

Low-rank matrix recovery under heavy-tailed errors

MYEONGHUN YU^{1,a}, QIANG SUN^{2,b} and WEN-XIN ZHOU^{3,c}

¹Department of Mathematics, University of California, San Diego, La Jolla, CA, 92093, USA, myyu@ucsd.edu

²Department of Statistical Sciences, University of Toronto, Toronto, ON M5G 1Z5, Canada,

^bqiang.sun@utoronto.ca

³Department of Information and Decision Sciences, University of Illinois at Chicago, Chicago, IL, 60607, USA,

^cwenxinz@uic.edu

This paper proposes convex relaxation based robust methods to recover approximately low-rank matrices in the presence of heavy-tailed and asymmetric errors, allowing for heteroscedasticity. We focus on three archetypal applications in matrix recovery: matrix compressed sensing, matrix completion and multitask regression. Statistically, we provide sub-Gaussian-type deviation bounds when the noise variables only have bounded variances in each aforementioned setting. Improving upon the earlier results in Fan, Wang and Zhu (*Ann. Statist.* **49** (2021) 1239–1266), the convergence rates of our estimators are proportional to the noise scale under matrix sensing and multitask regression settings, and thus diminish to 0 in the noiseless case. Computationally, we propose a matrix version of the local adaptive majorize-minimization algorithm, which is much faster than the alternating direction method of multiplier used in previous work and is scalable to large datasets. Numerical experiments demonstrate the advantage of our methods over their non-robust counterparts and corroborate the theoretical findings that the convergence rates are proportional to the noise scale.

Keywords: Heavy-tailed data; Huber loss; low-rank matrix recovery; nuclear norm; trace regression

1. Introduction

There has been a recent surge of interest in matrix recovery which aims to recover an unknown matrix from noisy observations. Matrix recovery has wide applications in practice, including collaborative filtering [14], multitask regression [1], quantum state tomography [15] and face recognition [24], to name a few. Statistically, it aims to estimate $\Theta^* \in \mathbb{R}^{d_1 \times d_2}$ based on n independently and identically distributed (i.i.d.) observations $\{(y_i, X_i)\}_{i=1}^n$ following the generative model

$$y_i = \text{tr}(X_i^T \Theta^*) + \epsilon_i =: \langle X_i, \Theta^* \rangle + \epsilon_i,$$

where $X_i \in \mathbb{R}^{d_1 \times d_2}$ is a random measurement matrix, and ϵ_i is an error variable satisfying $\mathbb{E}(\epsilon_i | X_i) = 0$ and $\mathbb{E}(\epsilon_i^2 | X_i) \leq \sigma_0^2$ for some $\sigma_0 > 0$. We consider matrix recovery in high dimensions, that is, $d_1 \times d_2$ can be much larger than the sample size n , making the problem ill-posed. It has been a common practice to assume that Θ^* is (approximately) low-rank, and the resulting problem is referred to as low-rank matrix recovery.

The problem of low-rank matrix recovery can be naturally formulated as a nonconvex empirical risk minimization problem subject to a rank constraint. To find local optima of such a rank-constrained program, commonly used methods are Riemannian gradient descent [41] and Burer-Monteiro type gradient descent [2, 7, 25, 38]. The former views the set of rank- r matrices as a smooth manifold, while the latter relies on the matrix factorization $\Theta = UV^T$, where $U \in \mathbb{R}^{d_1 \times r}$, $V \in \mathbb{R}^{d_2 \times r}$, and $r = \text{rank}(\Theta^*)$ is assumed to be known. To relax the restrictive assumption that the true rank r is known *a priori*, [22]

and [43] studied the gradient method for solving the reparameterized program in an over-parameterized regime where $\Theta = UV^T$ with $U \in \mathbb{R}^{d_1 \times r'}$ and $V \in \mathbb{R}^{d_2 \times r'}$, and $r' \geq r$ is an upper bound of the true rank.

Another line of research resorts to convex relaxation in order to obtain computationally feasible solutions. Similar to Lasso [37] in the context of sparse linear regression, convex low-rank matrix recovery methods are based on either constrained nuclear norm minimization or nuclear norm penalized least squares formulation. The nuclear norm of a matrix is defined as the sum of its singular values, and thus serves as a convex surrogate for its rank. We refer to [4–6, 28–31] and [18] for an unavoidably incomplete list of notable works on exact and noisy low-rank matrix recovery through convex relaxation. In the context of multitask learning, [23] introduced an approach that utilizes the group Lasso penalty when only a small number of rows in the matrix Θ^* are nonzero.

All the aforementioned methods, convex or nonconvex, are studied either in the noiseless setting or under a sub-Gaussian/sub-exponential assumption on the random error. However, both convex and nonconvex least squares estimators exhibit sub-optimal deviation bounds in the presence of heavy-tailed errors that only have a small number of finite moments. To make the estimator less sensitive to heavy-tailedness, a natural idea is to replace the ℓ_2 -loss with a robust loss function, such as the ℓ_1 -loss or the Huber loss [16]. For example, [10] proposed and studied nuclear norm penalized estimators using both the ℓ_1 -loss and the Huber loss; [35] considered robust sparse reduced rank regression by minimizing the empirical Huber loss plus a combination of the nuclear norm and entry-wise ℓ_1 -norm penalties. For methods that rely on nonconvex optimization with robust losses, [33] proposed a Riemannian sub-gradient method and proved the statistical properties of the iterates; [40] studied the statistical properties of vanilla gradient descent iterates for solving reparameterized (regularized) Huber loss minimization. In an over-parameterized regime, [26] showed that sub-gradient descent with the ℓ_1 -loss function converges to the ground truth at a near-linear rate in the presence of arbitrarily large outliers.

In this paper, we propose a robust approach to recover an approximately low-rank matrix in a trace regression model with heavy-tailed and asymmetric error, which complements the extant literature on low-rank matrix recovery via convex relaxation. Borrowing ideas from robust (sparse) linear regressions [11, 34], we adopt the Huber loss function with a diverging robustification parameter to achieve sub-Gaussian-type concentration bounds. We focus on three archetypal examples in matrix recovery: matrix compressed sensing, matrix completion and multitask regression. For each problem, we study the nonasymptotic deviation bounds of the nuclear norm penalized Huber estimator under both the Frobenius and nuclear norms, which match the minimax optimal rates. Our main contributions are as follows. First and foremost, we provide a comprehensive analysis of the nuclear norm penalized Huber regression estimator to gain robustness without compromising statistical efficiency. Our results either improve or complement those in [10, 13] and [35]. For example, [10] considered robust matrix completion under symmetric error distribution and also required a constant lower bound for the error density function. [35] examined the Huber-type estimator for sparse multitask regression but their analysis cannot be directly extended to the non-sparse setting. Secondly, we provide a unified algorithmic framework, which is a matrix variant of the local adaptive majorize-minimization (LAMM) algorithm [12], to solve the three problems (matrix sensing, matrix completion and multitask regression) all at once. By constructing an isotropic quadratic function that locally majorizes the empirical Huber loss, the solution to each proximal optimization problem has a closed form, which considerably facilitates the implementation. Compared to many other algorithms used in the literature, our algorithm is first-order and thus more scalable to large data sets.

1.1. Related work and paper organization

Our model setting is closely related to that in [13], but the proposed robust estimators provably achieve sharper convergence rates than those obtained in [13]. More specifically, [13] proposed a two-step procedure, which in step one applies shrinkage operators to the empirical average $(1/n) \sum_{i=1}^n y_i \mathbf{X}_i$. The truncation level on y_i 's, which appears in the final convergence rate, depends on the variance of y_i and thus is not proportional to the noise scale. In contrast, by employing the adaptive Huber loss as in [34] and the localized analysis developed by [12], we show that the convergence rates of our estimators are proportional to the noise scale for matrix sensing and multitask regression; see Theorem 3.3, Theorem 3.5 and the subsequent remarks for details. Moreover, [13] required ϵ_i to have bounded $(2k)$ -th moment for some $k > 1$, while our estimators enjoy optimal rates as long as ϵ_i 's have bounded variances. On the computational aspect, compared to the contractive Peaceman-Rachford splitting method and the alternating direction method of multiplier (ADMM) employed by [13] to solve the nuclear norm penalized programs in step two, the proposed matrix variant of the LAMM algorithm is first-order and has a lower computational cost per iteration. See Section 2.2 for a more detailed comparison of computational complexity.

Although the proposed estimators satisfy exponential deviation bounds when the noise distribution is asymmetric and has finite variance, this advantage is accompanied by a trade-off: they sacrifice a considerable level of robustness when facing adversarial contamination of the data. This is due to the use of a robustification parameter that increases with the sample size. In the case of adversarial contamination, recent studies have introduced robust estimators that showcase resistance to a small proportion of arbitrary outliers. For example, in sparse linear regression with Gaussian errors and adversarially corrupted labels, [9] demonstrated that the ℓ_1 -penalized Huber's M -estimator attains the optimal rate of convergence, up to a logarithmic factor. Moreover, several recent studies [8, 19, 21, 36] have specifically tackled the challenge of arbitrary outliers in the context of matrix sensing and matrix completion. For multitask regression, a robust multitask (reduced-rank) regression approach was introduced by [32] for simultaneous modeling and outlier detection. To address data contamination caused by arbitrary outliers, they formulated the problem as a regularized multivariate regression with a sparse mean-shift parametrization and developed a thresholding-based iterative procedure for optimization. It is worth noting that our methods and theory diverge from the conventional notion of robust statistics. While the aforementioned works assume sub-Gaussian or Gaussian noises, our work places emphasis on the distinct assumption of heavy-tailed errors rather than corruption by (arbitrary) outliers. The proposed methods and analysis therefore provide a useful complement to the current body of research on robust matrix completion and reduced-rank regression.

The rest of the paper proceeds as follows. In Section 2, we first review the trace regression model with three prototypical applications. Next, we introduce the nuclear norm penalized robust matrix estimator via the use of adaptive Huber loss, followed by a unified algorithm that applies to all three settings. We provide non-asymptotic high probability bounds for the proposed estimators case-by-case in Section 3. Section 4 presents our numerical experiments, conducted to demonstrate the advantage of our methods over their non-robust counterparts and to corroborate the theoretical findings that the convergence rates are proportional to the noise scale under the matrix sensing and multitask regression settings. All the proofs are relegated to the Appendix in the Supplementary Material [42].

