rewrite (3) into a more compact form

$$\min_{U \in \mathbb{R}^{r \times m}, V \in \mathbb{R}^{r \times n}} \mathcal{M}(b - \mathcal{A}(U^\top V)). \quad (4)$$

A. The case that \mathcal{M} is not smooth

The loss function $\mathcal{L}(\cdot)$, and hence $\mathcal{M}(\cdot)$, may be non-smooth in some M-estimation procedures, such as the L_1 loss. For this non-smooth case, we first employ Nesterov's smoothing method to obtain an optimal smooth approximation; see [12] (pp. 129-132).

Specifically, we first assume that the objective function \mathcal{M} has the following structure

$$\mathcal{M}(b - \mathcal{A}(U^\top V)) = \hat{\mathcal{M}}(b - \mathcal{A}(U^\top V)) + \max_u \left\{ \langle B(b - \mathcal{A}(U^\top V)), u \rangle_2 - \hat{\phi}(u) \right\}, \quad (5)$$

where $\hat{\mathcal{M}}(\cdot)$ is continuous and convex; see Eq. (2.2) in [12].

Then the objective function \mathcal{M} can be approximated by

$$\mathcal{M}_\pi(b - \mathcal{A}(U^\top V)) = \hat{\mathcal{M}}(b - \mathcal{A}(U^\top V)) + \max_u \left\{ \langle B(b - \mathcal{A}(U^\top V)), u \rangle_2 - \hat{\phi}(u) - \pi d_2(u) \right\}, \quad (6)$$

where π is a positive smoothness parameter. Note that $\mathcal{M}_\pi(\cdot)$ is smooth and can be optimized by many available gradient-based methods.

B. Alternating Minimization

Due to the bilinear form with respect to U and V in the objective functions Eq. (4) and Eq. (6), we consider an alternating minimization scheme. Specifically, in each iteration, we keep one of U and V fixed, optimize over the other, and then switch in the next iteration until the algorithm converges. Details can be found in Figure 1.

Moreover, direct optimization of Eq. (4) may lead to unstable solutions, and this paper solves this problem by adopting the regularization method via Procrustes flow in [13]. Specifically, rather than Eq. (4), we attempt to do the optimization below

$$\min_{U \in \mathbb{R}^{r \times m}, V \in \mathbb{R}^{r \times n}} \mathcal{M}(b - \mathcal{A}(U^\top V)) + \lambda \|UU^\top - VV^\top\|_F^2, \quad (7)$$

where $\|\cdot\|_F$ stands for the Frobenius norm, and the ad hoc choice of λ is $1/16$. Denote by $\mathcal{M}^\lambda(U, V)$ the regularized version (7), and all algorithms designed below are for this objective function. For the case with non-smooth \mathcal{M} , we can similarly define the regularized objective function, denoted by $\mathcal{M}_\pi^\lambda(U, V)$, and all algorithms for $\mathcal{M}^\lambda(U, V)$ can then be applied.

III. INITIALIZATION

When the starting values of U and V are orthogonal (or almost orthogonal) to the true space, an alternating minimization algorithm may never converge to true values, and hence an initialization procedure is needed to avoid this situation. Here we adopt the singular value projection (SVP), which originates from [14] and was later used by [13] to provide the initial values of U and V for the designed algorithm in the next section. The main difference between our method and the original one is summarized into Line 2 of Algorithm 1, where we use a loss function related to M-estimation,

$$\nabla_X \mathcal{M}(b - \mathcal{A}(X^t)) = -\frac{1}{p} \sum_{i=1}^p \dot{\mathcal{L}}(b - \mathcal{A}(X^t)) A_i,$$

where $\dot{\mathcal{L}}(\cdot)$ is the first derivative of $\mathcal{L}(\cdot)$.

Algorithm 1: Initialization by SVP algorithm

Input: \mathcal{A} , b , tolerance ϵ_1 , step size ξ_t with $t = 0, 1, \dots$, and $X^0 = 0_{m \times n}$
Output: X^{t+1}

- 1 **Repeat**
- 2 $Y^{t+1} \leftarrow X^t - \xi_t \nabla_X \mathcal{M}(b - \mathcal{A}(X^t))$
- 3 Compute top r singular vectors of Y^{t+1} : U_r, Σ_r, V_r
- 4 $X^{t+1} \leftarrow U_r \Sigma_r V_r$
- 5 $t \leftarrow t + 1$
- 6 **Until** $\|X^{t+1} - X^t\|_F \leq \epsilon_1$

It is noteworthy to point out that Algorithm 1 can be directly used to recover the matrix X^* if it is iterated for sufficient times. However, the singular value calculation here is time-consuming when the dimension of matrix X^* is large. In this initialization step, we do not need a very small tolerance ϵ_1 ,

i.e., a rough output is sufficient. Our simulation experiments show that, after several iterations of the SVP algorithm, the resulting values are close to the true ones, while it is not the case for random initialization.

Algorithm 1 can be rewritten into a compact form

$$X^{t+1} \leftarrow \mathcal{P}_r (X^t - \xi_t \nabla_X \mathcal{M}(b - \mathcal{A}(X^t))),$$

where \mathcal{P}_r denotes the operation of projection onto the space of rank- r matrices. Moreover, the original objective function \mathcal{M} , rather than the regularized one at Eq. (7), is used in Algorithm 1, and for the non-smooth case we will use \mathcal{M}_π at Eq. (6).

IV. ALGORITHM

There are two layers of iterations in our algorithm: the outer layer is the alternating minimization procedure; and the inner layer employs Nesterov's momentum algorithm to obtain updated values of U^{t+1} or V^{t+1} .

Algorithm 2: Nesterov's accelerate gradient (NAG) method

Input: U^t, V^t , momentum parameter γ , learning rate η , and tolerance ϵ_2
Output: U^{t+1}

- 1 **Repeat**
- 2 $\nu_{(i)}^t = \gamma \nu_{(i-1)}^t + \eta \nabla_U \mathcal{M}^\lambda(U_{(i-1)}^t - \gamma \nu_{(i-1)}^t, V^t)$
- 3 $U_{(i)}^t = U_{(i-1)}^t + \nu_{(i)}^t$
- 4 **Until** $\|U_{(i)}^t - U_{(i-1)}^t\|_F \leq \epsilon_2$

We first introduce the inner layer, where the Nesterov's momentum method is used to update the values of U^t and V^t to those of U^{t+1} and V^{t+1} . Algorithm 2 gives the details of updating the value of U^t , and we denote it by NAG_U for simplicity. Here the gradient $\nabla_U \mathcal{M}^\lambda(U, V)$ is defined as the gradient with respect to U , and similarly we can define $\nabla_V \mathcal{M}^\lambda(U, V)$. Moreover, $\nu_{(i)}^t$ stands for the momentum term, γ is the momentum parameter, and η is the learning rate. The value of γ is usually chosen to be around 0.9; see [15] and [16]. Similarly, we can give the detailed algorithm for updating the value of V^t , and it can be denoted by NAG_V .

