

# BRIDGING CONVEX AND NONCONVEX OPTIMIZATION IN ROBUST PCA: NOISE, OUTLIERS AND MISSING DATA

BY YUXIN CHEN<sup>1</sup>, JIANQING FAN<sup>2,\*</sup>, CONG MA<sup>3</sup> AND YULING YAN<sup>2,†</sup>

<sup>1</sup>Department of Electrical and Computer Engineering, Princeton University, [yuxin.chen@princeton.edu](mailto:yuxin.chen@princeton.edu)

<sup>2</sup>Department of Operations Research and Financial Engineering, Princeton University, [\\*jqfan@princeton.edu](mailto:*jqfan@princeton.edu);  
[†yulingy@princeton.edu](mailto:†yulingy@princeton.edu)

<sup>3</sup>Department of Statistics, University of Chicago, [congma@uchicago.edu](mailto:congma@uchicago.edu)

This paper delivers improved theoretical guarantees for the convex programming approach in low-rank matrix estimation, in the presence of (1) random noise, (2) gross sparse outliers and (3) missing data. This problem, often dubbed as *robust principal component analysis (robust PCA)*, finds applications in various domains. Despite the wide applicability of convex relaxation, the available statistical support (particularly the stability analysis in the presence of random noise) remains highly suboptimal, which we strengthen in this paper. When the unknown matrix is well conditioned, incoherent and of constant rank, we demonstrate that a principled convex program achieves near-optimal statistical accuracy, in terms of both the Euclidean loss and the  $\ell_\infty$  loss. All of this happens even when nearly a constant fraction of observations are corrupted by outliers with arbitrary magnitudes. The key analysis idea lies in bridging the convex program in use and an auxiliary nonconvex optimization algorithm, and hence the title of this paper.

**1. Introduction.** A diverse array of science and engineering applications (e.g., video surveillance, joint shape matching, graph clustering, covariance modeling, graphical models) involves estimation of low-rank matrices (Candès et al. (2011), Chandrasekaran, Parrilo and Willsky (2012), Chen, Guibas and Huang (2014), Chen et al. (2014), Chi, Lu and Chen (2019), Davenport and Romberg (2016), Fan, Liao and Mincheva (2013)). The imperfectness of data acquisition processes, however, presents several common yet critical challenges: (1) random noise: which reflects the uncertainty of the environment and/or the measurement processes; (2) outliers: which represent a sort of corruption that exhibits abnormal behavior and (3) incomplete data, namely, we might only get to observe a fraction of the matrix entries. Low-rank matrix estimation algorithms aimed at addressing these challenges have been extensively studied under the umbrella of *robust principal component analysis (Robust PCA)* (Candès et al. (2011), Chandrasekaran et al. (2011)), a terminology popularized by the seminal work (Candès et al. (2011)).

To formulate the above mentioned problem in a more precise manner, imagine that we seek to estimate an unknown low-rank matrix  $L^* \in \mathbb{R}^{n_1 \times n_2}$ . What we can obtain is a collection of partially observed and corrupted entries as follows:

$$(1.1) \quad M_{ij} = L_{ij}^* + S_{ij}^* + E_{ij}, \quad (i, j) \in \Omega_{\text{obs}},$$

where  $S^* = [S_{ij}^*]$  is a matrix consisting of outliers,  $E = [E_{ij}]$  represents the random noise, and we only observe entries over an index subset  $\Omega_{\text{obs}} \subseteq [n_1] \times [n_2]$  with  $[n] := \{1, 2, \dots, n\}$ . The current paper assumes that  $S^*$  is a relatively sparse matrix whose nonzero entries might have arbitrary magnitudes. This assumption has been commonly adopted in prior work to

---

Received January 2020; revised September 2020.

*MSC2020 subject classifications.* Primary 62F10; secondary 62B10.

*Key words and phrases.* Robust principal component analysis, convex relaxation,  $\ell_\infty$  guarantees, leave-one-out analysis.

model gross outliers, while enabling reliable disentanglement of the outlier component and the low-rank component (Candès et al. (2011), Chandrasekaran et al. (2011), Chen et al. (2013), Li (2013)). In addition, we suppose that the entries  $\{E_{ij}\}$  are independent zero-mean sub-Gaussian random variables, as commonly assumed in the statistics literature to model a large family of random noise. The aim is to reliably estimate  $\mathbf{L}^*$  given the grossly corrupted and possibly incomplete data (1.1). Ideally, this task should be accomplished without knowing the locations and magnitudes of the outliers  $\mathbf{S}^*$ .

1.1. *A principled convex programming approach.* Focusing on the noiseless case with  $\mathbf{E} = \mathbf{0}$ , the papers by Candès et al. (2011), Chandrasekaran et al. (2011) delivered a positive and somewhat surprising message: both the low-rank component  $\mathbf{L}^*$  and the sparse component  $\mathbf{S}^*$  can be efficiently recovered with absolutely no error by means of a principled convex program

$$(1.2) \quad \underset{\mathbf{L}, \mathbf{S} \in \mathbb{R}^{n_1 \times n_2}}{\text{minimize}} \quad \|\mathbf{L}\|_* + \tau \|\mathbf{S}\|_1 \quad \text{subject to} \quad \mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{L} + \mathbf{S} - \mathbf{M}) = \mathbf{0},$$

provided that certain “separation” and “incoherence” conditions on  $(\mathbf{L}^*, \mathbf{S}^*, \Omega_{\text{obs}})$  hold<sup>1</sup> and that the regularization parameter  $\tau$  is properly chosen. Here,  $\|\mathbf{L}\|_*$  denotes the nuclear norm (i.e., the sum of the singular values) of  $\mathbf{L}$ ,  $\|\mathbf{S}\|_1 = \sum_{i,j} |S_{ij}|$  denotes the usual entrywise  $\ell_1$  norm, and  $\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{M})$  denotes the Euclidean projection of a matrix  $\mathbf{M}$  onto the subspace of matrices supported on  $\Omega_{\text{obs}}$ . Given that the nuclear norm  $\|\cdot\|_*$  (resp., the  $\ell_1$  norm  $\|\cdot\|_1$ ) is the convex relaxation of the rank function  $\text{rank}(\cdot)$  (resp., the  $\ell_0$  counting norm  $\|\cdot\|_0$ ), the rationale behind (1.2) is rather clear: it seeks a decomposition  $(\mathbf{L}, \mathbf{S})$  of  $\mathbf{M}$  by promoting the low-rank structure of  $\mathbf{L}$  as well as the sparsity structure of  $\mathbf{S}$ .

Moving on to the more realistic noisy setting, a natural strategy is to solve the following regularized least-squares problem:

$$(1.3) \quad \underset{\mathbf{L}, \mathbf{S} \in \mathbb{R}^{n_1 \times n_2}}{\text{minimize}} \quad \frac{1}{2} \|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{L} + \mathbf{S} - \mathbf{M})\|_{\text{F}}^2 + \lambda \|\mathbf{L}\|_* + \tau \|\mathbf{S}\|_1.$$

With the regularization parameters  $\lambda, \tau > 0$  properly chosen, one hopes to strike a balance between enhancing the goodness of fit (by enforcing  $\mathbf{L} + \mathbf{S} - \mathbf{M}$  to be small) and promoting the desired low-complexity structures (by regularizing both the nuclear norm of  $\mathbf{L}$  and the  $\ell_1$  norm of  $\mathbf{S}$ ). A natural and important question comes into our mind:

*Where does the algorithm (1.3) stand in terms of its statistical performance in the presence of random noise, sparse outliers and missing data?*

Unfortunately, however simple this program (1.3) might seem, the existing theoretical support remains far from satisfactory, as we shall discuss momentarily.

1.2. *Theory-practice gaps under random noise.* To assess the tightness of prior statistical guarantees for (1.3), we find it convenient to first look at a simple setting where (i)  $n_1 = n_2 = n$ , (ii)  $\mathbf{E}$  consists of independent Gaussian components, namely,  $E_{ij} \sim \mathcal{N}(0, \sigma^2)$  and (iii) there is no missing data. This simple scenario is sufficient to illustrate the suboptimality of prior theory.

*Prior statistical guarantees.* The paper Zhou et al. (2010) was the first to derive a sort of statistical performance guarantees for the above convex program. Under mild conditions,

<sup>1</sup>Clearly, if the low-rank matrix  $\mathbf{L}^*$  is also sparse, one cannot possibly separate  $\mathbf{S}^*$  from  $\mathbf{L}^*$ . The same holds true if the matrix  $\mathbf{S}^*$  is simultaneously sparse and low-rank.

Zhou et al. (2010) demonstrated that any minimizer  $(\widehat{\mathbf{L}}, \widehat{\mathbf{S}})$  of (1.3) achieves<sup>2</sup>

$$(1.4) \quad \|\widehat{\mathbf{L}} - \mathbf{L}^*\|_{\text{F}} = O(n\|\mathbf{E}\|_{\text{F}}) = O(\sigma n^2)$$

with high probability, where we have substituted in the well-known high-probability bound  $\|\mathbf{E}\|_{\text{F}} = O(\sigma n)$  under i.i.d. Gaussian noise. While this theory corroborates the potential stability of convex relaxation against both additive noise and sparse outliers, it remains unclear whether the estimation error bound (1.4) reflects the true performance of the convex program in use. In what follows, we shall compare it with an oracle error bound and collect some numerical evidence.

*Comparisons with an oracle bound.* Consider an idealistic scenario where an oracle informs us of the outlier matrix  $\mathbf{S}^*$ . With the assistance of this oracle, the task of estimating  $\mathbf{L}^*$  reduces to a low-rank matrix denoising problem (Donoho and Gavish (2014)). By fixing  $\mathbf{S}$  to be  $\mathbf{S}^*$  in (1.3), we arrive at a simplified convex program

$$(1.5) \quad \underset{\mathbf{L} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{L} - (\mathbf{L}^* + \mathbf{E})\|_{\text{F}}^2 + \lambda \|\mathbf{L}\|_{*}.$$

It is known that (e.g., Chen et al. (2020a), Donoho and Gavish (2014)): under mild conditions and with a properly chosen  $\lambda$ , the estimation error of (1.5) satisfies

$$(1.6) \quad \|\widehat{\mathbf{L}} - \mathbf{L}^*\|_{\text{F}} = O(\sigma \sqrt{nr}),$$

where we abuse the notation and denote by  $\widehat{\mathbf{L}}$  the minimizer of (1.5). The large gap between the above two bounds (1.4) and (1.6) is self-evident; in particular, if  $r = O(1)$ , the gap between these two bounds can be as large as an order of  $n^{1.5}$ .