Notation. For a matrix $\mathbf{A} = (A_{jk})_{1 \leq j \leq d_1, 1 \leq k \leq d_2} \in \mathbb{R}^{d_1 \times d_2}$, its singular values are denoted as $\sigma_1(\mathbf{A}) \geq \sigma_2(\mathbf{A}) \geq \dots \geq \sigma_{\min(d_1, d_2)}(\mathbf{A})$. Define its operator norm $\|\mathbf{A}\|_2 = \sigma_1(\mathbf{A})$, its Frobenius norm $\|\mathbf{A}\|_F = \sum_{j=1}^{\min(d_1, d_2)} \sigma_j^2(\mathbf{A})$, its nuclear norm $\|\mathbf{A}\|_* = \sum_{j=1}^{\min(d_1, d_2)} \sigma_j(\mathbf{A})$ and its max norm $\|\mathbf{A}\|_\infty = \max_{1 \leq j \leq d_1} \max_{1 \leq k \leq d_2} |A_{jk}|$. For two matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d_1 \times d_2}$, let $\langle \mathbf{A}, \mathbf{B} \rangle$ be the matrix inner product

defined as $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^T \mathbf{B})$. We use $\text{vec}(\mathbf{A}) \in \mathbb{R}^{d_1 d_2}$ to denote the long vector obtained by stacking the columns of \mathbf{A} .

2. Robust matrix recovery via adaptive Huber loss

2.1. Model and methods

Suppose we have collected n i.i.d. data points $\{(y_i, \mathbf{X}_i)\}_{i=1}^n$ generated according to the following heteroscedastic trace regression model

$$y_i = \langle \mathbf{X}_i, \boldsymbol{\Theta}^* \rangle + \epsilon_i, \quad (2.1)$$

where $\mathbf{X}_i \in \mathbb{R}^{d_1 \times d_2}$'s are random measurement matrices, and ϵ_i 's are additive random noise variables satisfying $\mathbb{E}(\epsilon_i | \mathbf{X}_i) = 0$ and $\mathbb{E}(\epsilon_i^2 | \mathbf{X}_i) \leq \sigma_0^2$. Based on the noisy observations $\{(y_i, \mathbf{X}_i)\}_{i=1}^n$, we are interested in recovering the unknown matrix $\boldsymbol{\Theta}^* \in \mathbb{R}^{d_1 \times d_2}$ that is either exactly or approximately low-rank. More specifically, assume for some $0 \leq q \leq 1$ and $\rho > 0$ that

$$\boldsymbol{\Theta}^* \in \mathcal{B}_q(\rho) := \left\{ \boldsymbol{\Theta} \in \mathbb{R}^{d_1 \times d_2} : \sum_{j=1}^{\min(d_1, d_2)} \sigma_j(\boldsymbol{\Theta})^q \leq \rho \right\}. \quad (2.2)$$

In particular, $\mathcal{B}_0(\rho) = \{\boldsymbol{\Theta} \in \mathbb{R}^{d_1 \times d_2} : \text{rank}(\boldsymbol{\Theta}) \leq \rho\}$ denotes the set of matrices with rank at most ρ , and $\mathcal{B}_q(\rho)$ with $0 < q \leq 1$ is set of approximately low-rank matrices. Throughout the rest of the paper, we assume without loss of generality that $d_1 \geq d_2$.

The difficulty of recovering $\boldsymbol{\Theta}^*$ varies depending on the random structures of the measurement matrices \mathbf{X}_i . Below we list three prototypical applications of model (2.1), which will be the main focus of this work.

- (i) *Matrix sensing*: Matrix sensing often assumes that the entries of $\mathbf{X}_i \in \mathbb{R}^{d_1 \times d_2}$ are independently generated from the $\mathcal{N}(0, 1)$ distribution. More generally, $\text{vec}(\mathbf{X}_i)$'s are assumed to be zero-mean sub-Gaussian/sub-exponential random vectors.
- (ii) *Matrix completion*: In matrix completion, \mathbf{X}_i are randomly drawn from the set

$$\mathcal{X} = \{\mathbf{e}_j(d_1) \mathbf{e}_k^T(d_2), 1 \leq j \leq d_1, 1 \leq k \leq d_2\},$$

where $\mathbf{e}_1(d), \dots, \mathbf{e}_d(d)$ are the canonical basis vectors in \mathbb{R}^d .

- (iii) *Multitask regression*: The multitask (reduced-rank) regression assumes

$$\mathbf{y}_i = (\boldsymbol{\Theta}^*)^T \mathbf{x}_i + \boldsymbol{\epsilon}_i, \quad i = 1, \dots, n, \quad (2.3)$$

where $\mathbf{y}_i = (y_{i1}, \dots, y_{id_2})^T \in \mathbb{R}^{d_2}$ are observed response vectors, $\mathbf{x}_i \in \mathbb{R}^{d_1}$ are covariate vectors, $\boldsymbol{\Theta}^* \in \mathbb{R}^{d_1 \times d_2}$ is the target regression coefficient matrix, and $\boldsymbol{\epsilon}_i = (\epsilon_{i1}, \dots, \epsilon_{id_2})^T \in \mathbb{R}^{d_2}$ are independent zero-mean random noise vectors. For $i = 1, \dots, n$ and $k = 1, \dots, d_2$, define

$$y_{(i-1)d_2+k} = y_{ik}, \quad \mathbf{X}_{(i-1)d_2+k} = \mathbf{x}_i \mathbf{e}_k^T(d_2) \quad \text{and} \quad \epsilon_{(i-1)d_2+k} = \epsilon_{ik}. \quad (2.4)$$

Then the sample $\{(y_j, \mathbf{X}_j)\}_{j=1}^N$ with $N = nd_2$ satisfies model (2.1).

For matrix sensing and matrix completion with noisy measurements, a popular approach is the the following convex relaxation approach [4]

$$\widehat{\Theta}_\lambda \in \operatorname{argmin}_{\Theta \in C} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - \langle X_i, \Theta \rangle)^2 + \lambda \|\Theta\|_* \right\}, \quad (2.5)$$

where C is a convex feasible set of $\mathbb{R}^{d_1 \times d_2}$ and $\lambda > 0$ is a regularization parameter. When $C = \mathbb{R}^{d_1 \times d_2}$, $\widehat{\Theta}_\lambda$ is the matrix analog of the Lasso estimator for linear regression [37]. The statistical properties of $\widehat{\Theta}_\lambda$ in (2.5), mainly nonasymptotic deviation bounds under various matrix norms, have been studied in the literature when the additive noises ϵ_i are either Gaussian or sub-Gaussian. The performance of such a least-square-type estimator may break down quickly when the noise distribution is heavier-tailed. This is because outliers occur more frequently and the square loss is very sensitive to outliers. The impact of heavy-tailed errors on low-rank matrix recovery can be alleviated by replacing the ℓ_2 -loss with a more robust loss function, typified by the ℓ_1 -loss and the Huber loss [10]. When the error distribution is not only heavy-tailed but also asymmetric around zero, the use of ℓ_1 -loss or Huber loss with a fixed tuning parameter induces a bias that remains non-negligible as the number of measurements grows. For a better trade-off between robustness and bias, in the following we propose to use adaptive Huber loss [11,34] for robust low-rank matrix recovery, with a focus on the above three prototypical applications.

For matrix sensing and completion problems, i.e. applications (i) and (ii), we define the empirical loss function to be

$$\widehat{L}_\tau(\Theta) = \frac{1}{n} \sum_{i=1}^n \ell_\tau(y_i - \langle X_i, \Theta \rangle), \quad \Theta \in \mathbb{R}^{d_1 \times d_2}, \quad (2.6)$$

where $\ell_\tau(u) = \min\{u^2/2, \tau|u| - \tau^2/2\}$ denotes the adaptive Huber loss parameterized by $\tau = \tau_n > 0$, referred to as the robustification parameter in [34]. For any pre-specified convex subset C of $\mathbb{R}^{d_1 \times d_2}$, we consider the following nuclear norm penalized robust regression estimator

$$\widehat{\Theta}_{\tau,\lambda} \in \operatorname{argmin}_{\Theta \in C} \left\{ \widehat{L}_\tau(\Theta) + \lambda \|\Theta\|_* \right\}, \quad (2.7)$$

where $\tau = \tau_n > 0$ and $\lambda = \lambda_n > 0$ are the robustification and regularization parameters respectively.

For the multitask regression problem – Application (iii), recall that the vector-valued observations $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$ can be written as $\{(y_j, X_j)\}_{j=1}^N$ ($N = nd_2$) via (2.4) so model (2.1) can be used. The classical reduced-rank regression method is based on solving the rank-constrained problem [17]

$$\min_{\operatorname{rank}(\Theta) \leq r} \left\{ \sum_{i=1}^n \|\mathbf{y}_i - \Theta^T \mathbf{x}_i\|_2^2 = \sum_{j=1}^N (y_j - \langle X_j, \Theta \rangle)^2 \right\},$$

for which an analytic solution is available. To robustify this classical procedure, similarly to the formulation (2.7) one may naively apply the Huber loss to each residual $y_j - \langle X_j, \Theta \rangle$. This, however, is no longer plausible because X_j 's are now dependent random matrices. Moreover, since we do not impose independence on the entries of $\epsilon_i = (\epsilon_{i1}, \dots, \epsilon_{id_2})^T$, ϵ_j 's defined in (2.4) may also be highly correlated. We propose to replace the ℓ_2 -loss on $\|\mathbf{y}_i - \Theta^T \mathbf{x}_i\|_2$ with the Huber loss, leading to $\min_{\operatorname{rank}(\Theta) \leq r} \sum_{i=1}^n \ell_\tau(\|\mathbf{y}_i - \Theta^T \mathbf{x}_i\|_2)$, which is a highly nonconvex problem. Similarly to (2.7), we resort to convex relaxation and consider the following nuclear norm penalized estimator

$$\widehat{\Theta}_{\tau,\lambda} \in \operatorname{argmin}_{\Theta \in \mathbb{R}^{d_1 \times d_2}} \left\{ \frac{1}{n} \sum_{i=1}^n \ell_\tau(\|\mathbf{y}_i - \Theta^T \mathbf{x}_i\|_2) + \lambda \|\Theta\|_* \right\}. \quad (2.8)$$

In Section 3, we characterize the nonasymptotic statistical accuracy for the robust low-rank estimator $\widehat{\Theta}_{\tau,\lambda}$ defined in (2.7) and (2.8) when the noise variables only have bounded variances. The key is to seek a suitable choice of τ and λ jointly to trade off among robustness, bias and approximation errors.