The alternating minimization method is employed for the outer layer of iterations; see Algorithm 3 for details. The final solutions can be denoted by \hat{U} and \hat{V} , and we then can use $\hat{U}^\top \hat{V}$ to approximate the low-rank matrix X^* .

Algorithm 3: Alternating Minimization

Input: U^0, V^0
Output: \hat{U}, \hat{V}

- 1 **Repeat**
- 2 Update U^t with $U^{t+1} = \text{NAG}_U(U^t, V^t)$
- 3 Update V^t with $V^{t+1} = \text{NAG}_V(U^{t+1}, V^t)$
- 4 **Until** converge

V. CONVERGENCE RESULTS

Suppose that U and V are a pair of solutions, i.e. $X = U^\top V$. It then holds that, for an orthonormal matrix R satisfying $R^\top R = I_r$, $U^\dagger = RU$ and $V^\dagger = RV$ are another pair of solutions. To evaluate the performance of the proposed algorithm, we first define a distance between two matrices

$$\text{dist}(U, U^\dagger) = \min_{R \in \mathbb{R}^{r \times r}: R^\top R = I_r} \|U - RU^\dagger\|_F,$$

where $U, U^\dagger \in \mathbb{R}^{r \times m}$ with $m \geq r$; see [13].

Definition 1. (*Restricted Isometry Property (RIP)*) A linear map \mathcal{A} satisfies the r -RIP with constant δ_r , if

$$(1 - \delta_r)\|X\|_F^2 \leq \|\mathcal{A}(X)\|_2^2 \leq (1 + \delta_r)\|X\|_F^2$$

is satisfied for all matrices $X \in \mathbb{R}^{m \times n}$ of rank at most r .

Theorem 1. Let $X \in \mathbb{R}^{m \times n}$ be a rank r matrix, with singular values $\sigma_1(X) \geq \sigma_2(X) \geq \dots \geq \sigma_r(X) > 0$ and condition number $\kappa = \sigma_1(X)/\sigma_r(X)$. Denote by $X = A^\top \Sigma B$ the corresponding SVD decomposition. Let $U = A^\top \Sigma^{1/2} \in \mathbb{R}^{m \times r}$ and $V = B^\top \Sigma^{1/2} \in \mathbb{R}^{n \times r}$. Assume that \mathcal{A} satisfies a rank- $6r$ RIP condition with RIP constant $\sigma_{6r} < \frac{1}{25}$, $\xi_t = \frac{1}{p}$. Then using $T_0 \geq 3 \log(\sqrt{r}\kappa) + 5$ iterations in Algorithm 1 yields a solution U_0, V_0 obeying

$$\text{dist}\left(\begin{bmatrix} U_0 \\ V_0 \end{bmatrix}, \begin{bmatrix} U \\ V \end{bmatrix}\right) \leq \frac{1}{4}\sigma_r(U). \quad (8)$$

Furthermore, starting from any initial solution obeying (8), the t^{th} iterate of Algorithm 3 satisfies that

$$\text{dist}\left(\begin{bmatrix} U_t \\ V_t \end{bmatrix}, \begin{bmatrix} U \\ V \end{bmatrix}\right) \leq \frac{1}{4}(1 - \tilde{\tau}_1)^t \frac{\tilde{\mu}}{\xi} \frac{1 + \delta_r}{1 - \delta_r} \sigma_r(U) \quad (9)$$

Eq. (8) in the above theorem guarantees that, under the RIP assumption on linear measurements \mathcal{A} , our algorithm can achieve a good initialization. Eq. (9) states that, starting from a sufficiently accurate initialization, the algorithm exhibits a linear convergence rate. Moreover, the specific convergence rates of the initialization at Eq. (8) and the alternating minimization at Eq. (9) both depend on the RIP constant δ_{6r} .

Jain *et al.* [14] and Tu *et al.* [13] built the convergence result for alternating minimization under least square matrix factorization, however, their proving methods cannot be directly adopted to derive the results at Eq. (8). This paper extends [14]’s SVP method to more general objective functions based on M-estimation, and a detailed proof of the convergence of SVP initialization is also provided. For the linear convergence at Eq. (9), Tu *et al.* [13] gave a similar result, while the gradient descent method was used for the alternating minimization. This paper adopts the proof technique in [17] to establish the linear convergence of the alternating minimization based on Nesterov’s momentum method.

When \mathcal{M} is not smooth, the regularized objective function \mathcal{M}^λ , defined in Eq. (7), is also not smooth, while the function \mathcal{M}_π^λ is smooth. Denote $\{U^\dagger, V^\dagger\} = \min_{U, V} \mathcal{M}^\lambda(U, V)$ and $\{U^{\pi\dagger}, V^{\pi\dagger}\} = \min_{U, V} \mathcal{M}_\pi^\lambda(U, V)$.

Theorem 2. (*Convergence of optimal solution of smoothed objective function*) As $\pi \rightarrow 0^+$, we have $U^{\pi\dagger\top} V^{\pi\dagger} \rightarrow U^{\dagger\top} V^\dagger$.

The above theorem guarantees that the optimal solution of the smoothed object function converges to that of the non-smooth one as the smoothness parameter π tends to zero.

Theorem 2 is one of our important contributions. In the literature, most of the state-of-the-art smoothing methods were given without theoretical justifications; see, for example, [18]. Theorem 2 implies that the optimal solution for the smooth approximation indeed converges to the solution of the non-smooth objective function, i.e. the smooth and non-smooth objective functions lead to the same results under some regularity conditions. Under matrix factorization settings, Yang *et al.* [19] considered Nesterov’s smoothing method to obtain a smooth approximation, while it handled nonnegative matrices only, which is a much simpler problem compared with the one in this paper. Moreover, no strong theoretical convergence guarantee was provided there.

VI. NUMERICAL EXPERIMENTS

As a robust alternative to L_2 loss, L_1 loss brings computation difficulties since it is not smooth. Using the Nesterov’s smoothing method, we obtained the well-known Huber loss:

$$\mathcal{L}_\mu(a) = \begin{cases} \frac{1}{2\mu}|a|^2 & \text{for } |a| \leq \mu \\ |a| - \frac{\mu}{2} & \text{otherwise} \end{cases},$$

where μ is the smoothness parameter, and $\mathcal{L}_\mu(a)$ approaches the L_1 loss as μ decreases. In this experiments section, we compare the performance of Huber and L_2 loss.