*A numerical example without oracles.* One might naturally wonder whether the discrepancy between the two bounds (1.4) and (1.6) stems from the magical oracle information (i.e.,  $\mathbf{S}^*$ ) which (1.3) does not have the luxury to know. To demonstrate that this is not the case, we conduct some numerical experiments to assess the importance of such oracle information. Generate  $\mathbf{L}^* = \mathbf{X}^* \mathbf{Y}^{*\top}$ , where  $\mathbf{X}^*, \mathbf{Y}^* \in \mathbb{R}^{n \times r}$  are random orthonormal matrices. Each entry of  $\mathbf{S}^*$  is generated independently from a mixed distribution: with probability 1/10, the entry is drawn from  $\mathcal{N}(0, 10)$ ; otherwise, it is set to be zero. In other words, approximately 10% of the entries in  $\mathbf{L}^*$  are corrupted by large outliers. Throughout the experiments, we set  $\lambda = 5\sigma \sqrt{n}$  and  $\tau = 2\sigma \sqrt{\log n}$  with  $\sigma$  the standard deviation of each noise entry  $\{E_{ij}\}$ . Figure 1(a) fixes  $r = 5$ ,  $\sigma = 10^{-3}$  and examines the dependency of the Euclidean error  $\|\widehat{\mathbf{L}} - \mathbf{L}^*\|_{\text{F}}$  on the size  $\sqrt{n}$ . Similarly, Figure 1(b) fixes  $r = 5$ ,  $n = 1000$  and displays the estimation error  $\|\widehat{\mathbf{L}} - \mathbf{L}^*\|_{\text{F}}$  as the noise size  $\sigma$  varies in a log-log plot. As can be seen from Figure 1, the performance of the oracle-aided estimator (1.5) matches the theoretical prediction (1.6), namely, the numerical estimation error  $\|\widehat{\mathbf{L}} - \mathbf{L}^*\|_{\text{F}}$  is proportional to both  $\sqrt{n}$  and  $\sigma$ . Perhaps more intriguingly, even without the help of the oracle, the original estimator (1.3) performs quite well and behaves qualitatively similarly. In comparison with the bound (1.4) derived in the prior work (Zhou et al. (2010)), our numerical experiments suggest that the convex estimator (1.3) might perform much better than previously predicted.

All in all, there seems to be a large gap between the practical performance of (1.3) and the existing theoretical support. This calls for a new theory that better explains practice, which we pursue in the current paper. We remark in passing that statistical guarantees have been developed in Agarwal, Negahban and Wainwright (2012), Klopp, Lounici and Tsybakov (2017) for other convex estimators (i.e., the ones that are different from the convex estimator (1.3) considered herein). We shall compare our results with theirs later in Section 1.4.

<sup>2</sup>Mathematically, Zhou et al. (2010) investigated an equivalent constrained form of (1.3) and developed an upper bound on the corresponding estimation error.

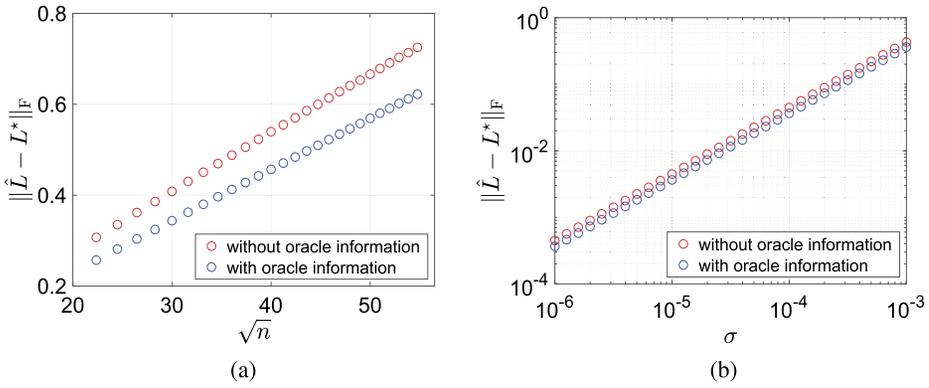


FIG. 1. (a) Euclidean estimation errors of (1.3) and (1.5) versus the problem size  $\sqrt{n}$ , where we fix  $r = 5$ ,  $\sigma = 10^{-3}$ ; (b) Euclidean estimation errors of (1.3) and (1.5) versus the noise level  $\sigma$  in a log-log plot, where we fix  $n = 1000$ ,  $r = 5$ . For both plots, we take  $\lambda = 5\sigma\sqrt{n}$  and  $\tau = 2\sigma\sqrt{\log n}$ . The results are averaged over 50 independent trials.

1.3. *Models, assumptions and notation.* As it turns out, the appealing empirical performance of the convex program (1.3) in the presence of both sparse outliers and zero-mean random noise can be justified in theory. Toward this end, we need to introduce several notations and model assumptions that will be used throughout. Let  $U^*\Sigma^*V^{*\top}$  be the singular value decomposition (SVD) of the unknown rank- $r$  matrix  $L^* \in \mathbb{R}^{n_1 \times n_2}$ , where  $U^* \in \mathbb{R}^{n_1 \times r}$  and  $V^* \in \mathbb{R}^{n_2 \times r}$  consist of orthonormal columns and  $\Sigma^* = \text{diag}\{\sigma_1^*, \dots, \sigma_r^*\}$  is a diagonal matrix. Here, we let

$$\sigma_{\max} := \sigma_1^* \geq \sigma_2^* \geq \dots \geq \sigma_r^* =: \sigma_{\min} \quad \text{and} \quad \kappa := \sigma_{\max}/\sigma_{\min}$$

represent the singular values and the condition number of  $L^*$ , respectively. We denote by  $\Omega^*$  the support set of  $S^*$ , that is,

$$(1.7) \quad \Omega^* := \{(i, j) \in \Omega_{\text{obs}} : S_{ij}^* \neq 0\}.$$

With this set of notation in place, we list below our key model assumptions.

ASSUMPTION 1.1 (Incoherence). The low-rank matrix  $L^*$  with SVD  $L^* = U^*\Sigma^*V^{*\top}$  is assumed to be  $\mu$ -incoherent in the sense that

$$(1.8) \quad \|U^*\|_{2,\infty} \leq \sqrt{\frac{\mu}{n_1}} \|U^*\|_F = \sqrt{\frac{\mu r}{n_1}} \quad \text{and} \quad \|V^*\|_{2,\infty} \leq \sqrt{\frac{\mu}{n_2}} \|V^*\|_F = \sqrt{\frac{\mu r}{n_2}}.$$

Here,  $\|U\|_{2,\infty}$  denotes the largest  $\ell_2$  norm of all rows of a matrix  $U$ .

ASSUMPTION 1.2 (Random sampling). Each entry is observed independently with probability  $p$ , namely,

$$(1.9) \quad \mathbb{P}\{(i, j) \in \Omega_{\text{obs}}\} = p.$$

ASSUMPTION 1.3 (Random locations of outliers). Each observed entry is independently corrupted by an outlier with probability  $\rho_s$ , namely,

$$(1.10) \quad \mathbb{P}\{(i, j) \in \Omega^* \mid (i, j) \in \Omega_{\text{obs}}\} = \rho_s,$$

where  $\Omega^* \subseteq \Omega_{\text{obs}}$  is the support of the outlier matrix  $S^*$ .

ASSUMPTION 1.4 (Random signs of outliers). The signs of the nonzero entries of  $\mathbf{S}^*$  are i.i.d. symmetric Bernoulli random variables (independent from the locations), namely,

$$(1.11) \quad \text{sign}(S_{ij}^*) \stackrel{\text{ind.}}{=} \begin{cases} 1, & \text{with probability } 1/2, \\ -1, & \text{else,} \end{cases} \quad \text{for all } (i, j) \in \Omega^*.$$

ASSUMPTION 1.5 (Random noise). The noise matrix  $\mathbf{E} = [E_{ij}]$  is composed of independent symmetric<sup>3</sup> zero-mean sub-Gaussian random variables with sub-Gaussian norm at most  $\sigma > 0$ , that is,  $\|E_{ij}\|_{\psi_2} \leq \sigma$  (see Vershynin (2012), Definition 5.7, for precise definitions).

We take a moment to expand on our model assumptions. Assumption 1.1 is standard in the low-rank matrix recovery literature (Candès and Recht (2009), Candès et al. (2011), Chen (2015), Chi, Lu and Chen (2019)). If  $\mu$  is small, then this assumption specifies that the singular spaces of  $\mathbf{L}^*$  is not sparse in the standard basis, thus ensuring that  $\mathbf{L}^*$  is not simultaneously low-rank and sparse. Assumption 1.3 requires the sparsity pattern of the outliers  $\mathbf{S}^*$  to be random, which precludes it from being simultaneously sparse and low-rank. In essence, Assumptions 1.1 and 1.3 are identifiability conditions, taken together as a sort of separation condition on  $(\mathbf{L}^*, \mathbf{S}^*)$ , which plays a crucial role in guaranteeing exact recovery in the noiseless case (i.e.,  $\mathbf{E} = \mathbf{0}$ ); see Candès et al. (2011) for more discussions on these conditions. Assumption 1.4 requires the signs of the outliers to be random, which has also been made in Wong and Lee (2017), Zhou et al. (2010).<sup>4</sup> We shall discuss in detail the crucial role of this random sign assumption (as opposed to deterministic sign patterns) in Section 1.6.