2.2. Algorithms

To solve the nuclear norm penalized optimization programs (2.7) and (2.8), in this section we present a unified algorithm by extending the local adaptive majorize-minimization (LMM) principle proposed in [12] to matrix settings. Recall that the proposed nuclear norm penalized Huber regression estimators have a general form

$$\widehat{\Theta}_{\tau,\lambda} \in \operatorname{argmin}_{\Theta \in C} \{\widehat{L}_{\tau}(\Theta) + \lambda \|\Theta\|_*\},$$

where $\widehat{L}_{\tau}(\Theta)$ is the empirical loss in (2.6) or (2.8), and C is taken to be $\mathbb{R}^{d_1 \times d_2}$ or $\{\Theta \in \mathbb{R}^{d_1 \times d_2} : \|\Theta\|_{\infty} \leq \alpha_0\}$ for some $\alpha_0 > 0$. The main idea of the LMM principle is to construct an isotropic quadratic function that locally majorizes the objective function at each iteration. In the matrix setting, given the previous iterate $\Theta^{(k-1)}$ at the k -th iteration, define the quadratic function

$$F(\Theta; \Theta^{(k-1)}, \phi_k) = \widehat{L}_{\tau}(\Theta^{(k-1)}) + \langle \nabla \widehat{L}_{\tau}(\Theta^{(k-1)}), \Theta - \Theta^{(k-1)} \rangle + \frac{\phi_k}{2} \|\Theta - \Theta^{(k-1)}\|_{\mathbb{F}}^2,$$

satisfying $F(\Theta^{(k-1)}; \Theta^{(k-1)}, \phi_k) = \widehat{L}_{\tau}(\Theta^{(k-1)})$, where $\phi_k > 0$ is a quadratic parameter. Next, define the k -th iterate as

$$\Theta^{(k)} \in \operatorname{argmin}_{\Theta \in C} \{F(\Theta; \Theta^{(k-1)}, \phi_k) + \lambda \|\Theta\|_*\}. \quad (2.9)$$

The parameter ϕ_k needs to be sufficiently large so that $\widehat{L}_{\tau}(\Theta^{(k)}) \leq F(\Theta^{(k)}; \Theta^{(k-1)}, \phi_k)$, which further implies

$$\begin{aligned} \widehat{L}_{\tau}(\Theta^{(k)}) + \lambda \|\Theta^{(k)}\|_* &\leq F(\Theta^{(k)}; \Theta^{(k-1)}, \phi_k) + \lambda \|\Theta^{(k)}\|_* \\ &\leq F(\Theta^{(k-1)}; \Theta^{(k-1)}, \phi_k) + \lambda \|\Theta^{(k-1)}\|_* \\ &= \widehat{L}_{\tau}(\Theta^{(k-1)}) + \lambda \|\Theta^{(k-1)}\|_*, \end{aligned}$$

where the second inequality is due to the optimality of $\Theta^{(k)}$. This ensures the descent of the objective function at each iteration. To choose a sufficiently large ϕ_k , we start from a small value, say $\phi_0 = 0.01$, and inflate it by a factor $\gamma > 1$, say $\gamma = 2$, until the local majorization requirement $\widehat{L}_{\tau}(\Theta^{(k)}) \leq F(\Theta^{(k)}; \Theta^{(k-1)}, \phi_k)$ is met. Since $F(\Theta; \Theta^{(k-1)}, \phi_k) \geq \widehat{L}_{\tau}(\Theta)$ when ϕ_k is no less than the largest eigenvalue of $\nabla^2 \widehat{L}_{\tau}(\Theta^{(k-1)})$, the iteration will stop after sufficiently many steps. Repeat the above steps until convergence (e.g., $\|\Theta^{(k)} - \Theta^{(k-1)}\|_{\mathbb{F}} \leq \epsilon$ for a sufficiently small $\epsilon > 0$) or until the maximum number of iterations is reached.

The main benefit of minimizing a penalized isotropic quadratic objective function is that the minimizer often has a closed form. For any matrix $\Theta \in \mathbb{R}^{d_1 \times d_2}$ with rank r , consider the singular value decomposition (SVD) $\Theta = U \Sigma V^T$, where U and V are, respectively, $d_1 \times r$ and $d_2 \times r$ matrices with orthonormal columns, and $\Sigma = \operatorname{diag}(\{\sigma_i\}_{1 \leq i \leq r})$ is an $r \times r$ diagonal matrix with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$.

Algorithm 1 LAMM algorithm for regularized adaptive Huber trace regression

Algorithm: $\{\Theta^{(k)}\}_{k=1}^\infty \leftarrow \text{LAMM}(\lambda, \Theta^{(0)}, \phi_0, \gamma, \epsilon)$
Input: $\lambda, \Theta^{(0)}, \phi_0, \gamma, \epsilon$
1: **for** $k = 1, 2, \dots$ **until** $\|\Theta^{(k)} - \Theta^{(k-1)}\|_2 \leq \epsilon$ **do**
2: **repeat**
3: $\Theta^{(k)} \leftarrow S(\Theta^{(k-1)} - \phi_k^{-1} \nabla \widehat{L}_\tau(\Theta^{(k-1)}), \phi_k^{-1} \lambda)$
4: $\Theta^{(k)} \leftarrow \Pi_C(\Theta^{(k)})$
5: **if** $F(\Theta^{(k)}; \Theta^{(k-1)}, \phi_k) < \widehat{L}_\tau(\Theta^{(k)})$ **then**
6: $\phi_k \leftarrow \gamma \cdot \phi_k$
7: **end if**
8: **until** $F(\Theta^{(k)}; \phi_k, \Theta^{(k-1)}) \geq \widehat{L}_\tau(\Theta^{(k)})$
9: **return** $\Theta^{(k)}$
10: **end for**
Output: $\widehat{\Theta} = \Theta^{(T)}$

For $\lambda > 0$, define the soft-thresholding operator $S(\Theta, \lambda) = \mathbf{U} \cdot \text{diag}(\{\max(\sigma_i - \lambda, 0)\}_{1 \leq i \leq r}) \cdot \mathbf{V}^T$. By Theorem 2.1 in [3], $\Theta^{(k)}$ given in (2.9) with $C = \mathbb{R}^{d_1 \times d_2}$ admits the closed-form expression

$$\Theta^{(k)} = T_{\lambda, \phi_k}(\Theta^{(k-1)}) := S(\Theta^{(k-1)} - \phi_k^{-1} \nabla \widehat{L}_\tau(\Theta^{(k-1)}), \phi_k^{-1} \lambda).$$

For a general convex subset $C \subseteq \mathbb{R}^{d_1 \times d_2}$, we can update $\Theta^{(k)}$ as

$$\Theta^{(k)} = \Pi_C(T_{\lambda, \phi_k}(\Theta^{(k-1)})),$$

where Π_C denotes Euclidean projection onto the subspace C . When $C = \{\Theta \in \mathbb{R}^{d_1 \times d_2} : \|\Theta\|_\infty \leq \alpha_0\}$, for example, $\Pi_C(\Theta) = (\max\{\min(\Theta_{jk}, \alpha_0), -\alpha_0\})_{1 \leq j \leq d_1, 1 \leq k \leq d_2}$. We summarize the key steps in Algorithm 1.

As a unified algorithm, Algorithm 1 applies to all three problems considered in this paper, matrix sensing, matrix completion and multitask regression. In terms of complexity, at each iteration $\nabla \widehat{L}_\tau(\Theta^{(k-1)})$ and $T_{\lambda, \phi_k}(\Theta^{(k-1)})$ can be computed in $O(nd_1 d_2)$ and $O(d_1 d_2^2)$ operations (assuming $d_1 \geq d_2$), respectively [39]. On the other hand, [13] employed the contractive Peaceman-Rachford splitting method for matrix sensing and multitask regression, and an ADMM-based algorithm for matrix completion. In addition to the operations described above, each ADMM iterate also involves computing the inverse of $2\mathbb{X}^T \mathbb{X} / n + \mathbf{I}_{d_1 d_2}$, where \mathbb{X} is an $n \times d_1 d_2$ matrix whose i -th row is $\text{vec}(X_i)$, and \mathbf{I}_k denotes the $k \times k$ identity matrix. By applying the Sherman–Morrison–Woodbury formula, this step can be implemented in $O(\min\{n, d_1 d_2\}^3)$ operations. Still, the computational complexity and storage cost (per iteration) of ADMM are much higher than the LAMM algorithm in the context of matrix completion, especially for large-scale datasets.

3. Theoretical guarantees

In this section, we establish the finite-sample statistical properties of the robust estimator $\widehat{\Theta}_{\tau, \lambda}$ for matrix sensing, matrix completion and multitask regression. Throughout, the noise variables ϵ_i in (2.1) and ϵ_i in (2.3) are assumed to have bounded variance only, and we do not require independence between ϵ_i and X_i or ϵ_i and x_i .

3.1. Matrix sensing

In the case of matrix compressed sensing, we set $C = \mathbb{R}^{d_1 \times d_2}$ in (2.7), and impose the following assumptions.

- (A1) $\Theta^* \in \mathcal{B}_q(\rho)$ for some $0 \leq q \leq 1$ and $\rho > 0$.
- (A2) The random matrix $X_i \in \mathbb{R}^{d_1 \times d_2}$ satisfies (i) $\mathbb{E}X_i = \mathbf{0}$, and (ii) $\text{vec}(X_i) \in \mathbb{R}^{d_1 d_2}$ is sub-exponential, that is, there exists a constant $\nu_0 \geq 1$ such that for any $A \in \mathbb{R}^{d_1 \times d_2}$ and $u \geq 0$,

$$\mathbb{P}(|\text{vec}(A)^T \text{vec}(X_i)| \geq \nu_0 \|A\|_F \cdot u) \leq 2e^{-u}.$$

Moreover, there exists a constant $c_l > 0$ such that $\lambda_{\min}(\mathbb{E} \text{vec}(X_i) \text{vec}(X_i)^T) \geq c_l$.

- (A3) The noise variables ϵ_i are such that $\mathbb{E}(\epsilon_i | X_i) = 0$ and $\mathbb{E}(\epsilon_i^2 | X_i) \leq \sigma_0^2$ (almost surely) for some constant $\sigma_0 > 0$.

Remark 1. The parameter ν_0 is often referred to as the sub-exponential parameter. For various well-behaved distributions on $\mathbb{R}^{d_1 \times d_2}$, the associated sub-exponential parameters are independent of the dimensions d_1 and d_2 . As prototypical examples, the distributions listed below satisfy Condition (A2) with dimension-free parameters ν_0 and c_l .