A. Synthetic data

We first generate a matrix X of size $m \times n$ by sampling each entry from the Gaussian distribution $\mathcal{N}(0, 1)$, and the true matrix X^* is obtained via truncated singular value decomposition (SVD) by keeping the first r largest singular values. The sampling matrices A_i , $i = 1, \dots, p$ are independently produced by sampling each entry from the Gaussian distribution $\mathcal{N}(0, 1)$. To evaluate the robustness of the proposed method, we consider eight distributions for the error term ϵ_i : (a) No error; (b) $\mathcal{N}(0, 2)$; (c) $\mathcal{N}(0, 10)$; (d) $\log\mathcal{N}(0, 1)$; (e) Cauchy; (f) $t(3)$; (g) Pareto(1, 1); (h) $0.9\mathcal{N}(0, 1) + 0.1\mathcal{N}(100, 1)$. Note that (e) and (g) are heavy-tailed with the first-order moment being infinite. We consider two loss functions for the recovery: Huber and L_2 loss, and there are 100 replications for each error distribution with $m = n = 100$, $r = 10$ and $p = 5000$.

We consider the following three metrics for evaluation: (1) relative error (RE): $\|X^* - \hat{U}^\top \hat{V}\|_F / \|X^*\|_F$; (2) recovery rate (RR): fraction of the elements whose element-wise relative error is smaller than 5%, where element-wise relative error is defined as $|(X_{ij}^* - (\hat{U}^\top \hat{V})_{ij}) / X_{ij}^*|$; (3) test error (MSE): a test set $\{A^*, b^*\}$ with $p^* = 100$ was generated using the same pipeline, and the test error is defined as $p^{*-1} \|A^*(\hat{U}^\top \hat{V}) - b^*\|^2$. For both RE and MSE, a smaller value is desired, while a value closer to one indicates better performance for RR. Table

I presents the results based on the average of 100 replications, and we highlight preferable figures by boldface.

When the data is very heavy-tailed (error (e) and (g)), the algorithm diverges for some cases. Table II shows the percentages of convergent repetitions. For the other six distributions, the algorithm converges for all replications. Moreover, we observe that scaling down the learning rate can make the algorithm converge, but this severely slows down the algorithm. When starting with a larger learning rate, the algorithm with L_2 loss tends to diverge, while that with the Huber loss converges faster.

Error	Loss	RE	RR	MSE
(a)	Huber	0.003	0.960	0.030
	L_2	0.003	0.969	0.022
(b)	Huber	0.030	0.650	2.885
	L_2	0.028	0.671	2.487
(c)	Huber	0.163	0.185	84.628
	L_2	0.141	0.211	63.040
(d)	Huber	0.028	0.667	2.617
	L_2	0.038	0.580	4.571
(e)	Huber	0.039	0.570	4.832
	L_2	1.083	0.037	4742
(f)	Huber	0.020	0.757	1.238
	L_2	0.024	0.711	1.862
(g)	Huber	0.068	0.399	14.68
	L_2	1.825	0.021	12451
(h)	Huber	0.025	0.704	1.920
	L_2	0.494	0.063	791.1

Table I: Simulation results for synthetic data with eight error distributions: (a) No error; (b) $\mathcal{N}(0, 2)$; (c) $\mathcal{N}(0, 10)$; (d) $\log\mathcal{N}(0, 1)$; (e) Cauchy; (f) $t(3)$; (g) Pareto(1, 1); (h) $0.9\mathcal{N}(0, 1) + 0.1\mathcal{N}(100, 1)$.

Error	Huber	L_2
(e)	99.4%	80.0%
(g)	100.0%	62.8%

Table II: Percentage of the replications where the algorithm converges for two error distributions: (e) Cauchy; (g) Pareto(1, 1).

When there is no error, Huber and L_2 loss are comparable, both reaching a recovery rate over 96%. In terms of the metrics presented, L_2 loss is better than Huber loss when the error is normally distributed. However, Huber loss is more robust when error introduces bias, skewness or outliers. Huber loss has a substantial advantage over L_2 loss especially when the error is skewed (error d) or heavy-tailed (error f). When the linear measurements are biased (error h) or contaminated by heavy-tailed outliers (errors e & g), matrix recovery using L_2 recovers only less than 7% of entries and has very large MSE, while employing Huber loss allows the procedure to recover at least 40% of the entries and the recovery error is less than 7% in general.

Figure 2 gives box plots of recovery rates for different error distributions and loss functions. It can be seen that recovery rates vary within a small scale, and the observations are consistent with those in Table I. Overall, the algorithm with Huber loss is comparable to that with L_2 loss in special cases (error a-c) and significantly better in general cases (error d-h).

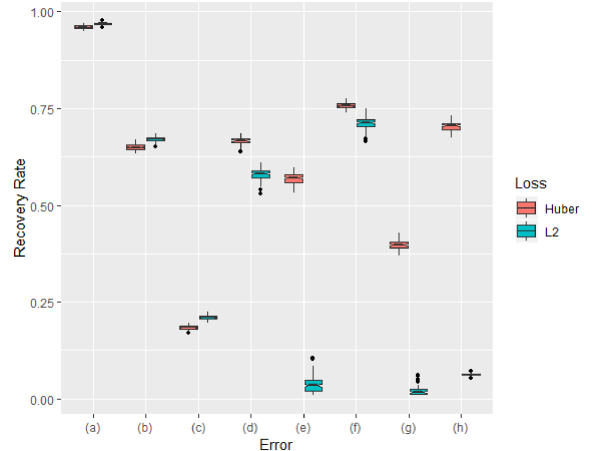


Figure 2: Box plots of recovery rates under Huber (left) and L_2 (right) loss functions with eight error distributions: (a) No error; (b) $\mathcal{N}(0, 2)$; (c) $\mathcal{N}(0, 10)$; (d) $\log\mathcal{N}(0, 1)$; (e) Cauchy; (f) $t(3)$; (g) Pareto(1, 1); (h) $0.9\mathcal{N}(0, 1) + 0.1\mathcal{N}(100, 1)$.

B. Real data

We further demonstrate the efficiency and robustness of our low-rank matrix recovery algorithm by applying it to two real examples.

The first is a chlorine concentration dataset, available in [20]. The dataset contains a matrix of chlorine concentration levels collected in a water distribution system. The observations in each row are collected at a certain location, and the columns correspond to observations at consecutive time points.

The second is compressed sensing of the MIT logo [14], [21], which is a 38×72 gray-scale image, see Figure 3.



Figure 3: MIT logo

1) *Chlorine concentration recovery*: We use a sub-matrix of size 120×180 as our ground truth matrix X^* and generated $p = 4200$ sensing matrices and linear measurements according to (2). We set $r = 6$ since the rank-6 truncated SVD of X^* achieves a relative low error of 0.069. In this experiment, as in [22], we use outliers to replace measurements b_i s, and the outliers are sampled from $\mathcal{N}(0, 10 \| X^* \|_F)$. The replacement happens with probability 1% or 5%, which is denoted by error distributions (i) and (j), respectively. As a comparison, we

Error	Loss	RE	RR
(a)	Huber	0.095	0.130
	L_2	0.102	0.121
(c)	Huber	0.448	0.032
	L_2	0.432	0.032
(i)	Huber	0.097	0.132
	L_2	12.587	0.002
(j)	Huber	0.114	0.109
	L_2	N.A.	