Finally, we introduce several notation to be used throughout. Denote by  $f(n) \lesssim g(n)$  or  $f(n) = O(g(n))$  the condition  $|f(n)| \leq Cg(n)$  for some constant  $C > 0$  when  $n$  is sufficiently large; we use  $f(n) \gtrsim g(n)$  to denote  $f(n) \geq C|g(n)|$  for some constant  $C > 0$  when  $n$  is sufficiently large; we also use  $f(n) \asymp g(n)$  to indicate that  $f(n) \lesssim g(n)$  and  $f(n) \gtrsim g(n)$  hold simultaneously. The notation  $f(n) \gg g(n)$  (resp.,  $f(n) \ll g(n)$ ) means that there exists a sufficiently large (resp., small) constant  $c_1 > 0$  (resp.,  $c_2 > 0$ ) such that  $f(n) \geq c_1g(n)$  (resp.,  $f(n) \leq c_2g(n)$ ). For any subspace  $T$ , we denote by  $\mathcal{P}_T(\mathbf{M})$  the Euclidean projection of a matrix  $\mathbf{M}$  onto the subspace  $T$ , and let  $\mathcal{P}_{T^\perp}(\mathbf{M}) := \mathbf{M} - \mathcal{P}_T(\mathbf{M})$ . For any index set  $\Omega$ , we denote by  $\mathcal{P}_\Omega(\mathbf{M})$  the Euclidean projection of a matrix  $\mathbf{M}$  onto the subspace of matrices supported on  $\Omega$ , and define  $\mathcal{P}_{\Omega^c}(\mathbf{M}) := \mathbf{M} - \mathcal{P}_\Omega(\mathbf{M})$ . For any matrix  $\mathbf{M}$ , we let  $\|\mathbf{M}\|$ ,  $\|\mathbf{M}\|_F$ ,  $\|\mathbf{M}\|_*$ ,  $\|\mathbf{M}\|_1$  and  $\|\mathbf{M}\|_\infty$  denote its spectral norm, Frobenius norm, nuclear norm, entrywise  $\ell_1$  norm and entrywise  $\ell_\infty$  norm, respectively.

1.4. *Main results.* Armed with the above model assumptions, we are positioned to present our improved statistical guarantees for convex relaxation (1.3) in the random noise setting. Without loss of generality, assume that

$$n_1 \geq n_2.$$

As we shall elucidate in Section 1.5 and Section 3, our theory is established by exploiting an intriguing and intimate connection between convex relaxation and nonconvex optimization, and hence the title of this paper.

For the sake of simplicity, we shall start by presenting our statistical guarantees when the rank  $r$ , the condition number  $\kappa$  and the incoherence parameter  $\mu$  of  $\mathbf{L}^*$  are all bounded by

<sup>3</sup>In fact, we only require  $E_{ij}$  to be symmetric for all  $(i, j) \in \Omega^*$ .

<sup>4</sup>Note that while the theorems in Wong and Lee (2017), Zhou et al. (2010) do not make explicit this random sign assumption, the proofs therein do rely on this assumption to guarantee the existence of certain approximate dual certificates.

some constants. Despite its simplicity, this setting subsumes as special cases a wide array of fundamentally important applications, including angular and phase synchronization (Singer (2011)) in computational biology, joint shape mapping problem (Chen, Guibas and Huang (2014), Huang and Guibas (2013)) in computer vision, and so on. All of these problems involve estimating a very well-conditioned matrix  $\mathbf{L}^*$  with a small rank.

**THEOREM 1.6.** *Suppose that Assumptions 1.1–1.5 hold, and that  $r, \kappa, \mu = O(1)$ . Take  $\lambda = C_\lambda \sigma \sqrt{n_1 p}$  and  $\tau = C_\tau \sigma \sqrt{\log n_2}$  in (1.3) for some large enough constants  $C_\lambda, C_\tau > 0$ . Assume that*

$$(1.12) \quad n_1 n_2 p \geq C_{\text{sample}} n_1 \log^6 n_1, \quad \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n_1}{p}} \leq \frac{c_{\text{noise}}}{\sqrt{\log n_1}} \quad \text{and} \quad \rho_s \leq \frac{c_{\text{outlier}}}{\log n_1}$$

for some sufficiently large constant  $C_{\text{sample}} > 0$  and some sufficiently small constants  $c_{\text{noise}}, c_{\text{outlier}} > 0$ . Then with probability exceeding  $1 - O(n_2^{-3})$ , the following hold:

1. Any minimizer  $(\mathbf{L}_{\text{cvx}}, \mathbf{S}_{\text{cvx}})$  of the convex program (1.3) obeys

$$(1.13a) \quad \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_{\text{F}} \leq C_{\text{err}} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n_1}{p}} \|\mathbf{L}^*\|_{\text{F}},$$

$$(1.13b) \quad \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_{\infty} \leq C_{\text{err}} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n_1 \log n_1}{p}} \|\mathbf{L}^*\|_{\infty},$$

$$(1.13c) \quad \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\| \leq C_{\text{err}} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n_1}{p}} \|\mathbf{L}^*\|$$

for some constant  $C_{\text{err}} > 0$ .

2. Letting  $\mathbf{L}_{\text{cvx},r} := \arg \min_{\mathbf{L}: \text{rank}(\mathbf{L}) \leq r} \|\mathbf{L} - \mathbf{L}_{\text{cvx}}\|_{\text{F}}$  be the best rank- $r$  approximation of  $\mathbf{L}_{\text{cvx}}$ , we have

$$(1.14) \quad \|\mathbf{L}_{\text{cvx},r} - \mathbf{L}_{\text{cvx}}\|_{\text{F}} \leq \frac{1}{n_2^5} \cdot \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n_1}{p}} \|\mathbf{L}^*\|_{\text{F}},$$

and the statistical guarantees (1.13) hold unchanged if  $\mathbf{L}_{\text{cvx}}$  is replaced by  $\mathbf{L}_{\text{cvx},r}$ .

Before we embark on interpreting our statistical guarantees, let us first parse the required conditions (1.12) in Theorem 1.6. For simplicity, we assume that  $n_1 = n_2 = n$ .

- *Missing data.* Theorem 1.6 accommodates the case where a dominant fraction of entries are unobserved (more precisely, the sample size can be as low as an order of  $n \text{poly} \log n$ ). This is an appealing result since, even when there is no noise and no outlier (i.e.,  $\mathbf{E} = \mathbf{0}$  and  $\rho_s = 0$ ), the minimal sample size required for exact matrix completion is at least on the order of  $n \log n$  (Candès and Tao (2010)). In comparison, prior theory on robust PCA with both sparse outliers and dense additive noise is either based on full observations (Agarwal, Negahban and Wainwright (2012), Zhou et al. (2010)), or assumes the sampling rate  $p$  exceeds some universal constant (Wong and Lee (2017)). In other words, these prior results require the number of observed entries to exceed the order of  $n^2$ . The only exception is Klopp, Lounici and Tsybakov (2017), which also allows a significant amount of missing data (i.e.,  $p \gtrsim (\text{poly} \log n)/n$ ).
- *Noise levels.* The noise condition, namely  $\sigma \sqrt{n \log n/p} \lesssim \sigma_{\min}$ , accommodates a wide range of noise levels. To see this, it is straightforward to check that this noise condition is equivalent to

$$\sigma \lesssim \sqrt{\frac{np}{\log n}} \|\mathbf{L}^*\|_{\infty}$$

as long as  $r, \mu, \kappa \asymp 1$ . In other words, the entrywise noise level  $\sigma$  is allowed to be significantly larger than the maximum magnitude of the entries in the low-rank matrix  $\mathbf{L}^*$ , as long as  $p \gg (\log n)/n$ .

- *Tolerable fraction of outliers.* The above theorem assumes that no more than a fraction  $\rho_s \lesssim 1/\log n$  of observations are corrupted by outliers. In words, our theory allows *nearly* a constant proportion (up to a logarithmic order) of the entries of  $\mathbf{L}^*$  to be corrupted with arbitrary magnitudes.

Next, we move on to the interpretation of our statistical guarantees. Note that we still assume that  $n_1 = n_2 = n$  for ease of presentation.

- *Near-optimal statistical guarantees.* Our first result (1.13a) gives an Euclidean estimation error bound of (1.3)

$$(1.15) \quad \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_F \lesssim \sigma \sqrt{\frac{n}{p}}.$$

This cannot be improved even when an oracle has informed us of the outliers  $\mathbf{S}^*$  and the tangent space of  $\mathbf{L}^*$ ; see Candès and Plan (2010), Section III.B. We remark that under similar model assumptions, the paper Wong and Lee (2017) derived an estimation error bound for a constrained version of the convex program (1.3), which asserts that this convex estimator  $\tilde{\mathbf{L}}_{\text{cvx}}$  satisfies<sup>5</sup>

$$(1.16) \quad \|\tilde{\mathbf{L}}_{\text{cvx}} - \mathbf{L}^*\|_F \lesssim \sigma n^{1.5},$$

with the proviso that  $p$  is at least on the constant order. The restriction on  $p$  arises from the dual certificate constructed in Candès et al. (2011), which is also used in the proof of Theorem 4 in Wong and Lee (2017). While this is suboptimal compared to our results in the setting considered herein, it is worth pointing out that the bound therein accommodates arbitrary noise matrix  $\mathbf{E}$  (e.g., deterministic, adversary), and here in (1.16) we specialize their result to the random noise setting, namely the noise  $\mathbf{E}$  obeys Assumption 1.5. In addition, under the full observation (i.e.,  $p = 1$ ) setting, the paper Agarwal, Negahban and Wainwright (2012) derived an estimation error bound for a convex program similar to (1.3), but with an additional constraint regularizing the spikiness of the low-rank component. Note that instead of imposing the incoherence condition as in Assumption 1.1, the prior work Agarwal, Negahban and Wainwright (2012) assumes a milder spikiness condition on  $\mathbf{L}^*$ , which only constrains the maximum entry in the matrix  $\mathbf{L}^*$  is not too large. When  $\{E_{ij}\}$  are i.i.d. drawn from  $\mathcal{N}(0, \sigma^2)$  and when there is no missing data (i.e.  $p = 1$ ), the Euclidean estimation error bound achievable by their estimator  $\mathbf{L}_{\text{cvx}}^{\text{ANW}}$  reads

$$(1.17) \quad \|\mathbf{L}_{\text{cvx}}^{\text{ANW}} - \mathbf{L}^*\|_F \lesssim \sigma \sqrt{n} \max\{1, \sqrt{n\rho_s \log n}\} + \|\mathbf{L}^*\|_\infty n \sqrt{\rho_s},$$

which is suboptimal compared to our results. In particular, (i) the bound (1.17) does not vanish even as the noise level decreases to zero, and (ii) it becomes looser as  $\rho_s$  grows (e.g., if  $\rho_s \asymp 1/\log n$ , the bound (1.17) is  $O(\sqrt{n})$  larger than our bound). Moreover, the work Agarwal, Negahban and Wainwright (2012) did not account for missing data. Similar to Agarwal, Negahban and Wainwright (2012) (but with an additional spikiness condition on  $\mathbf{S}^*$ ), the paper Klopp, Lounici and Tsybakov (2017) derived an estimation error bound

---

<sup>5</sup>More specifically, Wong and Lee (2017), Theorem 4, studies the following convex program  $\text{minimize}_{\mathbf{L}, \mathbf{S} \in \mathbb{R}^{n \times n}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1$  s.t.  $\|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{L} + \mathbf{S} - \mathbf{M})\|_F \leq \delta$ . Here, the quantity  $\delta$  needs to be larger than  $\|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{L} + \mathbf{S} - \mathbf{M})\|_F$ . Under our setting, the minimum level of  $\delta$  should be a high-probability upper bound on  $\|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{E})\|_F$ , which is on the order of  $\sigma n \sqrt{p}$ . With this choice of  $\delta$ , Wong and Lee (2017), Theorem 4, yields  $\|\tilde{\mathbf{L}}_{\text{cvx}} - \mathbf{L}^*\|_F \leq [2 + 8\sqrt{n}(1 + \sqrt{8/p})]\delta \lesssim \sigma n^{1.5}$ .