- (i) (Multivariate normal) $\text{vec}(X_i)$ follow $s \mathcal{N}(\mathbf{0}, \Sigma)$ with a positive-definite $\Sigma \in \mathbb{R}^{(d_1 d_2) \times (d_1 d_2)}$.
- (ii) (Uniform distribution on the Euclidean sphere) X_i follows the uniform distribution on the sphere centered at the origin with radius $(d_1 d_2)^{1/2}$, namely, $\{X \in \mathbb{R}^{d_1 \times d_2} : \|X\|_F = (d_1 d_2)^{1/2}\}$.
- (iii) (Uniform distribution on the ℓ_1 -ball) X_i follows the uniform distribution on the ℓ_1 -norm ball centered at the origin with radius $r \asymp d_1 d_2$, that is, $\{X \in \mathbb{R}^{d_1 \times d_2} : \sum_{j=1}^{d_1} \sum_{k=1}^{d_2} |X_{jk}| \leq r\}$.

Here, we note that the multivariate distributions in (i) and (ii) are not only sub-exponential but also sub-Gaussian.

To derive the convergence rate of $\widehat{\Theta}_{\tau, \lambda}$ under either the Frobenius norm or the nuclear norm, we first define a probability event that concerns the local restricted strong convexity (RSC) of the empirical loss $\widehat{L}_\tau(\cdot)$. For $s, l > 0$, define the Frobenius norm ball and trace norm cone

$$\mathbb{B}(s) = \{\Delta \in \mathbb{R}^{d_1 \times d_2} : \|\Delta\|_F \leq s\} \text{ and } \mathbb{C}(l) = \{\Delta \in \mathbb{R}^{d_1 \times d_2} : \|\Delta\|_* \leq l \|\Delta\|_F\}, \quad (3.1)$$

respectively.

Definition 3.1 (Local restricted strong convexity). Given radius parameters $s, l > 0$ and a curvature parameter $\kappa > 0$, define the event

$$\mathcal{E}(s, l, \kappa) = \left\{ \inf_{\Theta \in \Theta^* + \mathbb{B}(s) \cap \mathbb{C}(l)} \frac{\langle \nabla \widehat{L}_\tau(\Theta) - \nabla \widehat{L}_\tau(\Theta^*), \Theta - \Theta^* \rangle}{\|\Theta - \Theta^*\|_F^2} \geq \kappa \right\}, \quad (3.2)$$

which concerns the local restricted strong convexity of the empirical loss function.

We first provide a deterministic result on the convergence rate of $\widehat{\Theta}_{\tau, \lambda}$: for any choice of λ such that $\|\nabla \widehat{L}_\tau(\Theta^*)\|_2 \leq \lambda/2$, and conditioned on event $\mathcal{E}(s, l, \kappa)$ with suitably chosen (s, l) , we are guaranteed that

$$\|\widehat{\Theta}_{\tau, \lambda} - \Theta^*\|_F \lesssim \sqrt{\rho}(\lambda/\kappa)^{1-q/2} \quad \text{and} \quad \|\widehat{\Theta}_{\tau, \lambda} - \Theta^*\|_* \lesssim \rho(\lambda/\kappa)^{1-q}.$$

Theorem 3.2. Assume Condition (A1) holds. Let (λ, s, l, κ) satisfy $\lambda \geq 2\|\nabla \widehat{L}_\tau(\Theta^*)\|_2$,

$$s \geq 9.15\sqrt{\rho}(\lambda/\kappa)^{1-q/2} \quad \text{and} \quad l \geq 6.1\sqrt{\rho}(\kappa/\lambda)^{q/2}. \quad (3.3)$$

Conditioned on the event $\mathcal{E}(s, l, \kappa)$, the error matrix $\widehat{\Delta} := \widehat{\Theta}_{\tau, \lambda} - \Theta^*$ satisfies

$$\|\widehat{\Delta}\|_F \leq 9.15\sqrt{\rho} \left(\frac{\lambda}{\kappa}\right)^{1-q/2} \quad \text{and} \quad \|\widehat{\Delta}\|_* \leq 56\rho \left(\frac{\lambda}{\kappa}\right)^{1-q}.$$

In the following two propositions, we first derive an upper bound of $\|\nabla \widehat{L}_\tau(\Theta^*)\|_2$, and then establish the local RSC property of the empirical Huber loss function $\widehat{L}_\tau(\cdot)$. Together, these results show that with properly chosen λ, τ that depend on (n, d, s, l) along with the distributional parameters in Conditions (A2) and (A3), the event $\{\lambda \geq 2\|\nabla \widehat{L}_\tau(\Theta^*)\|_2\} \cap \mathcal{E}(s, l, c_l/4)$ occurs with high probability.

Proposition 3.1. Assume Conditions (A2) and (A3) hold. For any $\sigma \geq \sigma_0$ and $z > 0$, the empirical Huber loss $\widehat{L}_\tau(\cdot)$ with $\tau = \sigma\sqrt{n/(3d+z)}$ satisfies

$$\|\nabla \widehat{L}_\tau(\Theta^*)\|_2 \leq 10v_0 \cdot \sigma \sqrt{\frac{3d+z}{n}} \quad (3.4)$$

with probability at least $1 - e^{-z}$, where $d = d_1 + d_2$.

Proposition 3.2. Assume Conditions (A2) and (A3) hold. For any $s, l > 0$ and $z > 0$, let τ and n satisfy

$$\tau \geq 4v_0\sqrt{(2\sigma_0^2 + 96v_0^2s^2)/c_l} \quad \text{and} \quad n \geq C_1(\tau/s)^2(l^2d + z), \quad (3.5)$$

where $d = d_1 + d_2$ and $C_1 > 0$ is a constant depending only on v_0 and c_l . Then, the local RSC event $\mathcal{E}(s, l, \kappa)$ with $\kappa = c_l/4$ occurs with probability at least $1 - e^{-z}$.

Combining these high probability bounds with Theorem 3.2 leads to the convergence rate of $\widehat{\Theta}_{\tau, \lambda}$, as stated in the following theorem.

Theorem 3.3. Assume Conditions (A1)–(A3) hold. For any $z > 0$, the robust (approximately) low-rank matrix estimator $\widehat{\Theta}_{\tau, \lambda}$ defined in (2.7) with $C = \mathbb{R}^{d_1 \times d_2}$, $\tau \asymp \sigma_0\sqrt{n/(d+z)}$ and $\lambda \asymp \sigma_0\sqrt{(d+z)/n}$ satisfies

$$\|\widehat{\Theta}_{\tau, \lambda} - \Theta^*\|_F \lesssim \sigma_0^{1-q/2} \sqrt{\rho} \left(\frac{d+z}{n}\right)^{1/2-q/4} \quad \text{and} \quad \|\widehat{\Theta}_{\tau, \lambda} - \Theta^*\|_* \lesssim \sigma_0^{1-q} \rho \left(\frac{d+z}{n}\right)^{(1-q)/2}$$

with probability at least $1 - 2e^{-z}$ as long as $n \gtrsim \max\{(\rho/\sigma_0^q)^{2/(2-q)}, 1\}(d+z)$, where $d = d_1 + d_2$.

Remark 2. In the exact low-rank case, i.e. $q = 0$ and $\rho = r = \text{rank}(\Theta^*)$, the results in Theorem 3.3 imply that with high probability (over both the random sensing matrices X_i and noise variables ϵ_i), the robust estimator $\widehat{\Theta}_{\tau, \lambda}$ satisfies with high probability that

$$\|\widehat{\Theta}_{\tau, \lambda} - \Theta^*\|_F \lesssim \sigma_0 \sqrt{\frac{rd}{n}} \quad \text{as long as} \quad n \gtrsim rd.$$

Within a constant independent of (n, r, d_1, d_2) and noise scale σ_0 , this upper bound matches the information-theoretic lower bound established by [5] when ϵ_i are i.i.d. $\mathcal{N}(0, \sigma_0^2)$. The robustness manifests in two aspects. First, the noise distribution is only required to have bounded variance as opposed to sub-Gaussian tails. Secondly, we assume the random vector $\text{vec}(\mathbf{X}_i)$ is sub-exponential, whereas $\text{vec}(\mathbf{X}_i)$ is often assumed to have i.i.d. Gaussian/sub-Gaussian entries in the literature.

Remark 3. For matrix compressed sensing, based on a shrinkage principle [13] also proposed a robust low-rank estimator, which achieves near-optimal rate under heavy-tailed noise distributions. Its recovery guarantees (see Theorem 3 therein), however, depend on more stringent assumptions as needed in Theorem 3.3. In addition to Conditions (A2) and (A3), [13] assumed further that (i) $\text{vec}(\mathbf{X}_i)$ is sub-Gaussian, and (ii) $\mathbb{E}|y_i|^{2k} \leq M_k$ for some $k > 1$. Under these conditions and in the exact low-rank case (for brevity), their truncate/shrinkage estimator, denoted by $\tilde{\Theta}$, satisfies with high probability the bound

$$\|\tilde{\Theta} - \Theta^*\|_F \leq M_k^{1/(2k)} \sqrt{\frac{rd}{n}} \text{ as long as } n \gtrsim rd.$$

The above convergence rate is sub-optimal in terms of its dependence on the noise scale σ_0 . As the noise scale decreases, $M_k^{1/(2k)}$ remains to be bounded away from zero because

$$M_k^{1/k} > \mathbb{E}y_i^2 = \mathbb{E}\langle \mathbf{X}_i, \Theta^* \rangle^2 + \mathbb{E}(\epsilon_i^2).$$

Remark 4. The sample size requirement in Theorem 3.3 becomes more stringent as σ_0 goes to 0 when $q \neq 0$. This is an artifact of the technical argument used in the proof of the theorem. A similar sample size requirement, characterized by its inverse proportionality to the moment of the response, can be found in Theorem 3 of [13]. To modify the sample size requirement, we can choose

$$\tau \asymp \max(\sigma_0, 1)\sqrt{n/(d+z)} \text{ and } \lambda \asymp \max(\sigma_0, 1)\sqrt{(d+z)/n}$$

for a given $z > 0$. This results in a revised sample size requirement of $n \gtrsim \max(\rho^{2/(2-q)}, 1)(d+z)$, accompanied by an error bound

$$\|\hat{\Theta}_{\tau, \lambda} - \Theta^*\|_F \lesssim \max(\sigma_0^{1-q/2}, 1)\sqrt{\rho} \left(\frac{d+z}{n} \right)^{1/2-q/4}$$

with probability at least $1 - 2e^{-z}$. This error bound is similar to the deviation bound in Corollary 5 of [28], but it should be noted that they are not proportional to the noise level. Additionally, Theorem 3.5 in Section 3.3 also relies on a similar technical argument, necessitating a larger sample size as σ_0 approaches zero.

3.2. Matrix completion

This subsection investigates matrix completion under the following assumptions.