Table III: Performance of Chlorine concentration recovery with four error distributions: (a) No error; (c) $\mathcal{N}(0, 10)$; (i) 1% outliers with outliers sampled from $\mathcal{N}(0, 10\|X^*\|_F)$; (j) 5% outliers with outliers sampled from $\mathcal{N}(0, 10\|X^*\|_F)$.

also consider the cases with no error and the error distribution of $\mathcal{N}(0, 10)$, which correspond to error distributions (a) and (c) in the synthetic data, respectively.

Table III gives the sensing performance under different levels of outliers. The recovery rates are all low, and this possibly is due to the fact that X^* has many entries close to zero. As a result, the relative error might be a more informative metric here. The Huber and L_2 loss have comparable performance when there is no error (see also Figure 4) or normal error of a moderate scale (see also Figure 5). Replacing 1% or 5% of the measurements with $\mathcal{N}(0, 10\|X^*\|_F)$ outliers hardly affects the recovery if Huber loss is used. However, for L_2 loss, the result is severely affected by outliers. When there are 5% outliers, the algorithm is either too slow or divergent, thus the result is not available (N.A.).

Figure 4-7 visualize the element-wise comparison between the true matrix and the recovery results by plotting the second row of the matrices in a plot. Figure 4 shows that recovery is good when there is no error. Figure 5 shows that both procedures are affected by the $\mathcal{N}(0, 10)$ noise, but the reconstructed values still share the same pattern with the true values. When there are 1% outliers, using the L_2 loss results in large fluctuations and can not recover the matrix (see Figure 6). In contrast, using Huber loss allows us to recover the true X^* even when 5% of the observations are outliers, see Figure 7.

2) *Compressed sensing of MIT logo*: We use the gray-scale logo as the ground truth matrix, which can be well-approximated by a rank-4 matrix (RE = 0.0194 by truncated SVD). We set $m = 38, n = 72, r = 4$, and take $p = 1200$ measurements using (2). Figure 8 visually compares the sensing results under L_2 loss and Huber loss.

If Huber loss is used, the MIT logo can be recovered in all 4 scenarios examined. Recall that the $\mathcal{N}(0, 10)$ noise is the most adverse case for Huber loss among the eight error distributions tested using synthetic dataset. In this MIT logo experiment, the reconstructed image is still recognizable when error follows $\mathcal{N}(0, 10)$. For the other 3 types of errors, the recovery error is hardly noticeable.

On the other hand, using L_2 loss allows us to recover the logo when there is no or Normal error. Under Normal error,