TABLE 1  
*Comparison of our statistical guarantee and prior theory*

	Euclidean estimation error	Accounting for missing data
Zhou et al. (2010)	$\sigma n^2$	no
Agarwal, Negahban and Wainwright (2012)	$\sigma \sqrt{n} \max\{\sqrt{r}, \sqrt{n\rho_s \log n}\}$ $+ \ \mathbf{L}^*\ _{\infty} n \sqrt{\rho_s}$	no
Wong and Lee (2017)	$\sigma n^{1.5}$	yes ( $p \gtrsim 1$ )
Klopp, Lounici and Tsybakov (2017)	$\max\{\sigma, \ \mathbf{L}^*\ _{\infty}, \ \mathbf{S}^*\ _{\infty}\} \cdot$ $\sqrt{(n \log n)/p} \max\{1, \sqrt{np\rho_s}\}$	yes ( $p \gtrsim (\text{poly } \log n)/n$ )
This paper	$\sigma \sqrt{nr/p}$	yes ( $p \gtrsim \kappa^4 \mu^2 r^2 (\text{poly } \log n)/n$ )

for a constrained convex program, with a new constraint regularizing the spikiness of the sparse outliers. Their Euclidean estimation error bound reads

$$(1.18) \quad \|\mathbf{L}_{\text{cvx}}^{\text{KLT}} - \mathbf{L}^*\|_{\text{F}} \lesssim \max\{\sigma, \|\mathbf{L}^*\|_{\infty}, \|\mathbf{S}^*\|_{\infty}\} \sqrt{\frac{n \log n}{p}} \max\{1, \sqrt{np\rho_s}\},$$

which is also suboptimal compared to our results. In particular, (1) their error bound degrades as the magnitude  $\|\mathbf{S}^*\|_{\infty}$  of the outlier increases; (2) when there is no missing data (i.e.,  $p = 1$ ), their bound might be off by a factor as large as  $O(\sqrt{n})$ . It is worth emphasizing that the theory developed in these prior works is developed to accommodate a broader range of matrices. For example, both Agarwal, Negahban and Wainwright (2012) and Klopp, Lounici and Tsybakov (2017) study the set of entrywise bounded low-rank matrices (without assuming the incoherence condition); Agarwal, Negahban and Wainwright (2012) even allows  $\mathbf{L}^*$  to be approximately low rank. To ease comparison, Table 1 displays a summary of our results versus prior statistical guarantees when specialized to the settings considered herein.

- *Entrywise and spectral norm error control.* Moving beyond Euclidean estimation errors, our theory also provides statistical guarantees measured by two other important metrics: the entrywise  $\ell_{\infty}$  norm (cf. (1.13b)) and the spectral norm (cf. (1.13c)). In particular, our entrywise error bound (1.13b) in reads

$$(1.19) \quad \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_{\infty} \lesssim \sigma \sqrt{\frac{\log n}{np}}$$

as long as  $r, \kappa, \mu \asymp 1$ , which is about  $O(n)$  times small than the Euclidean loss (1.15) modulo some logarithmic factor. This uncovers an appealing “delocalization” behavior of the estimation errors, namely, the estimation errors of  $\mathbf{L}^*$  are fairly spread out across all entries. This can also be viewed as an “implicit regularization” phenomenon: the convex program automatically controls the spikiness of the low-rank solution, without the need of explicitly regularizing it (e.g., adding a constraint  $\|\mathbf{L}\|_{\infty} \leq \alpha$  as adopted in the prior work Agarwal, Negahban and Wainwright (2012), Klopp, Lounici and Tsybakov (2017)). See Figure 2 for the numerical evidence for the relative entrywise and spectral norm error of  $\mathbf{L}_{\text{cvx}}$ .

- *Approximate low-rank structure of the convex estimator  $\mathbf{L}_{\text{cvx}}$ .* Last but not least, Theorem 1.6 ensures that the convex estimate  $\mathbf{L}_{\text{cvx}}$  is nearly rank- $r$ , so that a rank- $r$  approximation of  $\mathbf{L}_{\text{cvx}}$  is extremely accurate. In other words, the convex program automatically

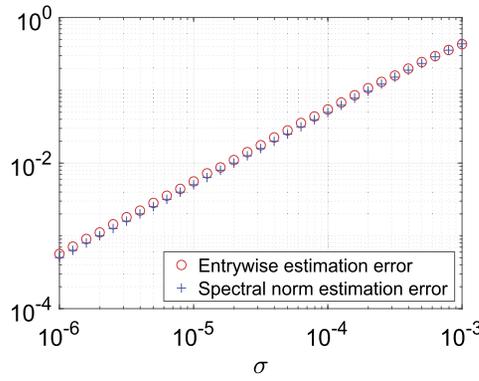


FIG. 2. The relative estimation error of  $\mathbf{L}_{\text{cvx}}$  measured by both  $\|\cdot\|_\infty$  (i.e.,  $\|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_\infty / \|\mathbf{L}^*\|_\infty$ ) and  $\|\cdot\|$  (i.e.,  $\|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\| / \|\mathbf{L}^*\|$ ) versus the standard deviation  $\sigma$  of the noise in a log-log plot. The results are reported for  $n = 1000$ ,  $r = 5$ ,  $p = 0.2$ ,  $\rho_s = 0.1$ ,  $\lambda = 5\sigma\sqrt{np}$ ,  $\tau = 2\sigma\sqrt{\log n}$ , and are averaged over 50 independent trials. In addition, the data generating process is similar to that in Figure 1.

adapts to the true rank of  $\mathbf{L}^*$  without having any prior knowledge about  $r$ . As we shall see shortly, this is a crucial observation underlying the intimate connection between convex relaxation and a certain nonconvex approach.

Moving beyond the setting with  $r, \kappa, \mu \asymp 1$ , we have developed theoretical guarantees that allow  $r, \kappa, \mu$  to grow with the problem dimension  $n_1, n_2$ . The result is this.

**THEOREM 1.7.** *Suppose that Assumptions 1.1–1.5 hold and that  $n_1 \geq n_2$ . Take  $\lambda = C_\lambda \sigma \sqrt{n_1 p}$  and  $\tau = C_\tau \sigma \sqrt{\log n_2}$  in (1.3) for some large enough constants  $C_\lambda, C_\tau > 0$ . Assume that*

$$(1.20) \quad \begin{aligned} n_1 n_2 p &\geq C_{\text{sample}} \kappa^4 \mu^2 r^2 n_1 \log^6 n_1, \\ \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n_1}{p}} &\leq \frac{c_{\text{noise}}}{\sqrt{\kappa^4 \mu r \log n_1}}, \quad \text{and} \quad \rho_s \leq \frac{c_{\text{outlier}}}{\kappa^3 \mu r \log n_1} \end{aligned}$$

for some sufficiently large constant  $C_{\text{sample}} > 0$  and some sufficiently small constants  $c_{\text{noise}}, c_{\text{outlier}} > 0$ . Then with probability exceeding  $1 - O(n_2^{-3})$ , the following hold:

1. Any minimizer  $(\mathbf{L}_{\text{cvx}}, \mathbf{S}_{\text{cvx}})$  of the convex program (1.3) obeys

$$(1.21a) \quad \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_{\text{F}} \leq C_{\text{err}} \kappa \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n_1}{p}} \|\mathbf{L}^*\|_{\text{F}},$$

$$(1.21b) \quad \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_\infty \leq C_{\text{err}} \sqrt{\kappa^3 \mu r} \cdot \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n_1 \log n_1}{p}} \|\mathbf{L}^*\|_\infty,$$

$$(1.21c) \quad \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\| \leq C_{\text{err}} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n_1}{p}} \|\mathbf{L}^*\|$$

for some constant  $C_{\text{err}} > 0$ .

2. Letting  $\mathbf{L}_{\text{cvx},r} := \arg \min_{\mathbf{L}: \text{rank}(\mathbf{L}) \leq r} \|\mathbf{L} - \mathbf{L}_{\text{cvx}}\|_{\text{F}}$  be the best rank- $r$  approximation of  $\mathbf{L}_{\text{cvx}}$ , we have

$$(1.22) \quad \|\mathbf{L}_{\text{cvx},r} - \mathbf{L}_{\text{cvx}}\|_{\text{F}} \leq \frac{1}{n_2^5} \cdot \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n_1}{p}} \|\mathbf{L}^*\|_{\text{F}},$$

and the statistical guarantees (1.21) hold unchanged if  $\mathbf{L}_{\text{cvx}}$  is replaced by  $\mathbf{L}_{\text{cvx},r}$ .

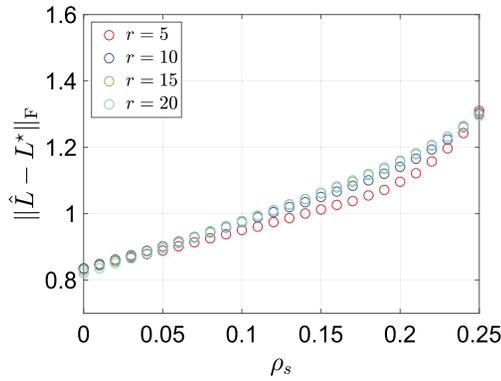


FIG. 3. Euclidean estimation errors of (1.3) versus  $\rho_s$  under four different ranks  $r = 5, 10, 15, 20$ . The results are reported for  $n = 1000$ ,  $p = 0.04r$ ,  $\sigma = 10^{-3}$ ,  $\lambda = 5\sigma\sqrt{np}$ ,  $\tau = 2\sigma\sqrt{\log n}$ , and are averaged over 50 independent trials. In addition, the data generating process is similar to that in Figure 1.