- (B1) $\Theta^* \in \mathcal{B}_q(\rho)$ and $\|\Theta^*\|_\infty \leq \alpha_0$ for some $\alpha_0 > 0$. We thus set $C = \{\Theta \in \mathbb{R}^{d_1 \times d_2} : \|\Theta\|_\infty \leq \alpha_0\}$ in (2.7) so that $\Theta^* \in C$.
- (B2) $\mathbf{X}_i \in \mathbb{R}^{d_1 \times d_2}$ is uniformly sampled from $\{\mathbf{e}_j(d_1)\mathbf{e}_k^T(d_2)\}_{1 \leq j \leq d_1, 1 \leq k \leq d_2}$, where $\{\mathbf{e}_j(d)\}_{j=1}^d$ are the canonical basis vectors in \mathbb{R}^d . Specifically, $\mathbb{P}\{\mathbf{X}_i = \mathbf{e}_j(d_1)\mathbf{e}_k^T(d_2)\} = (d_1 d_2)^{-1}$.
- (B3) $\mathbb{E}(\epsilon_i | \mathbf{X}_i) = 0$ and $\mathbb{E}(\epsilon_i^2 | \mathbf{X}_i) \leq \sigma_0^2$ (almost surely) for some constant $\sigma_0 > 0$.

Remark 5. In addition to the assumption that Θ^* is of (approximately) low-rank, we require in Condition (B1) that $\|\Theta^*\|_\infty \leq \alpha_0$ for some $\alpha_0 > 0$. Past works on noisy matrix completion also imposed the same or similar conditions. For instance, [18] and [27] assumed that $\|\Theta^*\|_\infty$ is bounded; [29] and [13] required the spikiness ratio $\|\Theta^*\|_\infty / \|\Theta^*\|_F$ to be bounded; [4] and [6] relied on matrix incoherence conditions. Without such extra conditions, the number of measurements should satisfy $n \asymp d_1 d_2$ in order to recover Θ^* in the worst case; see [6] and [29] for details.

Similarly to the matrix sensing case, the key steps to establish the convergence rate of $\widehat{\Theta}_{\tau, \lambda}$ are (i) an upper bound of $\|\nabla \widehat{L}_\tau(\Theta^*)\|_2$ as shown in Proposition 3.3 below, and (ii) a lower bound for

$$\langle \nabla \widehat{L}_\tau(\Theta) - \nabla \widehat{L}_\tau(\Theta^*), \Theta - \Theta^* \rangle$$

uniformly over Θ in a neighborhood of Θ^* . The sparsity of X_i in this case (see Condition (B2)) introduces more subtleties into the analysis of the latter, as we will see from Proposition 3.4.

Proposition 3.3. Assume Conditions (B2) and (B3) hold. For any $\sigma \geq \sigma_0$ and $z > 0$, the loss function $\widehat{L}_\tau(\cdot)$ with $\tau = \sigma \sqrt{n / \{d_2(z + \log d)\}}$ satisfies with probability at least $1 - e^{-z}$ that

$$\|\nabla \widehat{L}_\tau(\Theta^*)\|_2 \leq (3\sigma_0 + 2\sigma/3) \sqrt{\frac{z + \log d}{d_2 n}}, \quad (3.6)$$

where $d = d_1 + d_2$.

Proposition 3.4. Assume Conditions (B2) and (B3) hold. For any $s, l > 0$ and $z > 0$, let τ and n satisfy

$$\tau^2 \geq 16 \max[ns^2 / \{l^2 d_1^2 d_2(z + \log d)\}, \sigma_0^2] \quad \text{and} \quad n \geq d_2 \log d,$$

where $d = d_1 + d_2$. Define the constrain set

$$\mathbb{A}(s, l) = \left\{ \Delta \in \mathbb{B}(s) \cap \mathbb{C}(l) : \frac{\|\Delta\|_\infty^2}{\|\Delta\|_F^2 / (d_1 d_2)} \leq \frac{1}{8} \sqrt{\frac{n}{z + \log d}} \right\}, \quad (3.7)$$

where $\mathbb{B}(s)$ and $\mathbb{C}(l)$ are given in (3.1). Then, for all $\Theta \in \mathbb{R}^{d_1 \times d_2}$ with $\Delta := \Theta - \Theta^* \in \mathbb{A}(s, l)$, we have with probability at least $1 - e^{-z}$ that

$$\langle \nabla \widehat{L}_\tau(\Theta) - \nabla \widehat{L}_\tau(\Theta^*), \Theta - \Theta^* \rangle \geq \frac{1}{4d_1 d_2} \|\Delta\|_F^2 - C_0 l^2 \frac{d_1(z + \log d)}{n} \|\Delta\|_\infty^2,$$

where $C_0 > 1$ is an absolute constant.

With the above preparations, we now state the statistical guarantees for matrix completion under heavy-tailed noise.

Theorem 3.4. Assume Conditions (B1)–(B3) hold. For any $z > 0$, set

$$\tau \asymp \sigma \sqrt{\frac{n}{d_2(z + \log d)}} \quad \text{and} \quad \lambda \asymp \sigma \sqrt{\frac{z + \log d}{d_2 n}},$$

where $\sigma = \max\{\sigma_0, \alpha_0\}$ and $d = d_1 + d_2$. Then, the robust (approximately) low-rank matrix estimator $\widehat{\Theta}_{\tau, \lambda}$ defined in (2.7) with $C = \{\Theta \in \mathbb{R}^{d_1 \times d_2} : \|\Theta\|_\infty \leq \alpha_0\}$ satisfies

$$\frac{1}{d_1 d_2} \|\widehat{\Theta}_{\tau, \lambda} - \Theta^*\|_F^2 \lesssim \max \left[\sigma^{2-q} \rho \left(\frac{d_1(z + \log d)}{n} \right)^{1-q/2}, \alpha_0^2 \sqrt{\frac{z + \log d}{n}} \right] \quad (3.8)$$

and

$$\frac{1}{\sqrt{d_1 d_2}} \|\widehat{\Theta}_{\tau, \lambda} - \Theta^*\|_* \lesssim \max \left[\sigma^{1-q} \rho \left(\frac{d_1(z + \log d)}{n} \right)^{\frac{1-q}{2}}, \alpha_0^{\frac{2(1-q)}{2-q}} \frac{\rho^{\frac{1}{2-q}}}{(d_1 d_2)^{\frac{q}{2(2-q)}}} \left(\frac{z + \log d}{n} \right)^{\frac{1-q}{2(2-q)}} \right]$$

with probability at least $1 - 2e^{-z}$ whenever $n \gtrsim d_2(z + \log d)$.

Remark 6. In the context of matrix completion, one is interested in recovering a large low-rank data matrix from a highly incomplete subset of its entries. A natural assumption is $n \leq d_1 d_2$, which in turn implies $\sqrt{\log(d)/n} \leq d_1 \log(d)/n$, where $d = d_1 + d_2$ and $d_1 \geq d_2$. Therefore, taking $z = \log n$, the maximum in (3.8) is often given by its first term. In the exact low-rank case, the general results in Theorem 3.4 imply that the proposed robust estimator $\widehat{\Theta}_{\tau, \lambda}$ with $\tau \asymp \sigma_0 \sqrt{n/(d_2 \log d)}$ and $\lambda \asymp \sigma_0 \sqrt{\log(d)/(d_2 n)}$ satisfies the bound

$$\frac{1}{d_1 d_2} \|\widehat{\Theta}_{\tau, \lambda} - \Theta^*\|_F^2 \lesssim \max\{\alpha_0^2, \sigma_0^2\} \frac{r d_1 \log d}{n} \quad (3.9)$$

with high probability as long as $n \gtrsim d_2 \log d$. Under our notations, Theorem 6 in [20] shows that when $\epsilon_i \sim N(0, \sigma_0^2)$ is independent of X_i , there exist absolute constants $\beta \in (0, 1)$ and $c > 0$ such that

$$\inf_{\widehat{\Theta}} \sup_{\text{rank}(\Theta^*) \leq r, \|\Theta^*\|_\infty \leq \alpha_0} \mathbb{P} \left\{ \frac{1}{d_1 d_2} \|\widehat{\Theta} - \Theta^*\|_F^2 > c \min(\sigma_0^2, \alpha_0^2) \frac{r d_1}{n} \right\} \geq \beta,$$

where $\inf_{\widehat{\Theta}}$ is the infimum over all estimators $\widehat{\Theta} \in \mathbb{R}^{d_1 \times d_2}$. Therefore, the rate derived in Theorem 3.4 is minimax optimal up to a logarithmic factor and a trailing term.

Remark 7 (Comparison to existing work on robust (noisy) matrix completion). In the context of matrix completion with heavy-tailed noise, several robust estimators have been proposed and studied. [27] proposed a two-step method that computes a truncation-type matrix estimator, denoted by $\widetilde{\Theta}$, in step one and then solves the nuclear norm penalized optimization $\|\Theta - \widetilde{\Theta}\|_F^2/(d_1 d_2) + \lambda \|\Theta\|_*$. In the exact low-rank case, this two-step estimator satisfies a high probability bound, which is similar to (3.8) with $q = 0$, when ϵ_i is independent of X_i and has bounded variance. The independence assumption can be removed by slightly modifying the proof in [27]. For matrix sensing and multitask regression, it is unclear whether such a two-step procedure will also lead to robust estimates that satisfy sharp error bounds proportional only to the noise scale. Concurrently, [13] considered a similar two-step estimator, but their theoretical result requires a slightly stronger moment condition, i.e. $\mathbb{E}\{\mathbb{E}(\epsilon_i^2 | X_i)^k\} \leq M_k$ for some $k > 1$. Our proposal is more relevant to [10], who also used the Huber loss for matrix completion in the presence of heavy-tailed errors. Their results, however, depend on stronger assumptions on the error distribution. In addition to Conditions (B1) and (B2), they further assumed that (i) the distribution of ϵ_i is symmetric around 0, and (ii) there exists a constant $C_1 > 0$ such that the cumulative distribution function $F(\cdot)$ of ϵ_i satisfies

$$F(x + \tau) - F(x - \tau) \geq 1/C_1^2 \quad \text{for all } |x| \leq 2\alpha_0 \text{ and } \tau \leq 2\alpha_0.$$

Under these conditions and in the exact low-rank case, [10] proved that the nuclear norm penalized Huber regression estimator, denoted by $\check{\Theta}$, satisfies

$$\frac{1}{d_1 d_2} \|\check{\Theta} - \Theta^*\|_F^2 = O_{\mathbb{P}} \left\{ \max(\alpha_0^2, \tau^2) C_1^4 \frac{r d_1 \log(d_1 + d_2)}{n} \right\}$$

under the sample size requirement $n \gtrsim d_2 \log(d_2) \log(d_1 + d_2)$.