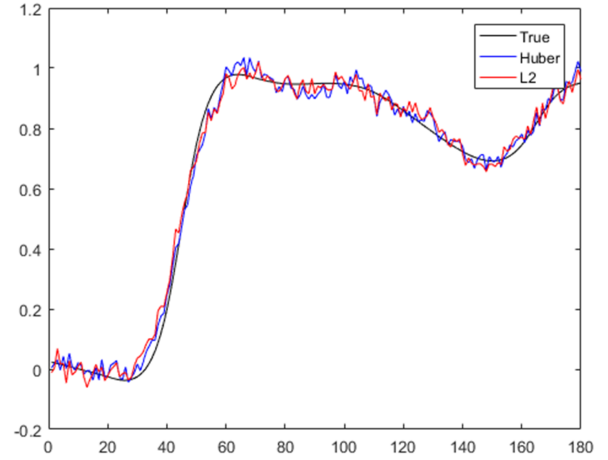


Figure 4: Chlorine concentration recovery: no error

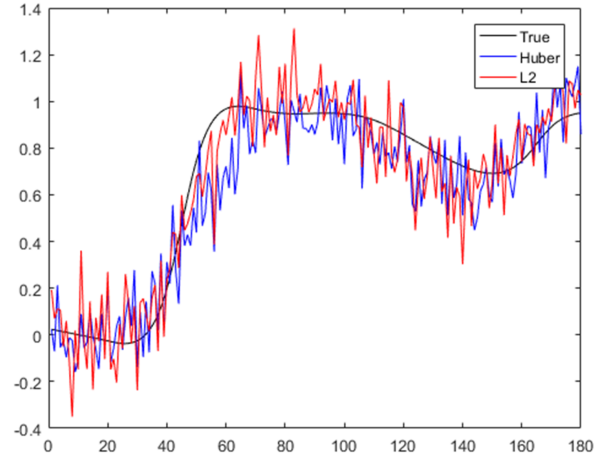


Figure 5: Chlorine concentration recovery: $\mathcal{N}(0, 10)$ error

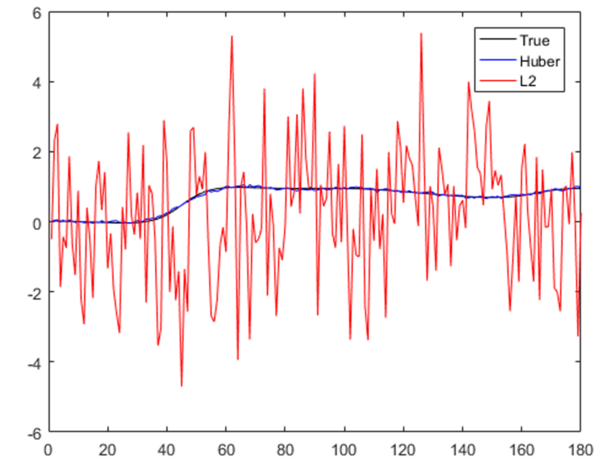


Figure 6: Chlorine concentration recovery: 1% outlier

L_2 loss has a slight advantage over Huber loss, as the outcome

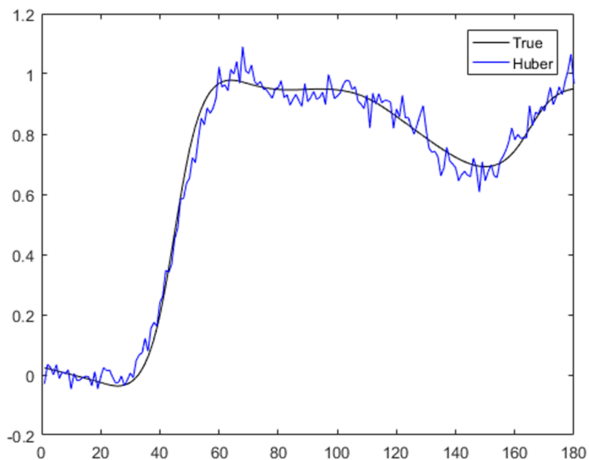


Figure 7: Chlorine concentration recovery: 5% outlier

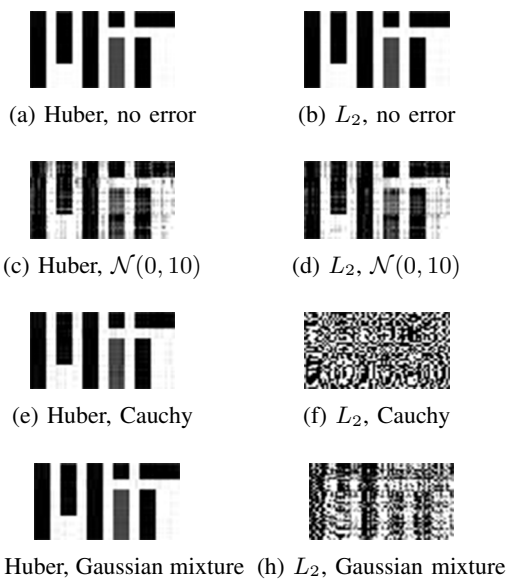


Figure 8: Compressed sensing of MIT logo. (a) No error; (b) $\mathcal{N}(0, 2)$; (c) $\mathcal{N}(0, 10)$; (d) $\log\mathcal{N}(0, 1)$; (e) Cauchy; (f) $t(3)$; (g) Pareto(1, 1); (h) $0.9\mathcal{N}(0, 1) + 0.1\mathcal{N}(100, 1)$.

is a bit clearer. The reconstructed image is very blurred when Gaussian mixture error presents, and can not be recognized under Cauchy error. Numerical metrics are supplied in Table IV.

VII. CONCLUSION

This paper proposes to use M-estimation in the problem of matrix factorization for low-rank matrix recovery. The proposed framework is very general, and it includes the popular L_1 , L_2 and quantile loss functions as special cases. Especially the Huber M-estimation enjoys good robustness property, and the loss function is smooth. The alternating minimization method is applied at each step, and Nesterov's momentum method is employed to accelerate the convergence. The numerical experiments on both synthetic data and real data

Error	Loss	RE	RR
(a)	Huber	0.024	0.547
	L_2	0.026	0.538
(c)	Huber	0.275	0.105
	L_2	0.235	0.113
(e)	Huber	0.063	0.380
	L_2	9.795	0.002
(h)	Huber	0.043	0.487
	L_2	0.951	0.035

Table IV: Recovery of MIT logo with four error distributions: (a) No error; (c) $\mathcal{N}(0, 10)$; (e) Cauchy; (h) $0.9\mathcal{N}(0, 1) + 0.1\mathcal{N}(100, 1)$.

demonstrate that the Huber loss function provides more robust performance for skewed and/or heavy-tailed data comparing with L_2 loss.

APPENDIX

Appendix A: Proof Sketch of Theorem 1

The proof of the convergence shares the same nature with the combined analysis of [13], [23] and [17], thus we only provide a road map.

We first present theorems where \mathcal{M} is smooth. The following assumptions are made on the function $\mathcal{M}(\cdot)$:

- 1 Assume that $\mathcal{M}(\mathcal{A}(Y) - b) - \mathcal{M}(\mathcal{A}(X) - b) \geq \langle \nabla \mathcal{M}(\mathcal{A}(X) - b), Y - X \rangle + \xi \|\mathcal{A}(X) - \mathcal{A}(Y)\|_2^2$. This assumption is the ξ -strong convexity assumption in [24].
- 2 Assume that $\mathcal{M}(\mathcal{A}(Y) - b) - \mathcal{M}(\mathcal{A}(X) - b) \leq \langle \nabla \mathcal{M}(\mathcal{A}(X) - b), Y - X \rangle + \eta \|\mathcal{A}(X) - \mathcal{A}(Y)\|_2^2$.
- 3 Without loss of generality, assume X^* is the optimal matrix, then $\mathcal{M}(\mathcal{A}(X^*) - b) = 0$ and $\mathcal{M}(X) \geq 0$.
- 4 Assume that \mathcal{M} also defines a matrix norm when acts on a matrix X , and $c_1 \|X\|_2^2 \leq \mathcal{M}(X) \leq c_2 \|X\|_2^2$.

Remark: Assumption 1-2 defines the condition number of the function $\mathcal{M}(\cdot)$, similar assumptions also appear in [25]. Usually $\kappa = \eta/\xi$ is called the condition number of a function f [26].

Lemma A.1. [27] Let \mathcal{A} satisfy $2r$ -RIP with constant δ_{2r} . Then for all matrices X, Y of rank at most r , we have

$$|\langle \mathcal{A}(X), \mathcal{A}(Y) \rangle - \langle X, Y \rangle| \leq \delta_{2r} \|X\|_F \|Y\|_F.$$

The next lemma characterizes the convergence rate of the initialization procedure:

Lemma A.2. [28] Let $X \in \mathbb{R}^{m \times n}$ be an arbitrary matrix of rank r . Let $b = \mathcal{A}(X) \in \mathbb{R}^p$ be p linear measurements. Consider the iterative updates

$$Y^{t+1} \leftarrow \mathcal{P}_r(Y^t - \xi_t \nabla_X \mathcal{M}(\mathcal{A}(Y^t) - b)),$$

where Y^i are $m \times n$ matrices. Then

$$\|Y^t - X\|_F \leq \psi(\mathcal{A})^t \|Y^0 - X\|_F$$

holds, where $\psi(\mathcal{A})$ is defined as

$$\psi(\mathcal{A}) = 2 \sup_{\|X\|_F = \|Y\|_F = 1, \text{rank}(X) \leq 2r, \text{rank}(Y) \leq 2r} |\langle \mathcal{A}(X), \mathcal{A}(Y) \rangle - \langle X, Y \rangle|.$$

We can prove this lemma by using the results of Theorem A.1.

First, we prove that the initialization procedure indeed converges to the true value of X .

A. Proof of Eq. (8) in Theorem 1

Lemma A.3 (Initialization). *Assume that $\xi \frac{1+\delta_{2k}}{1-\delta_{2k}} - \eta > 0$. Denote $\tilde{\mathcal{M}}(X) = \mathcal{M}(\mathcal{A}(X) - b)$. Let X^* be an optimal solution and let X^t be the value obtained by Algorithm 1 at t^{th} iteration. Then*

$$\tilde{\mathcal{M}}(X^{t+1}) \leq \tilde{\mathcal{M}}(X^t) + \left(\xi \frac{1+\delta_{2k}}{1-\delta_{2k}} - \eta \right) \|\mathcal{A}(X^* - X^t)\|_F^2.$$