Similar to Theorem 1.6, our general theory (i.e., Theorem 1.7) provides the estimation error of the convex estimator  $\mathbf{L}_{\text{cvx}}$  in three different norms (i.e., the Euclidean, entrywise and operator norms), and reveals the near low-rankness of the convex estimator (cf. (1.22)) as well as the implicit regularization phenomenon (cf. (1.21b)).

Finally, we make note of several aspects of our general theory that call for further improvement. For instance, when there is no missing data and  $n_1 = n_2 = n$ , the rank  $r$  of the unknown matrix  $\mathbf{L}^*$  needs to satisfy  $r \lesssim \sqrt{n}$ . On the positive side, our result allows  $r$  to grow with the problem dimension  $n$ . However, prior results in the noiseless case (Candès et al. (2011), Li (2013)) allow  $r$  to grow almost linearly with  $n$ . This unsatisfactory aspect arises from the suboptimal analysis (in terms of the dependency on  $r$ ) of a tightly related nonconvex estimation algorithm (to be elaborated on later), which to the best of our knowledge, has not been resolved in the nonconvex low-rank matrix recovery literature (Chen, Liu and Li (2020), Ma et al. (2020)). See Section 2 for more discussions about this point. Moreover, when  $\mathbf{E} = \mathbf{0}$ , it is known that  $\rho_s$  can be as large as a constant even when the rank  $r$  is allowed to grow with the dimension  $n$  (Chen et al. (2013), Li (2013)). Our current theory, however, fails to cover the case with  $\rho_s \asymp 1$  in the presence of noise. We demonstrate through numerical experiments that the dependence of  $\rho_s$  on  $r$  might indeed be suboptimal in our current theory. More specifically, Figure 3 depicts the numerical Euclidean estimation errors w.r.t. the corruption probability  $\rho_s$  as we vary the rank while fixing the sampling ratio. It can be seen that the estimation error curves corresponding to different ranks align very well with each other, thus suggesting the capability of convex relaxation in tolerating a constant fraction  $\rho_s$  of outliers.

1.5. *A peek at our technical approach.* Before delving into the proof details, we immediately highlight our key technical ideas and novelties. For simplicity, we assume  $n_1 = n_2 = n$  throughout this section.

*Connections between convex and nonconvex optimization.* Instead of directly analyzing the convex program (1.3), we turn attention to a seemingly different, but in fact closely related, nonconvex program

$$(1.23) \quad \underset{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}, \mathbf{S} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \frac{1}{2} \|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{X}\mathbf{Y}^\top + \mathbf{S} - \mathbf{M})\|_{\text{F}}^2 + \frac{\lambda}{2} (\|\mathbf{X}\|_{\text{F}}^2 + \|\mathbf{Y}\|_{\text{F}}^2) + \tau \|\mathbf{S}\|_1.$$

This idea is inspired by an interesting numerical finding (cf. Figure 4) that the solution to the convex program (1.3), and an estimate obtained by attempting to solve the nonconvex formulation (1.23), are exceedingly close in our experiments. If such an intimate connection

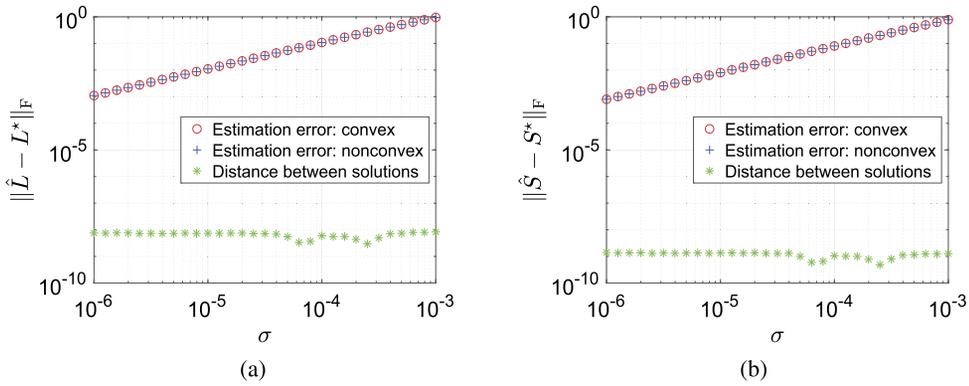


FIG. 4. (a) The relative estimation errors of both  $L_{\text{cvx}}$  (the convex estimator ((1.3))) and  $L_{\text{ncvx}}$  (the estimate returned by the nonconvex approach tailored to ((1.23))) and the relative distance between them versus the standard deviation  $\sigma$  of the noise. (b) The relative estimation errors of both  $S_{\text{cvx}}$  (the convex estimator in ((1.3))) and  $S_{\text{ncvx}}$  (the estimate returned by the nonconvex approach tailored to ((1.23))) and the relative distance between them versus the standard deviation  $\sigma$  of the noise. The results are reported for  $n = 1000$ ,  $r = 5$ ,  $p = 0.2$ ,  $\rho_S = 0.1$ ,  $\lambda = 5\sigma\sqrt{np}$ ,  $\tau = 2\sigma\sqrt{\log n}$  and are averaged over 50 independent trials.

can be formalized, then it suffices to analyze the statistical performance of the nonconvex approach instead.<sup>6</sup> Fortunately, recent advances in nonconvex low-rank factorization (see Chi, Lu and Chen (2019) for an overview) provide powerful tools for analyzing nonconvex low-rank estimation, allowing us to derive the desired statistical guarantees that can then be transferred to the convex approach. Of course, this is merely a high-level picture of our proof strategy, and we defer the details to Section 3.

It is worth emphasizing that our key idea—that is, bridging convex and nonconvex optimization—is drastically different from previous technical approaches for analyzing convex estimators (e.g., (1.3)). As it turns out, these prior approaches, which include constructing dual certificates and/or exploiting restricted strong convexity, have their own deficiencies in analyzing (1.3) and fall short of explaining the effectiveness of (1.3) in the random noise setting. For instance, constructing dual certificates in the noisy case is notoriously challenging given that we do not have closed-form expressions for the primal solutions (so that it is difficult to invoke the powerful dual construction strategies like the golfing scheme (Gross (2011)) developed for the noiseless case). If we directly utilize the dual certificates constructed for the noiseless case, we would end up with an overly conservative bound like (1.4), which is exactly why the results in Wong and Lee (2017), Zhou et al. (2010) are suboptimal. On the other hand, while it is viable to show certain strong convexity of (1.3) when restricted to some highly local sets and directions, it is unclear how (1.3) forces its solution to stay within the desired set and follow the desired directions, without adding further (and often unnecessary) constraints to (1.3).

*Nonconvex low-rank estimation with nonsmooth loss functions.* It is worth noting that a similar connection between convex and nonconvex optimization has been pointed out by Chen et al. (2020a) toward understanding the power of convex relaxation for noisy matrix completion. Due to the absence of sparse outliers in the noisy matrix completion problem, the nonconvex loss function considered therein is smooth in nature, which greatly simplifies

<sup>6</sup>On the surface, the convex program (1.3) and the nonconvex one (1.23) are closely related: the convex solution ( $L_{\text{cvx}}, S_{\text{cvx}}$ ) coincides with that of the nonconvex program (1.23) if  $L_{\text{cvx}}$  is rank- $r$ . This is an immediate consequence of the algebraic identity  $\|Z\|_* = \inf_{X, Y \in \mathbb{R}^{n \times r}, XY^T = Z} (\|X\|_F^2 + \|Y\|_F^2)$  (Mazumder, Hastie and Tibshirani (2010), Srebro and Shraibman (2005)). However, it is difficult to know *a priori* the rank of the convex solution. Hence such a connection does not prove useful in establishing the statistical properties of the convex estimator.

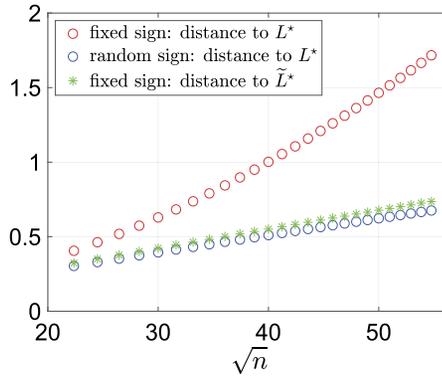


FIG. 5. The red (resp., blue) line displays the Euclidean estimation error of ((1.3)) versus  $\sqrt{n}$  under fixed (resp., random) sign patterns of  $\mathbf{S}^*$ . The green line displays the Euclidean distance between  $\mathbf{L}_{\text{cvx}}$  and  $\tilde{\mathbf{L}}^*$  under fixed sign patterns of  $\mathbf{S}^*$ . The results are reported for  $r = 5$ ,  $p = 1$ ,  $\sigma = 10^{-3}$  and  $\rho_s = 1/\log n$ , with  $\lambda = 5\sigma\sqrt{np}$  and  $\tau = 2\sigma\sqrt{\log n}$  and are averaged over 50 independent trials. For the random sign setting, the nonzero entries of  $\mathbf{S}^*$  are independently generated as  $z \cdot 5\sigma$ , where  $z$  follows a Rademacher distribution. For the fixed sign setting, each nonzero entry of  $\mathbf{S}^*$  equals to  $5\sigma$ .

both the algorithmic and theoretical development. By contrast, the nonsmoothness inherent in (1.23) makes it particularly challenging to achieve the two desiderata mentioned above, namely, connecting the convex and nonconvex solutions and establishing the optimality of the nonconvex solution. In fact, to establish the connection between convex and nonconvex solutions, we put forward a novel two-step analysis strategy. Specifically, we first develop a crude upper bound on the Euclidean estimation error leveraging the idea of approximate dual certificates; see Theorem 3.1. While this crude upper bound is far from optimal, it serves as an important starting point toward formalizing the intimate relation between the convex solution  $(\mathbf{L}_{\text{cvx}}, \mathbf{S}_{\text{cvx}})$  and the nonconvex solution  $(\mathbf{X}, \mathbf{Y}, \mathbf{S})$ , since it is challenging to establish  $\mathbf{X}\mathbf{Y} \approx \mathbf{L}_{\text{cvx}}$  and  $\mathbf{S} \approx \mathbf{S}_{\text{cvx}}$  simultaneously without the aid of a crude bound. Second, in establishing the optimality of the nonconvex solution, the nonsmoothness nature of the nonconvex loss prevents us from applying the vanilla gradient descent scheme (as has been done in Chen et al. (2020a)). To address this issue, we develop an alternating minimization scheme—which alternates between gradient updates on  $(\mathbf{X}, \mathbf{Y})$  and minimization of  $\mathbf{S}$ —aimed at minimizing the nonsmooth nonconvex loss function (1.23); see Algorithm 1 for details. As it turns out, such a simple algorithm allows us to track the proximity of the convex and nonconvex solutions and establish the optimality of the nonconvex solution all at once.