3.3. Multitask regression

In this section, we establish the statistical properties of the robust low-rank multitask (reduced-rank) regression estimator $\hat{\Theta}_{\tau, \lambda}$ (2.8). With slight abuse of notation, we write

$$\hat{L}_{\tau}(\Theta) = \frac{1}{n} \sum_{i=1}^n \ell_{\tau}(\|y_i - \Theta^T x_i\|_2), \quad \Theta \in \mathbb{R}^{d_1 \times d_2}, \quad (3.10)$$

where $\tau > 0$ is the robustification parameter.

- (C1) $\Theta^* \in \mathcal{B}_q(\rho)$ for some $0 \leq q \leq 1$ and $\rho > 0$.
- (C2) $x_i \in \mathbb{R}^{d_1}$ are i.i.d. zero-mean sub-Gaussian vectors, that is, there exists a (dimension-free) constant $v_0 \geq 1$ such that

$$\mathbb{E} e^{u^T x_i} \leq e^{v_0^2 \|u\|_2^2 / 2}, \quad \text{valid for any } u \in \mathbb{R}^{d_1}.$$

Moreover, there exists a constant $c_l > 0$ such that $\lambda_{\min}(\mathbb{E} x_i x_i^T) \geq c_l$.

- (C3) The noise vectors $\epsilon_i \in \mathbb{R}^{d_2}$ are such that $\mathbb{E}(\epsilon_i | X_i) = \mathbf{0}$ and $\lambda_{\max}(\mathbb{E}(\epsilon_i \epsilon_i^T | X_i)) \leq \sigma_0^2$ (almost surely).

Proposition 3.5. Assume Conditions (C2) and (C3) hold. For any $\sigma \geq \sigma_0$ and $z > 0$, choose $\tau = \sigma \sqrt{n/(z + \log d)}$ with $d = d_1 + d_2$. Then, it holds with probability at least $1 - e^{-z}$ that

$$\|\nabla \hat{L}_{\tau}(\Theta^*)\|_2 \leq C v_0 \sigma \sqrt{\frac{d(z + \log d)}{n}},$$

where $C > 0$ is a universal constant.

Proposition 3.6. Assume Conditions (C2) and (C3) hold, and let $s, \tau > 0$ and $z > 0$ satisfy

$$\tau \geq \max \left\{ 4\sigma_0 \sqrt{d_2}, 2v_0 s \sqrt{2d_1 + 3z + 3 \log n} \right\}. \quad (3.11)$$

Then, with probability at least $1 - 2e^{-z}$,

$$\langle \nabla \hat{L}_{\tau}(\Theta) - \nabla \hat{L}_{\tau}(\Theta^*), \Theta - \Theta^* \rangle \geq \frac{c_l}{2} \|\Theta - \Theta^*\|_F^2 \quad \text{for all } \Theta \in \Theta^* + \mathbb{B}(s),$$

provided that $n \gtrsim v_0^4 c_l^{-2} (d_1 + z)$.

Theorem 3.5. Assume Conditions (C1)-(C3) hold. For any $z > 0$, the robust (approximately) low-rank matrix estimator $\widehat{\Theta}_{\tau,\lambda}$ defined in (2.8) with $\tau \asymp \sigma_0 \sqrt{n/(z + \log d)}$ and $\lambda \asymp \sigma_0 \sqrt{d(z + \log d)/n}$ ($d = d_1 + d_2$) satisfies the bounds

$$\|\widehat{\Theta}_{\tau,\lambda} - \Theta^*\|_F \lesssim \sigma_0^{1-q/2} \sqrt{\rho} \left\{ \frac{d(z + \log d)}{n} \right\}^{\frac{1}{2}-\frac{q}{4}} \quad \text{and} \quad \|\widehat{\Theta}_{\tau,\lambda} - \Theta^*\|_* \lesssim \sigma_0^{1-q} \rho \left\{ \frac{d(z + \log d)}{n} \right\}^{\frac{1-q}{2}}$$

with probability at least $1 - 3e^{-z}$ as long as

$$n \gtrsim \max\{(\rho/\sigma_0^q)^{2/(4-q)}, 1\} \cdot (d + z + \log n)(z + \log d).$$

Remark 8. Again, in the exact low-rank case where $\rho = r = \text{rank}(\Theta^*)$ and $q = 0$, Theorem 3.5 shows that for an arbitrary accuracy $\varepsilon > 0$, we have $\|\widehat{\Theta}_{\tau,\lambda} - \Theta^*\|_F \leq \varepsilon$ with an overwhelming probability provided that the number of measurements satisfies

$$n \gtrsim \sigma_0^2 \frac{rd \log d}{\varepsilon^2}. \quad (3.12)$$

This result improves Theorem 5 in [13] in several aspects. Under the multitask regression model (2.3), they assumed that

$$\lambda_{\max}(\mathbb{E}(\mathbf{y}_i \mathbf{y}_i^T)) \leq R < \infty, \quad \max_{1 \leq i \leq n, 1 \leq k \leq d_2} \mathbb{E}\{\mathbb{E}(\epsilon_{ik}^2 | \mathbf{x}_i)^k\} \leq M_k < \infty \quad \text{for some } k > 1,$$

and for each i , $\epsilon_{i1}, \dots, \epsilon_{id_2}$ are pairwise (conditionally) independent given \mathbf{x}_i . In contrast, Condition (C3) only assumes bounded variances and allows arbitrary dependency between ϵ_{ik} 's. For an arbitrary accuracy $\varepsilon > 0$, the truncated/shrinkage matrix estimator $\widetilde{\Theta}$ proposed by [13] satisfies $\|\widetilde{\Theta} - \Theta^*\|_F \leq \varepsilon$ with high probability provided

$$n \gtrsim (R + M_k^{1/k}) \frac{rd \log d}{\varepsilon^2}. \quad (3.13)$$

Here the term $R + M_k^{1/k}$ can be much larger than σ_0^2 in (3.12). More importantly, as the noise scale σ_0 decays, R stays away from zero because

$$R \geq \lambda_{\max}(\mathbb{E} \mathbf{y}_i \mathbf{y}_i^T) = \lambda_{\max}((\Theta^*)^T \Sigma \Theta^* + \mathbb{E} \boldsymbol{\varepsilon}_i \boldsymbol{\varepsilon}_i^T) \geq \lambda_{\max}((\Theta^*)^T \Sigma \Theta^*).$$

4. Numerical studies

4.1. Finite-sample performance

In this section, we perform simulation studies to assess the finite-sample performance of the nuclear norm penalized adaptive Huber trace regression method (Nuclear-AH) in all three problems. As a benchmark, we implement the nuclear norm penalized least squares (Nuclear-LS) estimator also via the LAMM algorithm.

In addition to the regularization parameter λ , the use of an adaptive Huber loss also involves a robustification parameter τ that changes with data scales. We set $\tau = c_\tau \cdot a_{n,d}$ and $\lambda = c_\lambda \cdot b_{n,d}$, where c_τ and c_λ are positive constants that are independent of (n, d) but depend on the noise scale, and $a_{n,d}$ and $b_{n,d}$ are determined by the theoretical results in Section 3, as follows:

Matrix Sensing	Normal	t	Pareto
Nuclear-LS	0.227 (0.010)	0.173 (0.052)	0.169 (0.093)
Nuclear-AH	0.227 (0.010)	0.132 (0.008)	0.107 (0.007)
Matrix Completion	Normal	t	Pareto
Nuclear-LS	0.424 (0.021)	0.280 (0.047)	0.315 (0.041)
Nuclear-AH	0.445 (0.022)	0.223 (0.023)	0.252 (0.022)
Multitask Regression	Normal	t	Pareto
Nuclear-LS	0.228 (0.005)	0.213 (0.112)	0.237 (0.181)
Nuclear-AH	0.228 (0.005)	0.148 (0.003)	0.120 (0.003)

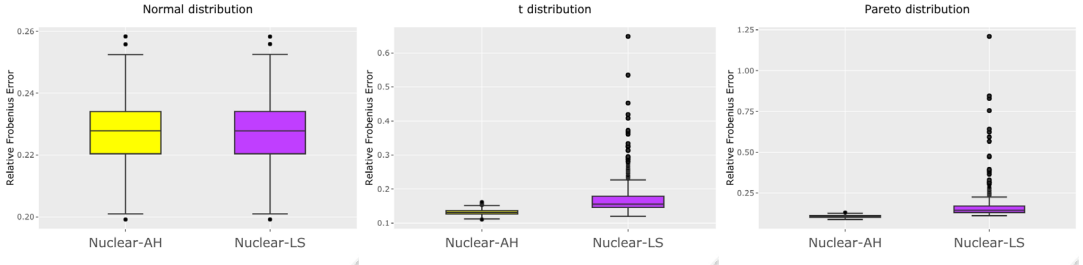
Table 1. Mean relative Frobenius error $\|\widehat{\Theta} - \Theta^*\|_F / \|\Theta^*\|_F$ (with standard deviations in parentheses), averaged over 500 replications, under the matrix sensing, matrix completion and multitask regression settings.

- (a) For matrix sensing, we choose $a_{n,d} = \sqrt{n/d}$ and $b_{n,d} = \sqrt{d/n}$.
- (b) For matrix completion, we choose $a_{n,d} = \sqrt{n/(d \log d)}$ and $b_{n,d} = \sqrt{\log(d)/(dn)}$.
- (c) For multitask learning, we choose $a_{n,d} = \sqrt{n/\log d}$ and $b_{n,d} = \sqrt{d \log(d)/n}$.

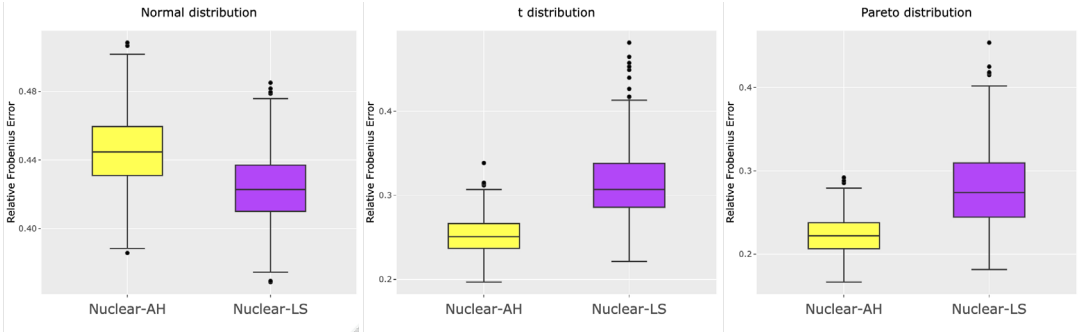
Then we follow the following steps to tune c_τ and c_λ .