Proof. From assumption, we have

$$\begin{aligned} & \tilde{\mathcal{M}}(X^{t+1}) - \tilde{\mathcal{M}}(X^t) \\ & \leq \langle \nabla \tilde{\mathcal{M}}(X^t), X^{t+1} - X^t \rangle + \xi \|\mathcal{A}(X^{t+1}) - \mathcal{A}(X^t)\|_F^2 \\ & \leq \langle \nabla \tilde{\mathcal{M}}(X^t), X^{t+1} - X^t \rangle + \xi(1 + \delta_{2k}) \|X^{t+1} - X^t\|_F^2 \end{aligned}$$

where the last inequality comes from RIP. Let $Y^{t+1} = X^t - \frac{1}{2\xi(1+\delta_{2k})} \nabla \tilde{\mathcal{M}}(X^t)$, and

$$f_t(X) = \langle \nabla \tilde{\mathcal{M}}(X^t), X - X^t \rangle + \xi(1 + \delta_{2k}) \|X - X^t\|_F^2.$$

Then

$$f_t(X) = \xi(1 + \delta_{2k}) \left[\|X - Y^{t+1}\|_F^2 - \frac{1}{4\xi^2(1 + \delta_{2k})^2} \|\nabla \tilde{\mathcal{M}}(X^t)\|_F^2 \right]$$

By definition, $\mathcal{P}_k(Y^{t+1}) = X^{t+1}$, then $f_t(X^{t+1}) \leq f_t(X^*)$. Thus

$$\begin{aligned} & \tilde{\mathcal{M}}(X^{t+1}) - \tilde{\mathcal{M}}(X^t) \leq f_t(X^{t+1}) \leq f_t(X^*) \\ & = \langle \nabla \tilde{\mathcal{M}}(X^t), X^* - X^t \rangle + \xi(1 + \delta_{2k}) \|X^* - X^t\|_F^2 \\ & \leq \nabla \tilde{\mathcal{M}}(X^t), X^* - X^t + \xi \frac{1 + \delta_{2k}}{1 - \delta_{2k}} \|\mathcal{A}(X^*) - \mathcal{A}(X^t)\|_F^2 \\ & \leq \nabla \tilde{\mathcal{M}}(X^t), X^* - X^t + \eta \|\mathcal{A}(X^*) - \mathcal{A}(X^t)\|_F^2 \\ & \quad + \left(\xi \frac{1 + \delta_{2k}}{1 - \delta_{2k}} - \eta \right) \|\mathcal{A}(X^*) - \mathcal{A}(X^t)\|_F^2 \\ & \leq \tilde{\mathcal{M}}(X^*) - \tilde{\mathcal{M}}(X^t) \\ & \quad + \left(\xi \frac{1 + \delta_{2k}}{1 - \delta_{2k}} - \eta \right) \|\mathcal{A}(X^*) - \mathcal{A}(X^t)\|_F^2 \end{aligned}$$

□

Theorem A.1. *Let $b = \mathcal{A}(X^*) + e$ for rank k matrix X^* for an error vector $e \in \mathbb{R}^p$, $D = \frac{1}{C^2} + \left(\xi \frac{1+\delta_{2k}}{1-\delta_{2k}} - \eta \right) \left(\frac{2}{C^2} + \sqrt{\frac{2}{c_1}} \frac{1}{C} + \frac{1}{c_1} \right)$. Then, under the assumption that $D < 1$, Algorithm 1 with step size $\eta_t = \frac{1}{2\xi(1+\delta_{2k})}$ outputs a matrix X of rank at most k such that $\mathcal{M}(\mathcal{A}(X) - b) \leq (C^2 + \varepsilon) \frac{\|e\|_2^2}{2}$, where $\varepsilon \geq 0$, in at most $\left\lceil \frac{1}{\log D} \log \frac{(C^2 + \varepsilon) \|e\|_2^2}{2c_2 \|b\|_2^2} \right\rceil$ iterations.*

Proof. Let the current solution X^t satisfy $\mathcal{M}(X^t) \geq \frac{C^2 \|e\|_2^2}{2}$. By lemma A.3 and $b - \mathcal{A}(X^*) = e$, we have

$$\begin{aligned} \mathcal{M}(X^{t+1}) & \leq \frac{\|e\|_2^2}{2} + \left(\xi \frac{1 + \delta_{2k}}{1 - \delta_{2k}} - \eta \right) \|b - \mathcal{A}(X^t) - e\|_2^2 \\ & \leq \frac{\|e\|_2^2}{2} + \left(\xi \frac{1 + \delta_{2k}}{1 - \delta_{2k}} - \eta \right) (\|b - \mathcal{A}(X^t)\|_2^2 - 2e^\top (b - \mathcal{A}(X^t)) + \|e\|_2^2) \\ & \leq \frac{\mathcal{M}(X^t)}{C^2} + \left(\xi \frac{1 + \delta_{2k}}{1 - \delta_{2k}} - \eta \right) \left(\frac{2\mathcal{M}(X^t)}{C^2} + \frac{\mathcal{M}(X^t)}{c_1} + \frac{\sqrt{2}\mathcal{M}(X^t)}{C\sqrt{c_1}} \right) \\ & = D\mathcal{M}(X^t). \end{aligned}$$

Since $D < 1$, combining the fact that $\mathcal{M}(X^0) \leq c_2 \|b\|_2^2$, by taking $t = \left\lceil \frac{1}{\log D} \log \frac{(C^2 + \varepsilon) \|e\|_2^2}{2c_2 \|b\|_2^2} \right\rceil$, we complete the proof. □

The following lemma is adapted from [13]:

Lemma A.4. *Let $X_1, X_2 \in \mathbb{R}^{m \times n}$ be two rank r matrices with SVD decomposition $X_1 = A_1^\top \Sigma_1 B_1$, $X_2 = A_2^\top \Sigma_2 B_2$. For $l = 1, 2$, define $U_l = A_l^\top \Sigma_l^{1/2} \in \mathbb{R}^{m \times r}$, $V_l = B_l^\top \Sigma_l^{1/2} \in \mathbb{R}^{n \times r}$. Assume X_1, X_2 obey $\|X_2 - X_1\| \leq \frac{1}{2} \sigma_r(X_1)$. Then*

$$\text{dist}^2 \left(\begin{bmatrix} U_2 \\ V_2 \end{bmatrix}, \begin{bmatrix} U_1 \\ V_1 \end{bmatrix} \right) \leq \frac{2}{\sqrt{2} - 1} \frac{\|X_2 - X_1\|_F^2}{\sigma_r(X_1)}.$$

Combine Lemma A.1 and A.4, and follow a similar route of [13], we can prove that using more than $3 \log(\sqrt{r}\kappa) + 5$ iterations of Algorithm 1, we can obtain

$$\text{dist} \left(\begin{bmatrix} U_0 \\ V_0 \end{bmatrix}, \begin{bmatrix} U \\ V \end{bmatrix} \right) \leq \frac{1}{4} \sigma_r(U).$$

Thus we finished the proof of convergence in the initialization procedure. Next, we prove the linear convergence of the main algorithm on the penalized objective function.