1.6. *Random signs of outliers.* The careful reader might wonder whether it is possible to remove the random sign assumption on  $\mathbf{S}^*$  (namely, Assumption 1.4) without compromising our statistical guarantees. After all, the results of Candès et al. (2011), Chandrasekaran et al. (2011), Li (2013) derived for the noise-free case do not rely on such a random sign assumption at all.<sup>7</sup> Unfortunately, removal of such a condition might be problematic in general, as illustrated by the following example.

*An example with nonrandom signs.* Suppose that (i)  $n_1 = n_2 = n$ , (ii) each nonzero entry of  $\mathbf{S}^*$  obeys  $S_{ij}^* = c_0\sigma$ , (iii)  $\rho_s = c_1/\log n$  for some sufficiently small constant  $c_1 > 0$ , and (iv) there is no missing data (i.e.,  $p = 1$ ). In such a scenario, the data matrix can be decomposed

<sup>7</sup>Notably, in the noisy setting, prior theory (Wong and Lee (2017), Zhou et al. (2010)) also implicitly assumes this random sign condition, while Agarwal, Negahban and Wainwright (2012), Klopp, Lounici and Tsybakov (2017) do not require this condition.

as

$$M = L^* + S^* + E = \underbrace{L^* + \mathbb{E}[S^*]}_{=: \tilde{L}^*} + \underbrace{S^* - \mathbb{E}[S^*]}_{=: \tilde{E}} + E.$$

Two observations are worth noting: (1) given that  $\mathbb{E}[S^*] = c_0 \rho_s \sigma \mathbf{1}\mathbf{1}^\top$  with  $\mathbf{1}$  the all-one vector, the rank of the matrix  $\tilde{L}^* = L^* + \mathbb{E}[S^*]$  is at most  $r + 1$ ; (2)  $\tilde{E}$  is a zero-mean random matrix consisting of independent entries with sub-Gaussian norm  $O(\sigma)$ . In other words, the decomposition  $M = \tilde{L}^* + \tilde{E}$  corresponds to a case with random noise but no outliers. Consequently, we can invoke Theorem 1.6 to conclude that (assuming  $r = O(1)$  and  $\tilde{L}^*$  is incoherent with condition number  $O(1)$ ): any minimizer  $(L_{\text{cvx}}, S_{\text{cvx}})$  of (1.3) obeys

$$\begin{aligned} \|L_{\text{cvx}} - L^* - \rho_s \sigma \mathbf{1}\mathbf{1}^\top\|_F &= \|L_{\text{cvx}} - \tilde{L}^*\|_F \lesssim \frac{\sigma}{\sigma_{\min}(\tilde{L}^*)} \sqrt{n} \|\tilde{L}^*\|_F \\ &\lesssim \sigma \sqrt{nr} \frac{\sigma_{\max}(\tilde{L}^*)}{\sigma_{\min}(\tilde{L}^*)} \lesssim \sigma \sqrt{n} \end{aligned}$$

with high probability. Here, the last step follows since  $\tilde{L}^*$  is of constant rank and condition number. This, however, leads to a lower bound on the estimation error

$$\begin{aligned} \|L_{\text{cvx}} - L^*\|_F &\geq \|c_0 \rho_s \sigma \mathbf{1}\mathbf{1}^\top\|_F - \|L_{\text{cvx}} - L^* - \rho_s \sigma \mathbf{1}\mathbf{1}^\top\|_F = \sigma (c_0 \rho_s n - O(\sqrt{n})) \\ &= (1 - o(1)) \frac{c_0 c_1 \sigma n}{\log n}, \end{aligned}$$

which can be  $O(\sqrt{n}/\log n)$  times larger than the desired estimation error  $O(\sigma \sqrt{n})$ . Numerical experiments under the above setting (with  $c_0 = 5$  and  $c_1 = 1$ ) also suggest that (i) the estimation error under the fixed sign setting might be orderwise larger than that under the random sign setting; and (ii) under the fixed sign setting, the estimator (1.3) approximately recovers  $\tilde{L}^*$  instead of  $L^*$ ; see Figure 5.

The take-away message is this: when the entries of  $S^*$  are of nonrandom signs, it might sometimes be possible to decompose  $S^*$  into (1) a low-rank bias component with a large Euclidean norm, and (2) a random fluctuation component whose typical size does not exceed that of  $E$ . If this is the case, then the convex program (1.3) might mistakenly treat the bias component as a part of the low-rank matrix  $L^*$ , thus dramatically hampering its estimation accuracy.

**2. Prior art.** Principal component analysis (PCA) (Fan et al. (2018), Jolliffe (1986), Pearson (1901)) is one of the most widely used statistical methods for dimension reduction in data analysis. However, PCA is known to be quite sensitive to adversarial outliers—even a single corrupted data point can make PCA completely off. This motivated the investigation of robust PCA, which aims at making PCA robust to gross adversarial outliers. As formulated in Candès et al. (2011), Chandrasekaran et al. (2011), this is closely related to the problem of disentangling a low-rank matrix  $L^*$  and a sparse outlier matrix  $S^*$  (with unknown locations and magnitudes) from a superposition of them. Consequently, robust PCA can be viewed as an outlier-robust extension of the low-rank matrix estimation/completion tasks (Candès and Recht (2009), Chi, Lu and Chen (2019), Keshavan, Montanari and Oh (2010)). In a similar vein, robust PCA has also been extensively studied in the context of structured covariance estimation under approximate factor models (Fan, Fan and Lv (2008), Fan, Liao and Mincheva (2013), Fan, Wang and Zhong (2017), Fan, Wang and Zhong (2019)), where the population covariance of certain random sample vectors is a mixture of a low-rank matrix and a sparse matrix, corresponding to the factor component and the idiosyncratic component, respectively.

Focusing on the convex relaxation approach, Candès et al. (2011), Chandrasekaran et al. (2011) started by considering the noiseless case with no missing data (i.e.,  $\mathbf{E} = \mathbf{0}$  and  $p = 1$ ) and demonstrated that, under mild conditions, convex relaxation succeeds in exactly decomposing both  $\mathbf{L}^*$  and  $\mathbf{S}^*$  from the data matrix  $\mathbf{L}^* + \mathbf{S}^*$ . More specifically, Chandrasekaran et al. (2011) adopted a deterministic model without assuming any probabilistic structure on the outlier matrix  $\mathbf{S}^*$ . As shown in Chandrasekaran et al. (2011) and several subsequent works (Chen et al. (2013), Hsu, Kakade and Zhang (2011)), convex relaxation is guaranteed to work as long as the fraction of outliers in each row/column does not exceed  $O(1/r)$ . In contrast, Candès et al. (2011) proposed a random model by assuming that  $\mathbf{S}^*$  has random support (cf. Assumption 1.3); under this model, exact recovery is guaranteed even if a constant fraction of the entries of  $\mathbf{S}^*$  are nonzero with arbitrary magnitudes. Following the random location model proposed in Candès et al. (2011), the paper Ganesh et al. (2010) showed that, in the absence of noise, convex programming can provably tolerate a dominant fraction of outliers, provided that the signs of the nonzero entries of  $\mathbf{S}^*$  are randomly generated (cf. Assumption 1.4). Later, the papers Chen et al. (2013), Li (2013) extended these results to the case when most entries of the matrix are unseen; even in the presence of highly incomplete data, convex relaxation still succeeds when a constant proportion of the observed entries are arbitrarily corrupted. It is worth noting that the results of Chen et al. (2013) accommodated both models proposed in Chandrasekaran et al. (2011) and Candès et al. (2011), while the results of Li (2013) focused on the latter model.

The literature on robust PCA with not only sparse outliers but also dense noise—namely, when the measurements take the form  $\mathbf{M} = \mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{L}^* + \mathbf{S}^* + \mathbf{E})$ —is relatively scarce. Agarwal, Negahban and Wainwright (2012), Zhou et al. (2010) were among the first to present a general theory for robust PCA with dense noise, which was further extended in Klopp, Lounici and Tsybakov (2017), Wong and Lee (2017). As we mentioned before, the first three (Agarwal, Negahban and Wainwright (2012), Wong and Lee (2017), Zhou et al. (2010)) accommodated arbitrary noise with the last one (Klopp, Lounici and Tsybakov (2017)) focusing on the random noise. As we have discussed in Section 1.4, the statistical guarantees provided in these papers are highly suboptimal when it comes to the random noise setting considered herein. The paper Chen and Chi (2014) extended the robust PCA results to the case where the truth is not only low-rank but also of Hankel structure. The results therein, however, suffered from the same suboptimality issue.