- (i) First, choose the constant c_λ in the Nuclear-LS method via five-fold CV with the absolute median loss as the criterion. In particular, we use the “one-standard-error” rule, which yields the most parsimonious model within one standard error of the minimum CV error.
- (ii) Next, let $\{r_i\}_{i=1}^n$ be the Nuclear-LS residuals with c_λ selected via CV as in Step (i). As a rule-of-thumb, we set c_τ as the median absolute deviation of $\{r_i\}$, i.e. $\text{median}\{|r_i - \text{median}(r_i)|\}/0.6745$.
- (iii) With c_τ determined after Step (ii), we choose the constant c_λ in the Nuclear-AH method again via five-fold CV under the one-standard-error rule.

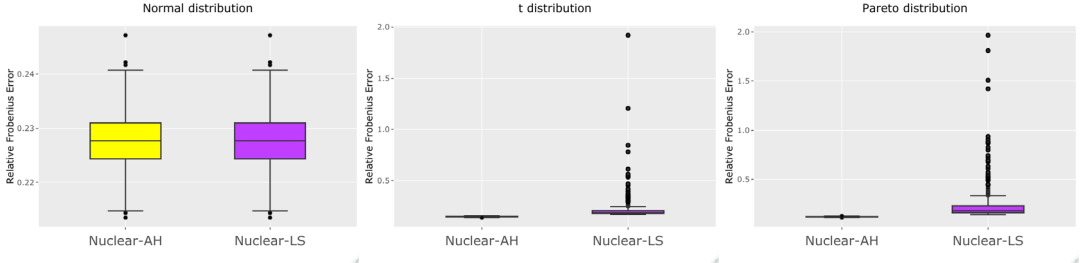
Under the matrix sensing and matrix completion setups, the data $\{(y_i, X_i)\}_{i=1}^n$ are generated from $y_i = \langle X_i, \Theta^* \rangle + \epsilon_i$, where ϵ_i follows one of the following three distributions: (i) $\mathcal{N}(0, 0.5^2)$ —centered normal distribution with standard deviation 0.5 (lighted-tailed and symmetric), (ii) $t_{2.1}/8$ —scaled t -distribution with 2.1 degrees of freedom (heavy-tailed and symmetric), and (iii) $\text{Par}(2, 1)/8$ —scaled Pareto distribution with scale parameter 1 and shape parameter 2 (heavy-tailed and asymmetric). For matrix sensing, we set $(d_1, d_2, n) = (50, 50, 1500)$, $\Theta^* \in \mathbb{R}^{d_1 \times d_2}$ is such that $\text{rank}(\Theta^*) = 5$ and all nonzero singular values of Θ^* are 1, and the design matrix X_i consists of i.i.d. standard normal entries. For matrix completion, we set $(d_1, d_2, n) = (50, 50, 2000)$, $\Theta^* \in \mathbb{R}^{d_1 \times d_2}$ is such that $\|\Theta^*\|_F = \sqrt{d_1 d_2}$ and $\text{rank}(\Theta^*) = 5$, and X_i is uniformly sampled from $\{e_j(d_1)e_k^T(d_2)\}_{1 \leq j \leq d_1, 1 \leq k \leq d_2}$. To implement LAMM, we use the initial estimates $\Theta^{(0)} = \mathbf{0}$ and $\Theta^{(0)} = (d_1 d_2 / n) \sum_{i=1}^n y_i X_i$, respectively, under the two setups. In the case of multitask regression, the data vectors $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$ are generated from $y_i = (\Theta^*)^T \mathbf{x}_i + \epsilon_i$, where $\Theta^* \in \mathbb{R}^{d_1 \times d_2}$ is the same as in the matrix sensing setting, $\mathbf{x}_i \in \mathbb{R}^{d_1}$ are i.i.d. standard normal and



(a) Matrix sensing setting with $(d_1, d_2, n) = (50, 50, 1500)$



(b) Matrix completion setting with $(d_1, d_2, n) = (50, 50, 2000)$



(c) Multitask regression setting with $(d_1, d_2, n) = (80, 80, 2000)$

Figure 1. Boxplots of relative Frobenius errors $\|\hat{\Theta} - \Theta^*\|_F / \|\Theta^*\|_F$ (based on 500 repetitions) for the Nuclear-AH and Nuclear-LS estimators under the matrix sensing, matrix completion and multitask regression settings.

$\epsilon_i = (\epsilon_{i1}, \dots, \epsilon_{id_2})^T$ consists of i.i.d. entries following one of the above three errors distributions. In this case we set $d_1 = d_2 = 80$ and $n = 2000$.

Simulation results on the relative Frobenius error $\|\hat{\Theta} - \Theta^*\|_F / \|\Theta^*\|_F$, averaged over 500 repetitions, are presented in Table 1. To better demonstrate the robustness property of Nuclear-AH, Figure 1 shows the boxplots of (relative) Frobenius errors for the cross-validated Nuclear-LS and Nuclear-AH estimators under three error distributions. We see that Nuclear-LS and Nuclear-AH have almost identical performance when errors have symmetric and light tails, while the latter achieves considerably better performance under all three settings in the presence of heavy-tailed and/or asymmetric errors.

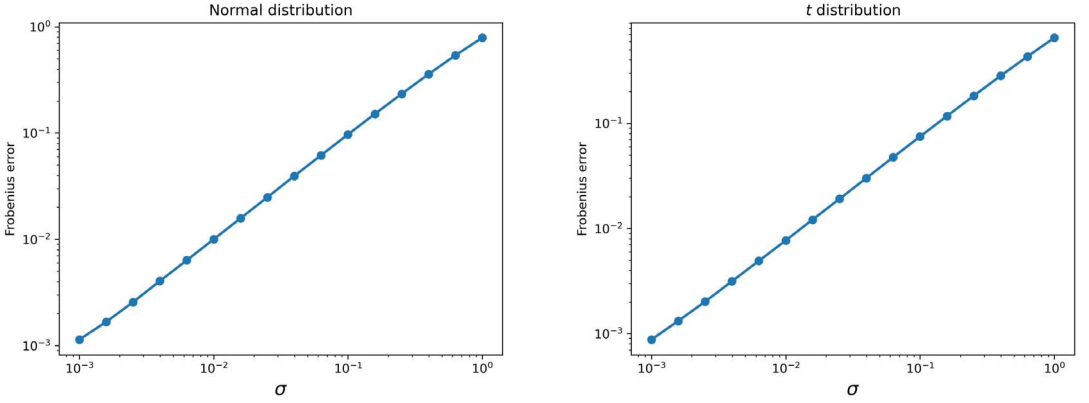
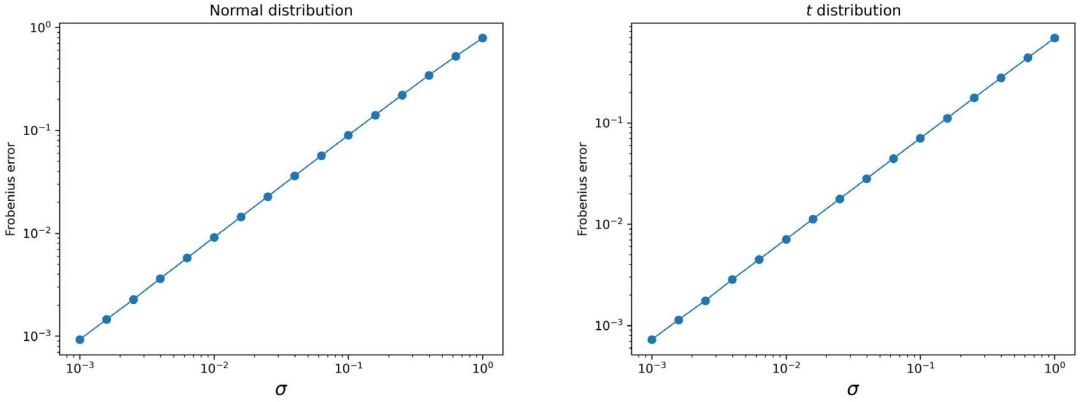
(a) Matrix sensing setting with $(d_1, d_2, n) = (50, 50, 2000)$ (b) Multitask regression setting with $(d_1, d_2, n) = (50, 50, 2000)$

Figure 2. Plots of Frobenius error $\|\hat{\Theta} - \Theta^*\|_F$ versus noise scale based on 200 simulations under the matrix sensing and multitask regression settings.

4.2. Convergence rate versus noise scale

In this section, we numerically examine the dependence of $\|\hat{\Theta} - \Theta^*\|_F$ on the noise scale under the matrix sensing and multitask regression settings. Our theoretical results, Theorem 3.3 and Theorem 3.5, indicate that in the exact low-rank case, $\|\hat{\Theta} - \Theta^*\|_F$ should be proportional to the noise scale σ , where $\sigma^2 = \mathbb{E}(\epsilon_i^2)$ or $\sigma^2 = \lambda_{\max}(\mathbb{E}(\epsilon_i \epsilon_i^T))$. To verify this, given a sequence of σ values ranging from 10^{-3} to 1, we generate ϵ_i from either $\sigma \cdot \mathcal{N}(0, 1)$ or $\sigma \cdot t_{2,1}/16$. The specifications of Θ^* and X_i or x_i are the same as in Section 4.1.

Under the matrix sensing setting, we set $(d_1, d_2, n) = (50, 50, 2000)$ and choose $\tau = 2\sigma\sqrt{n/d}$ and $\lambda = \sigma\sqrt{d/n}$ with $d = d_1 + d_2$. For multitask regression, we set $(d_1, d_2, n) = (100, 100, 3000)$ and choose $\tau = \sigma\sqrt{n/\log d}$ and $\lambda = 0.5\sigma\sqrt{d \log(d)/n}$. Figure 2 shows the plots of the Frobenius error versus noise scale, based on 200 replications, under these two settings and two error distributions. Consistent with the predictions of Theorems 3.3 and 3.5, we observe a nearly perfect linear growth in all four plots.

Acknowledgments

We thank the Editor, an Associate Editor, and two anonymous reviewers for their constructive comments and valuable suggestions, which have significantly helped us improve the quality of this work.

Funding

MY and WZ are supported in part by the NSF Grant DMS-2113409. QS is partially supported by Natural Sciences and Engineering Research Council of Canada (Grant RGPIN-2018-06484), a New Frontiers in Research Fund NFRFE-2019-00603, and a Data Sciences Institute Catalyst Grant.

Supplementary Material

Supplement to “Low-rank matrix recovery under heavy-tailed errors” (DOI: [10.3150/23-BEJ1675SUPP](https://doi.org/10.3150/23-BEJ1675SUPP); .pdf). The supplementary material [42] contains the proofs of the theoretical results established in Section 3.