We first rewrite the objective function (7) by using the up-lifting technique, so that it is easier to simultaneously consider \mathcal{M} and the regularization term. To see this, consider rank r matrix $X \in \mathbb{R}^{m \times n}$ with SVD decomposition $X = U^\top \Sigma V$. Define $\text{Sym} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{(m+n) \times (m+n)}$ as

$$\text{Sym}(X) = \begin{bmatrix} 0_{m \times m} & X \\ X^\top & 0_{n \times n} \end{bmatrix}.$$

Given the block matrix $\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$ with $\mathbf{A}_{11} \in \mathbb{R}^{m \times m}$, $\mathbf{A}_{12} \in \mathbb{R}^{m \times n}$, $\mathbf{A}_{21} \in \mathbb{R}^{n \times m}$, $\mathbf{A}_{22} \in \mathbb{R}^{n \times n}$, define $\mathcal{P}_{\text{diag}}(\mathbf{A}) = \begin{bmatrix} \mathbf{A}_{11} & 0_{m \times n} \\ 0_{n \times m} & \mathbf{A}_{22} \end{bmatrix}$ and $\mathcal{P}_{\text{off}}(\mathbf{A}) = \begin{bmatrix} 0_{m \times m} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & 0_{n \times n} \end{bmatrix}$. Define $\mathcal{B} : \mathbb{R}^{(m+n) \times (m+n)} \rightarrow \mathbb{R}^p$ as the uplified version of the \mathcal{A} operator

$$\mathcal{B}(X)_k = \langle B_k, X \rangle, \text{ where } B_k = \text{Sym}(A_k).$$

Define $W = [U^\top, V^\top]$. As a result, we can rewrite object function (7) as:

$$\begin{aligned} g(W) &= g(U, V) \\ &= \mathcal{M}(b - \mathcal{A}(U^\top V)) + \lambda \|UU^\top - VV^\top\|_F^2 \\ &= \frac{1}{2} \mathcal{M}(\mathcal{B}(\text{Sym}(U^\top V)) - \text{Sym}(X)) \\ &\quad + \frac{1}{2} \lambda \|\text{Sym}(U^\top V) - \text{Sym}(X)\|_F^2. \end{aligned} \quad (10)$$

From equation (10) we can see that two parts of the penalized objective function have similar structures.

As a result, although we made Assumptions 1, 2, 4 on the original function \mathcal{M} , we can see from (10) that the penalized objective function still retains a similar property. As for Assumption 3, we can also use some location transform techniques to make the penalized objective function satisfies this assumption. Thus, we can deal with the penalized objective function or the unpenalized objective function in a similar way.

Then, similar to [13], the alternating minimization together with Nesterov's momentum algorithm with respect to U and V (sub-matrices of W) can be written as a NAG algorithm applied to $g(W)$ with respect to W .

For the convergence analysis of the Nesterov's momentum algorithm, we employ the following lemma, which is a theorem in [17].

Lemma A.5. *For minimization problem $\min_{x \in \mathcal{X}} f(x)$, where x is a vector, using the Lyapunov function*

$$\tilde{V}_k = f(y_k) + \xi \|z_k - x^*\|^2,$$

it can be shown that

$$\tilde{V}_{k+1} - \tilde{V}_k = -\tau_k \tilde{V}_k + \varepsilon_{k+1} \quad (11)$$

where the error is expressed as

$$\varepsilon_{k+1} = \left(\frac{\tau_k^2}{4\xi} \right) \|\nabla f(x_k)\|^2 + \left(\tau_k \eta - \frac{\xi}{\tau_k} \right) \|x_k - y_k\|^2,$$

τ_k is the step size in Nesterov's momentum algorithm, usually equals $1/\sqrt{k}$, $y_{k+1} = x_k - \frac{1}{2\eta} \nabla f(x_k)$, $x_{k+1} = \frac{1}{1+\tau_k} y_k + \frac{\tau_k}{1+\tau_k} z_k$, $z_{k+1} = z_k + \tau_k \left(x_{k+1} - z_k - \frac{1}{2\xi} \nabla f(x_{k+1}) \right)$.

Assume that $\tau_0 = 0$, $\tau_1 = \tau_2 = \dots = \tilde{\tau}$, and $\varepsilon_1, \dots, \varepsilon_{k+1}$ has a common upper bound $\tilde{\varepsilon}$, then (11) implies:

$$\begin{aligned} |\tilde{V}_{k+1}| &= |(1 - \tilde{\tau})^{k+1} \tilde{V}_0 + \sum_{i=1}^{k+1} (1 - \tilde{\tau})^{i-1} \varepsilon_{k+2-i}| \\ &\leq (1 - \tilde{\tau})^{k+1} |\tilde{V}_0| + \frac{\tilde{\varepsilon} - \tilde{\varepsilon}(1 - \tilde{\tau})^k}{\tilde{\tau}} \end{aligned}$$

Substitute x_k with $W_{k+1} = [U^{k+1}, V^{k+1}]^\top$, f with g . To deal with the convergence analysis with respect to $W_t - W^*$ and $W_0 - W^*$, we need to handle two parts. The first part is the error part with respect to $\tilde{\varepsilon}$. This can be solved by choosing an initial estimate close to the true value. As a result, $\|\nabla f(x_1)\|$ can be arbitrary close to 0. For notational simplicity, assume

that $\frac{\tilde{\varepsilon} - \tilde{\varepsilon}(1 - \tilde{\tau})^k}{\tilde{\tau}} \leq \varepsilon^\dagger$. Since \tilde{V}_k still satisfies Assumption 1 and 2, without loss of generality, assume the corresponding parameter are $\tilde{\xi}$ and $\tilde{\eta}$.

Next we want to study the relation between $W_t - W^*$ and \tilde{V}_t . This involves Assumption 1 and 2 as well as Lemma A.1. With rough handling of the gradient part in Assumption 1 and 2, we can obtain

$$\tilde{\xi} \|\mathcal{A}(W_t) - \mathcal{A}(W^*)\|_2^2 \leq (1 - \tilde{\tau})^t \tilde{\mu} \|\mathcal{A}(W_0) - \mathcal{A}(W^*)\|_2^2 + \tilde{\varepsilon}.$$

Notice that $\tilde{\varepsilon}$ can be made arbitrary small so that

$$\tilde{\xi} \|\mathcal{A}(W_t) - \mathcal{A}(W^*)\|_2^2 \leq (1 - \tilde{\tau}_1)^t \tilde{\mu} \|\mathcal{A}(W_0) - \mathcal{A}(W^*)\|_2^2$$

and $1 - \tilde{\tau}_1$ can still be larger than 0 smaller than 1. Employ the RIP property, we have

$$\tilde{\xi}(1 - \delta_r) \|W_t - W^*\|_2^2 \leq (1 - \tilde{\tau}_1)^t \tilde{\mu}(1 + \delta_r) \|W_0 - W^*\|_2^2.$$

Thus

$$\begin{aligned} &\text{dist} \left(\begin{bmatrix} U_t \\ V_t \end{bmatrix}, \begin{bmatrix} U \\ V \end{bmatrix} \right) \\ &\leq (1 - \tilde{\tau}_1)^t \frac{\tilde{\mu}}{\tilde{\xi}} \frac{1 + \delta_r}{1 - \delta_r} \text{dist} \left(\begin{bmatrix} U_0 \\ V_0 \end{bmatrix}, \begin{bmatrix} U \\ V \end{bmatrix} \right) \\ &\leq \frac{1}{4} (1 - \tilde{\tau}_1)^t \frac{\tilde{\mu}}{\tilde{\xi}} \frac{1 + \delta_r}{1 - \delta_r} \sigma_r(U). \end{aligned}$$

Theorem 1 is proved.