Moving beyond convex relaxation methods, another line of work proposed nonconvex approaches for robust PCA (Cai, Cai and Wei (2019), Cherapanamjeri, Gupta and Jain (2017), Drusvyatskiy, Ioffe and Lewis (2021), Gu, Wang and Liu (2016), Li et al. (2019), Netrapalli et al. (2014), Yi et al. (2016), Zhang, Wang and Gu (2018)), largely motivated by the recent success of nonconvex methods in low-rank matrix factorization (Cai and Wei (2018), Candès, Li and Soltanolkotabi (2015), Charisopoulos et al. (2019), Chen and Candès (2017), Chen and Candès (2018), Chen and Wainwright (2015), Chen et al. (2019a), Chi, Lu and Chen (2019), Jain, Netrapalli and Sanghavi (2013), Keshavan, Montanari and Oh (2010), Ma et al. (2020), Netrapalli, Jain and Sanghavi (2015), Sun and Luo (2016), Wang, Giannakis and Eldar (2018), Wei et al. (2016), Zhang, Chi and Liang (2016), Zheng and Lafferty (2016)). Following the deterministic model of Chandrasekaran et al. (2011), the paper Netrapalli et al. (2014) proposed an alternating projection/minimization scheme to seek a low-rank and sparse decomposition of the observed data matrix. In the noiseless setting, that is,  $\mathbf{E} = \mathbf{0}$ , this alternating minimization scheme provably disentangles the low-rank and sparse matrix from their superposition under mild conditions. In addition, Netrapalli et al. (2014) extended their result to the arbitrary noise case where the size of the noise is extremely small, namely,  $\|\mathbf{E}\|_{\infty} \ll \sigma_{\min}/n$ . When the noise  $\{E_{ij}\} \sim \mathcal{N}(0, \sigma^2)$ , this is equivalent to the condition  $\sigma \ll \sigma_{\min}/(n\sqrt{\log n})$ . Comparing this with our noise condition  $\sigma \ll \sigma_{\min}/(\sqrt{n\log n})$

(cf. (1.12)) when  $r, \mu, \kappa \asymp 1$ , one sees that our theoretical guarantees cover a wider range of noise levels. Similarly, Yi et al. (2016) applied regularized gradient descent on a smooth nonconvex loss function which enjoys provable convergence guarantees to  $(L^*, S^*)$  under the noiseless and partial observation setting. A recent paper Drusvyatskiy, Ioffe and Lewis (2021) considered the nonsmooth nonconvex formulation for robust PCA and established rigorously the convergence of subgradient-type methods in the rank-1 setting, that is,  $r = 1$ . However, the extension to more general rank remains out of reach.

It is worth noting that noisy matrix completion problem (Candès and Plan (2010), Chen et al. (2020a)) is subsumed as a special case by the model studied in this paper (namely, it is a special case with  $S^* = \mathbf{0}$ ). Statistical optimality under the random noise setting (cf. Assumption 1.5)—including the convex relaxation approach (Chen et al. (2020a), Klopp (2014), Koltchinskii, Lounici and Tsybakov (2011), Negahban and Wainwright (2012)) and the non-convex approach (Chen, Liu and Li (2020), Ma et al. (2020))—has been extensively studied. Focusing on arbitrary deterministic noise, Candès and Plan (2010) established the stability of the convex approach, whose resulting estimation error bound is similar to the one established for robust PCA with noise in Zhou et al. (2010)) (see (1.4)). The paper Krahmer and Stöger (2021) later confirmed that the estimation error bound established in Candès and Plan (2010) is the best one can hope for in the arbitrary noise setting for matrix completion, although it might be highly suboptimal if we restrict attention to random noise.

Finally, there is also a large literature considering robust PCA under different settings and/or from different perspectives. For instance, the computational efficiency in solving the convex optimization problem (1.3) and its variants has been studied in the optimization literature (e.g., Goldfarb, Ma and Scheinberg (2013), Ma and Aybat (2018), Shen, Wen and Zhang (2014), Tao and Yuan (2011)). The problem has also been investigated under a streaming/online setting (Feng, Xu and Yan (2013), Guo, Qiu and Vaswani (2014), Qiu and Vaswani (2010), Qiu et al. (2014), Vaswani and Narayanamurthy (2018), Zhan et al. (2016)). These are beyond the scope of the current paper.

**3. Architecture of the proof.** In this section, we give an outline for proving Theorem 1.7. The proof of Theorem 1.6 follows immediately as it is a special case of Theorem 1.7. For simplicity of presentation, our proof sets  $n_1 = n_2 = n$ . It is straightforward to obtain the proof for the general rectangular case via minor modification.

The main ingredient of the proof lies in establishing an intimate link between convex and nonconvex optimization. Unless otherwise noted, we shall set the regularization parameters as

$$(3.1) \quad \lambda = C_\lambda \sigma \sqrt{np} \quad \text{and} \quad \tau = C_\tau \sigma \sqrt{\log n}$$

throughout. In addition, the soft thresholding operator at level  $\tau$  is defined such that

$$(3.2) \quad \mathcal{S}_\tau(x) := \text{sign}(x) \max(|x| - \tau, 0).$$

For any matrix  $X$ , the matrix  $\mathcal{S}_\tau(X)$  is obtained by applying the soft thresholding operator  $\mathcal{S}_\tau(\cdot)$  to each entry of  $X$  separately. Additionally, we define the true low-rank factors as follows:

$$(3.3) \quad X^* := U^*(\Sigma^*)^{1/2} \quad \text{and} \quad Y^* := V^*(\Sigma^*)^{1/2},$$

where  $U^* \Sigma^* V^{*\top}$  is the SVD of the true low-rank matrix  $L^*$ .

3.1. *Crude estimation error bounds for convex relaxation.* We start by delivering a crude upper bound on the Euclidean estimation error, built upon the (approximate) duality certificate previously constructed in Chen et al. (2013). The proof is postponed to Appendix 4.

**THEOREM 3.1.** *Consider any given  $\lambda > 0$  and set  $\tau \asymp \lambda\sqrt{(\log n)/np}$ . Suppose that Assumptions 1.1–1.4 hold, and that*

$$n^2 p \geq C\mu^2 r^2 n \log^6 n \quad \text{and} \quad \rho_s \leq c$$

*hold for some sufficiently large (resp. small) constant  $C > 0$  (resp.,  $c > 0$ ). Then with probability at least  $1 - O(n^{-10})$ , any minimizer  $(\mathbf{L}_{\text{cvx}}, \mathbf{S}_{\text{cvx}})$  of the convex program (1.3) satisfies*

$$(3.4) \quad \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_{\text{F}}^2 + \|\mathbf{S}_{\text{cvx}} - \mathbf{S}^*\|_{\text{F}}^2 \lesssim \lambda^2 n^5 \log^3 n + \frac{n}{\lambda^2} \|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{E})\|_{\text{F}}^4.$$

It is worth noting that the above theorem holds true for an arbitrary noise matrix  $\mathbf{E}$ . When specialized to the case with independent sub-Gaussian noise, this crude bound admits a simpler expression as follows.

**COROLLARY 3.2.** *Take  $\lambda = C_\lambda \sigma \sqrt{np}$  and  $\tau = C_\tau \sigma \sqrt{\log n}$  for some universal constant  $C_\lambda, C_\tau > 0$ . Under the assumptions of Theorem 3.1 and Assumption 1.5, we have—with probability exceeding  $1 - O(n^{-10})$ —that*

$$(3.5) \quad \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_{\text{F}} \lesssim \sigma n^3 \log^{3/2} n \quad \text{and} \quad \|\mathbf{S}_{\text{cvx}} - \mathbf{S}^*\|_{\text{F}} \lesssim \sigma n^3 \log^{3/2} n.$$

**PROOF.** This corollary follows immediately by combining Theorem 3.1 and Lemma 3.3 below.  $\square$

**LEMMA 3.3.** *Suppose that Assumption 1.5 holds and that  $n^2 p > C_1 n \log^2 n$  for some sufficiently large constant  $C_1 > 0$ . Then with probability exceeding  $1 - O(n^{-10})$ , one has*

$$\|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{E})\| \lesssim \sigma \sqrt{np} \quad \text{and} \quad \|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{E})\|_{\text{F}} \lesssim \sigma n \sqrt{p}.$$

While the above results often lose a polynomial factor in  $n$ , namely, the optimal error bound, it serves as an important starting point that paves the way for subsequent analytical refinement.

3.2. *Approximate stationary points of the nonconvex formulation.* Instead of analyzing the convex estimator directly, we take a detour by considering the following nonconvex optimization problem

$$(3.6) \quad \begin{aligned} & \underset{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}, \mathbf{S} \in \mathbb{R}^{n \times n}}{\text{minimize}} && F(\mathbf{X}, \mathbf{Y}, \mathbf{S}) \\ & := \underbrace{\frac{1}{2p} \|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{X}\mathbf{Y}^\top + \mathbf{S} - \mathbf{M})\|_{\text{F}}^2 + \frac{\lambda}{2p} \|\mathbf{X}\|_{\text{F}}^2 + \frac{\lambda}{2p} \|\mathbf{Y}\|_{\text{F}}^2 + \frac{\tau}{p} \|\mathbf{S}\|_1}_{=: f(\mathbf{X}, \mathbf{Y}; \mathbf{S})}. \end{aligned}$$

Here,  $f(\mathbf{X}, \mathbf{Y}; \mathbf{S})$  is a function of  $\mathbf{X}$  and  $\mathbf{Y}$  with  $\mathbf{S}$  frozen, which contains the smooth component of the loss function  $F(\mathbf{X}, \mathbf{Y}, \mathbf{S})$ . As it turns out, the solution to convex relaxation (1.3) is exceedingly close to an estimate  $(\mathbf{X}, \mathbf{Y}, \mathbf{S})$  obtained by a nonconvex algorithm aimed at solving (3.6)—to be detailed in Section 3.3. This fundamental connection between the two algorithmic paradigms provides a powerful framework that allows us to understand convex relaxation by studying nonconvex optimization.

In what follows, we set out to develop the aforementioned intimate connection. Before proceeding, we first state the following conditions concerned with the interplay between the noise size, the estimation accuracy of the nonconvex estimate  $(\mathbf{X}, \mathbf{Y}, \mathbf{S})$ , and the regularization parameters.

CONDITION 3.4. *The regularization parameters  $\lambda$  and  $\tau \asymp \lambda\sqrt{(\log n)/np}$  satisfy*

- $\|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{E})\| < \lambda/16$  and  $\|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{E})\|_{\infty} \leq \tau/4$ ;
- $\|\mathbf{S} - \mathbf{S}^*\| < \lambda/16$  and  $\|\mathbf{X}\mathbf{Y}^{\top} - \mathbf{L}^*\|_{\infty} \leq \tau/4$ ;
- $\|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{X}\mathbf{Y}^{\top} - \mathbf{L}^*) - p(\mathbf{X}\mathbf{Y}^{\top} - \mathbf{L}^*)\| < \lambda/8$ .