References

- [1] Argyriou, A., Evgeniou, T. and Pontil, M. (2008). Convex multi-task feature learning. *Mach. Learn.* **73** 243–272.
- [2] Burer, S. and Monteiro, R.D.C. (2003). A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization. *Math. Program.* **95** 329–357. [MR1976484 https://doi.org/10.1007/s10107-002-0352-8](https://doi.org/10.1007/s10107-002-0352-8)
- [3] Cai, J.-F., Candès, E.J. and Shen, Z. (2010). A singular value thresholding algorithm for matrix completion. *SIAM J. Optim.* **20** 1956–1982. [MR2600248 https://doi.org/10.1137/080738970](https://doi.org/10.1137/080738970)
- [4] Candès, E.J. and Plan, Y. (2009). Matrix completion with noise. *Proc. IEEE* **98** 925–936.
- [5] Candès, E.J. and Plan, Y. (2011). Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Trans. Inf. Theory* **57** 2342–2359. [MR2809094 https://doi.org/10.1109/TIT.2011.2111771](https://doi.org/10.1109/TIT.2011.2111771)
- [6] Candès, E.J. and Recht, B. (2009). Exact matrix completion via convex optimization. *Found. Comput. Math.* **9** 717–772. [MR2565240 https://doi.org/10.1007/s10208-009-9045-5](https://doi.org/10.1007/s10208-009-9045-5)
- [7] Chen, J., Liu, D. and Li, X. (2020). Nonconvex rectangular matrix completion via gradient descent without $\ell_{2,\infty}$ regularization. *IEEE Trans. Inf. Theory* **66** 5806–5841. [MR4158648 https://doi.org/10.1109/TIT.2020.2992234](https://doi.org/10.1109/TIT.2020.2992234)
- [8] Chen, Y., Jalali, A., Sanghavi, S. and Caramanis, C. (2013). Low-rank matrix recovery from errors and erasures. *IEEE Trans. Inf. Theory* **59** 4324–4337.
- [9] Dalalyan, A. and Thompson, P. (2019). Outlier-robust estimation of a sparse linear model using ℓ_1 -penalized Huber’s M -estimator. In *Advances in Neural Information Processing Systems* **32** 13188–13198.
- [10] Elsener, A. and van de Geer, S. (2018). Robust low-rank matrix estimation. *Ann. Statist.* **46** 3481–3509. [MR3852659 https://doi.org/10.1214/17-AOS1666](https://doi.org/10.1214/17-AOS1666)
- [11] Fan, J., Li, Q. and Wang, Y. (2017). Estimation of high dimensional mean regression in the absence of symmetry and light tail assumptions. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **79** 247–265. [MR3597972 https://doi.org/10.1111/rssb.12166](https://doi.org/10.1111/rssb.12166)
- [12] Fan, J., Liu, H., Sun, Q. and Zhang, T. (2018). I-LAMM for sparse learning: Simultaneous control of algorithmic complexity and statistical error. *Ann. Statist.* **46** 814–841. [MR3782385 https://doi.org/10.1214/17-AOS1568](https://doi.org/10.1214/17-AOS1568)
- [13] Fan, J., Wang, W. and Zhu, Z. (2021). A shrinkage principle for heavy-tailed data: High-dimensional robust low-rank matrix recovery. *Ann. Statist.* **49** 1239–1266. [MR4298863 https://doi.org/10.1214/20-aos1980](https://doi.org/10.1214/20-aos1980)

- [14] Goldberg, D., Nichols, D., Oki, B.M. and Terry, D. (1992). Using collaborative filtering to weave an information tapestry. *Commun. ACM* **35** 61–70.
- [15] Gross, D., Liu, Y.-K., Flammia, S.T., Becker, S. and Eisert, J. (2010). Quantum state tomography via compressed sensing. *Phys. Rev. Lett.* **105** 150401. <https://doi.org/10.1103/PhysRevLett.105.150401>
- [16] Huber, P.J. (1973). Robust regression: Asymptotics, conjectures and Monte Carlo. *Ann. Statist.* **1** 799–821. [MR0356373](https://doi.org/10.2307/2346178)
- [17] Izenman, A.J. (1975). Reduced-rank regression for the multivariate linear model. *J. Multivariate Anal.* **5** 248–264. [MR0373179 https://doi.org/10.1016/0047-259X\(75\)90042-1](https://doi.org/10.1016/0047-259X(75)90042-1)
- [18] Klopp, O. (2014). Noisy low-rank matrix completion with general sampling distribution. *Bernoulli* **20** 282–303. [MR3160583 https://doi.org/10.3150/12-BEJ486](https://doi.org/10.3150/12-BEJ486)
- [19] Klopp, O., Lounici, K. and Tsybakov, A.B. (2017). Robust matrix completion. *Probab. Theory Related Fields* **169** 523–564. [MR3704775 https://doi.org/10.1007/s00440-016-0736-y](https://doi.org/10.1007/s00440-016-0736-y)
- [20] Koltchinskii, V., Lounici, K. and Tsybakov, A.B. (2011). Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *Ann. Statist.* **39** 2302–2329. [MR2906869 https://doi.org/10.1214/11-AOS894](https://doi.org/10.1214/11-AOS894)
- [21] Li, X. (2013). Compressed sensing and matrix completion with constant proportion of corruptions. *Constr. Approx.* **37** 73–99. [MR3010211 https://doi.org/10.1007/s00365-012-9176-9](https://doi.org/10.1007/s00365-012-9176-9)
- [22] Li, Y., Ma, T. and Zhang, H. (2018). Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference on Learning Theory* 2–47.
- [23] Lounici, K., Pontil, M., van de Geer, S. and Tsybakov, A.B. (2011). Oracle inequalities and optimal inference under group sparsity. *Ann. Statist.* **39** 2164–2204. [MR2893865 https://doi.org/10.1214/11-AOS896](https://doi.org/10.1214/11-AOS896)
- [24] Luan, X., Fang, B., Liu, L., Yang, W. and Qian, J. (2014). Extracting sparse error of robust PCA for face recognition in the presence of varying illumination and occlusion. *Pattern Recognit.* **47** 495–508.
- [25] Ma, C., Wang, K., Chi, Y. and Chen, Y. (2020). Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution. *Found. Comput. Math.* **20** 451–632. [MR4099988 https://doi.org/10.1007/s10208-019-09429-9](https://doi.org/10.1007/s10208-019-09429-9)
- [26] Ma, J. and Fattahi, S. (2023). Global convergence of sub-gradient method for robust matrix recovery: Small initialization, noisy measurements, and over-parameterization. *J. Mach. Learn. Res.* **24** Paper No. [96], 84. [MR4582518](https://doi.org/10.48550/jmlr.2023.24.96)
- [27] Minsker, S. (2018). Sub-Gaussian estimators of the mean of a random matrix with heavy-tailed entries. *Ann. Statist.* **46** 2871–2903. [MR3851758 https://doi.org/10.1214/17-AOS1642](https://doi.org/10.1214/17-AOS1642)
- [28] Negahban, S. and Wainwright, M.J. (2011). Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *Ann. Statist.* **39** 1069–1097. [MR2816348 https://doi.org/10.1214/10-AOS850](https://doi.org/10.1214/10-AOS850)
- [29] Negahban, S. and Wainwright, M.J. (2012). Restricted strong convexity and weighted matrix completion: Optimal bounds with noise. *J. Mach. Learn. Res.* **13** 1665–1697. [MR2930649](https://doi.org/10.48550/jmlr.2012.13.1665)
- [30] Recht, B., Fazel, M. and Parrilo, P.A. (2010). Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Rev.* **52** 471–501. [MR2680543 https://doi.org/10.1137/070697835](https://doi.org/10.1137/070697835)
- [31] Rohde, A. and Tsybakov, A.B. (2011). Estimation of high-dimensional low-rank matrices. *Ann. Statist.* **39** 887–930. [MR2816342 https://doi.org/10.1214/10-AOS860](https://doi.org/10.1214/10-AOS860)
- [32] She, Y. and Chen, K. (2017). Robust reduced-rank regression. *Biometrika* **104** 633–647. [MR3694587 https://doi.org/10.1093/biomet/asx032](https://doi.org/10.1093/biomet/asx032)
- [33] Shen, Y., Li, J., Cai, J. and Xia, D. (2022). Computationally efficient and statistically optimal robust low-rank matrix estimation. [arXiv:2203.00953](https://arxiv.org/abs/2203.00953).
- [34] Sun, Q., Zhou, W.-X. and Fan, J. (2020). Adaptive Huber regression. *J. Amer. Statist. Assoc.* **115** 254–265. [MR4078461 https://doi.org/10.1080/01621459.2018.1543124](https://doi.org/10.1080/01621459.2018.1543124)
- [35] Tan, K.M., Sun, Q. and Witten, D. (2022). Sparse reduced rank Huber regression in high dimensions. *J. Amer. Statist. Assoc.* To appear.
- [36] Thompson, P. (2020). Outlier-robust sparse/low-rank least-squares regression and robust matrix completion. [arXiv:2012.06750](https://arxiv.org/abs/2012.06750).
- [37] Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B* **58** 267–288. [MR1379242](https://doi.org/10.2307/2346178)

- [38] Tong, T., Ma, C. and Chi, Y. (2021). Accelerating ill-conditioned low-rank matrix estimation via scaled gradient descent. *J. Mach. Learn. Res.* **22** Paper No. 150, 63. [MR4318506](#) <https://doi.org/10.1080/15502287.2020.1856971>
- [39] Trefethen, L.N. and Bau, D. III (1997). *Numerical Linear Algebra*. Philadelphia, PA: SIAM. [MR1444820](#) <https://doi.org/10.1137/1.9780898719574>
- [40] Wang, B. and Fan, J. (2022). Robust matrix completion with heavy-tailed noise. [arXiv:2206.04276](#).
- [41] Wei, K., Cai, J.-F., Chan, T.F. and Leung, S. (2016). Guarantees of Riemannian optimization for low rank matrix recovery. *SIAM J. Matrix Anal. Appl.* **37** 1198–1222. [MR3543156](#) <https://doi.org/10.1137/15M1050525>
- [42] Yu, M., Sun, Q. and Zhou, W.-X. (2024). Supplement to “Low-rank matrix recovery under heavy-tailed errors.” <https://doi.org/10.3150/23-BEJ1675SUPP>
- [43] Zhang, J., Fattahi, S. and Zhang, R. (2021). Preconditioned gradient descent for over-parameterized nonconvex matrix factorization. In *Advances in Neural Information Processing Systems* **34** 5985–5996.

Received November 2022 and revised September 2023