Remark on (10): From (10), we provide a guideline with respect to the selection of λ compared with [13], by combining (10) and Assumption 4.

Appendix B: Proof Sketch of Theorem 2

Proof. Employ Theorem 1 in [29] and Lemma 3 in [30] we know that

$$\mathcal{M}_\pi^\lambda(b - \mathcal{A}(U^{\pi* \top} V^{\pi*})) \rightarrow \mathcal{M}^\lambda(b - \mathcal{A}(U^{*\top} V^*))$$

Giving the fact that \mathcal{M}_π^λ and \mathcal{M}^λ are convex, as well as the restricted isometry property, we finish the proof of Theorem 2. \square

REFERENCES

- [1] Y. Koren, R. Bell, and C. Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, no. 8, pp. 30–37, 2009.
- [2] Z. Lin, C. Xu, and H. Zha, "Robust matrix factorization by majorization minimization," *IEEE transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 1, pp. 208–220, 2018.
- [3] R. Basri, D. Jacobs, and I. Kemelmacher, "Photometric stereo with general, unknown lighting," *International Journal of Computer Vision*, vol. 72, no. 3, pp. 239–257, 2007.
- [4] M. A. Davenport and J. Romberg, "An overview of low-rank matrix recovery from incomplete observations," *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 4, pp. 608–622, 2016.
- [5] D. Bindel, "Matrix factorizations for computer network tomography," in *Householder Symposium XVIII on Numerical Linear Algebra*, p. 27, 2011.
- [6] S. C. Ponz, *Machine Learning based Models for Matrix Factorization*. PhD thesis, 2017.
- [7] R. Zhu, D. Niu, L. Kong, and Z. Li, "Expectile matrix factorization for skewed data analysis," in *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [8] R. He, T. Tan, and L. Wang, "Robust recovery of corrupted low-rank matrix by implicit regularizers," *IEEE transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 4, pp. 770–783, 2014.

- [9] J. Tukey, "A survey of sampling from contaminated distributions," in *Contributions to probability and statistics* (I. Olkin, ed.), pp. 448–485, Stanford University Press, 1960.
- [10] F. Hampel, *Contribution to the theory of robust estimation*. Phd thesis, University of California, Berkeley, 1968.
- [11] P. J. Huber and E. M. Ronchetti, *Robust Statistics*. Wiley, 2009.
- [12] Y. Nesterov, "Smooth minimization of non-smooth functions," *Mathematical Programming*, vol. 103, no. 1, pp. 127–152, 2005.
- [13] S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, and B. Recht, "Low-rank solutions of linear matrix equations via procrustes flow," in *International Conference on Machine Learning*, pp. 964–973, 2016.
- [14] P. Jain, R. Meka, and I. S. Dhillon, "Guaranteed rank minimization via singular value projection," in *Advances in Neural Information Processing Systems*, pp. 937–945, 2010.
- [15] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *International Conference on Machine Learning*, pp. 1139–1147, 2013.
- [16] J. Zhang and I. Mitliagkas, "Yellowfin and the art of momentum tuning," *arXiv preprint arXiv:1706.03471*, 2017.
- [17] R. Brooks, "Convergence analysis of deterministic and stochastic methods for convex optimization," Master's thesis, University of Waterloo, 2017.
- [18] A. Y. Aravkin, A. Kambadur, A. C. Lozano, and R. Luss, "Sparse quantile huber regression for efficient and robust estimation," *arXiv preprint arXiv:1402.4624*, 2014.
- [19] Z. Yang, Y. Zhang, W. Yan, Y. Xiang, and S. Xie, "A fast non-smooth nonnegative matrix factorization for learning sparse representation," *IEEE access*, vol. 4, pp. 5161–5168, 2016.
- [20] S. Papadimitriou, J. Sun, and C. Faloutsos, "Streaming pattern discovery in multiple time-series," in *Proceedings of the 31st International Conference on Very Large Data Bases*, pp. 697–708, VLDB Endowment, 2005.
- [21] B. Recht, M. Fazel, and P. A. Parrilo, "Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization," *SIAM review*, vol. 52, no. 3, pp. 471–501, 2010.
- [22] Y. Li, Y. Chi, H. Zhang, and Y. Liang, "Nonconvex low-rank matrix recovery with arbitrary outliers via median-truncated gradient descent," *arXiv preprint arXiv:1709.08114*, 2017.
- [23] E. J. Candes, X. Li, and M. Soltanolkotabi, "Phase retrieval via wirtinger flow: Theory and algorithms," *IEEE Transactions on Information Theory*, vol. 61, no. 4, pp. 1985–2007, 2015.
- [24] D. Park, A. Kyrillidis, C. Caramanis, and S. Sanghavi, "Finding low-rank solutions to matrix problems, efficiently and provably," *arXiv preprint arXiv:1606.03168*, 2016.
- [25] W. Su, S. Boyd, and E. Candes, "A differential equation for modeling nesterov's accelerated gradient method: Theory and insights," in *Advances in Neural Information Processing Systems*, pp. 2510–2518, 2014.
- [26] S. Bhojanapalli, A. Kyrillidis, and S. Sanghavi, "Dropping convexity for faster semi-definite optimization," in *Conference on Learning Theory*, pp. 530–582, 2016.
- [27] E. J. Candes, "The restricted isometry property and its implications for compressed sensing," *Comptes rendus mathematique*, vol. 346, no. 9-10, pp. 589–592, 2008.
- [28] S. Oymak, B. Recht, and M. Soltanolkotabi, "Sharp time–data tradeoffs for linear inverse problems," *IEEE Transactions on Information Theory*, vol. 64, no. 6, pp. 4129–4158, 2018.
- [29] N. Guan, D. Tao, Z. Luo, and J. Shawe-Taylor, "Mahnmf: Manhattan non-negative matrix factorization," *arXiv preprint arXiv:1207.3438*, 2012.
- [30] O. Fercoq and P. Richtárik, "Smooth minimization of nonsmooth functions with parallel coordinate descent methods," *arXiv preprint arXiv:1309.5885*, 2013.