As an interpretation, the above condition says that: (1) the regularization parameters are not too small compared to the size of the noise, so as to ensure that we enforce a sufficiently large degree of regularization; (2) the estimate represented by the point  $(\mathbf{X}\mathbf{Y}^{\top}, \mathbf{S})$  is sufficiently close to the truth. At this point, whether this condition is meaningful or not remains far from clear; we shall return to justify its feasibility shortly.

In addition, we need another condition concerning the injectivity of  $\mathcal{P}_{\Omega^*}$  w.r.t. a certain tangent space. For a rank- $r$  matrix  $\mathbf{L}$  with singular value decomposition  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\top}$  where  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times r}$ , the tangent space of the set of rank- $r$  matrices at the point  $\mathbf{L}$  is given by

$$\{\mathbf{U}\mathbf{A}^{\top} + \mathbf{B}\mathbf{V}^{\top} \mid \mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times r}\}.$$

Again, the validity of this condition will be discussed momentarily.

CONDITION 3.5 (Injectivity). *Let  $T$  be the tangent space of the set of rank- $r$  matrices at the point  $\mathbf{X}\mathbf{Y}^{\top}$ . Assume that there exist a constants  $c_{\text{inj}} > 0$  such that for all  $\mathbf{H} \in T$ , one has*

$$p^{-1}\|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{H})\|_{\text{F}}^2 \geq \frac{c_{\text{inj}}}{\kappa}\|\mathbf{H}\|_{\text{F}}^2 \quad \text{and} \quad p^{-1}\|\mathcal{P}_{\Omega^*}(\mathbf{H})\|_{\text{F}}^2 \leq \frac{c_{\text{inj}}}{4\kappa}\|\mathbf{H}\|_{\text{F}}^2.$$

With the above conditions in place, we are ready to make precise the intimate link between convex relaxation and a candidate nonconvex solution. The proof is deferred to Appendix 5.

THEOREM 3.6. *Suppose that  $n \geq \kappa$  and  $\rho_s \leq c/\kappa$  for some sufficiently small constant  $c > 0$ . Assume that there exists a triple  $(\mathbf{X}, \mathbf{Y}, \mathbf{S})$  such that*

$$(3.7) \quad \|\nabla f(\mathbf{X}, \mathbf{Y}; \mathbf{S})\|_{\text{F}} \leq \frac{1}{n^{20}} \frac{\lambda}{p} \sqrt{\sigma_{\min}}, \quad \text{and} \quad \mathbf{S} = \mathcal{P}_{\Omega_{\text{obs}}}(\mathcal{S}_{\tau}(\mathbf{M} - \mathbf{X}\mathbf{Y}^{\top})).$$

*Further, assume that any singular value of  $\mathbf{X}$  and  $\mathbf{Y}$  lies in  $[\sqrt{\sigma_{\min}}/2, \sqrt{2\sigma_{\max}}]$ . If the solution  $(\mathbf{L}_{\text{cvx}}, \mathbf{S}_{\text{cvx}})$  to the convex program (1.3) admits the following crude error bound*

$$(3.8) \quad \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_{\text{F}} \lesssim \sigma n^4,$$

*then under Conditions 3.4–3.5 we have*

$$\|\mathbf{X}\mathbf{Y}^{\top} - \mathbf{L}_{\text{cvx}}\|_{\text{F}} \lesssim \frac{\sigma}{n^5} \quad \text{and} \quad \|\mathbf{S} - \mathbf{S}_{\text{cvx}}\|_{\text{F}} \lesssim \frac{\sigma}{n^5}.$$

This theorem is a deterministic result, focusing on some sort of “approximate stationary points” of  $F(\mathbf{X}, \mathbf{Y}, \mathbf{S})$ . To interpret this, observe that in view of (3.7), one has  $\nabla f(\mathbf{X}, \mathbf{Y}; \mathbf{S}) \approx \mathbf{0}$ , and  $\mathbf{S}$  minimizes  $F(\mathbf{X}, \mathbf{Y}, \cdot)$  for any fixed  $\mathbf{X}$  and  $\mathbf{Y}$ . If one can identify such an approximate stationary point that is sufficiently close to the truth (so that it satisfies Condition 3.4), then under mild conditions our theory asserts that

$$\mathbf{X}\mathbf{Y}^{\top} \approx \mathbf{L}_{\text{cvx}} \quad \text{and} \quad \mathbf{S} \approx \mathbf{S}_{\text{cvx}}.$$

This would in turn formalize the intimate relation between the solution to convex relaxation and an approximate stationary point of the nonconvex formulation. The existence of such approximate stationary points will be verified shortly in Section 3.3.

The careful reader might immediately remark that this theorem does not say anything explicit about the minimizer of the nonconvex optimization problem (3.6); rather, it only pays

---

**Algorithm 1** Alternating minimization method for solving the nonconvex problem (3.6)

---

**Suitable initialization:**  $X^0 = X^*, Y^0 = Y^*, S^0 = S^*$ .

**Gradient updates:** for  $t = 0, 1, \dots, t_0 - 1$  do

$$(3.9a) \quad X^{t+1} = X^t - \eta \nabla_X f(X^t, Y^t; S^t) = X^t - \frac{\eta}{p} [\mathcal{P}_{\Omega_{\text{obs}}}(X^t Y^{t\top} + S^t - M) Y^t + \lambda X^t];$$

$$(3.9b) \quad Y^{t+1} = Y^t - \eta \nabla_Y f(X^t, Y^t; S^t) = Y^t - \frac{\eta}{p} \{[\mathcal{P}_{\Omega_{\text{obs}}}(X^t Y^{t\top} + S^t - M)]^\top X^t + \lambda Y^t\};$$

$$(3.9c) \quad S^{t+1} = \mathcal{S}_\tau[\mathcal{P}_{\Omega_{\text{obs}}}(M - X^{t+1} Y^{t+1\top})].$$


---

attention to a special class of approximate stationary points of the nonconvex formulation. This arises mainly due to a technical consideration: it seems more difficult to analyze the nonconvex optimizer directly than to study certain approximate stationary points. Fortunately, our theorem indicates that any approximate stationary point obeying the above conditions serves as an extremely tight approximation of the convex estimate and, therefore, it suffices to identify and analyze any such points.

3.3. *Constructing an approximate stationary point via nonconvex algorithms.* By virtue of Theorem 3.6, the key to understanding convex relaxation is to construct an approximate stationary point of the nonconvex problem (3.6) that enjoys desired statistical properties. For this purpose, we resort to the iterative algorithm (Algorithm 1) to solve the nonconvex program (3.6).

In a nutshell, Algorithm 1 alternates between one iteration of gradient updates (w.r.t. the decision matrices  $X$  and  $Y$ ) and optimization of the nonsmooth problem w.r.t.  $S$  (with  $X$  and  $Y$  frozen).<sup>8</sup> For the sake of simplicity, we initialize this algorithm from the ground truth  $(X^*, Y^*, S^*)$ , but our analysis framework might be extended to accommodate other more practical initialization (e.g., the one obtained by a spectral method [Chen et al. \(2020b\)](#)).

The following theorem makes precise the statistical guarantees of the above nonconvex optimization algorithm; the proof is deferred to Appendix 6. Here and throughout, we define

$$(3.10) \quad H^t := \arg \min_{R \in \mathcal{O}^{r \times r}} (\|X^t R - X^*\|_F^2 + \|Y^t R - Y^*\|_F^2)^{1/2},$$

where  $\mathcal{O}^{r \times r}$  denotes the set of  $r \times r$  orthonormal matrices.

**THEOREM 3.7.** *Instate the assumptions of Theorem 1.7 and define*

$$\delta_n := \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}}.$$

Take  $t_0 = n^{47}$  and  $\eta \asymp 1/(n\kappa^3 \sigma_{\max})$  in Algorithm 1. With probability at least  $1 - O(n^{-3})$ , the iterates  $\{(X^t, Y^t, S^t)\}_{0 \leq t \leq t_0}$  of Algorithm 1 satisfy

$$(3.11a) \quad \max\{\|X^t H^t - X^*\|_F, \|Y^t H^t - Y^*\|_F\} \lesssim \delta_n \|X^*\|_F,$$

$$(3.11b) \quad \max\{\|X^t H^t - X^*\|, \|Y^t H^t - Y^*\|\} \lesssim \delta_n \|X^*\|,$$

---

<sup>8</sup>Note that for any given  $X$  and  $Y$ , the solution to minimize  $\mathcal{S}F(X, Y, S)$  is given precisely by  $\mathcal{S}_\tau(\mathcal{P}_{\Omega_{\text{obs}}}(M - XY^\top))$ .

$$(3.11c) \quad \begin{aligned} & \max\{\|X^t H^t - X^*\|_{2,\infty}, \|Y^t H^t - Y^*\|_{2,\infty}\} \\ & \lesssim \kappa \sqrt{\log n \delta_n} \max\{\|X^*\|_{2,\infty}, \|Y^*\|_{2,\infty}\}, \end{aligned}$$

$$(3.11d) \quad \|S^t - S^*\| \lesssim \sigma \sqrt{np}.$$

In addition, with probability at least  $1 - O(n^{-3})$ , one has

$$(3.12) \quad \min_{0 \leq t < t_0} \|\nabla f(X^t, Y^t; S^t)\|_F \leq \frac{1}{n^{20}} \frac{\lambda}{p} \sqrt{\sigma_{\min}}.$$

In short, the bounds (3.11a)–(3.11c) reveal that the entire sequence  $\{X^t, Y^t\}_{t=0}^{t_0}$  stays sufficiently close to the truth (measured by  $\|\cdot\|_F$ ,  $\|\cdot\|$ , and more importantly,  $\|\cdot\|_{2,\infty}$ ), the inequality (3.11d) demonstrates the goodness of fit of  $\{S^t\}_{0 \leq t \leq t_0}$  in terms of the spectral norm accuracy, whereas the last bound (3.12) indicates that there is at least one point in the sequence  $\{X^t, Y^t, S^t\}_{0 \leq t \leq t_0}$  that can serve as an approximate stationary point of the nonconvex formulation.

We shall also gather a few immediate consequences of Theorem 3.7 as follows, which contain basic properties that will be useful throughout.

**COROLLARY 3.8.** *Instate the assumptions of Theorem 3.7. Suppose that the sample size obeys  $n^2 p \gg \kappa^4 \mu^2 r^2 n \log^4 n$ , the noise satisfies  $\delta_n \ll 1/\sqrt{\kappa^4 \mu r \log n}$ , the outlier fraction satisfies  $\rho_s \ll 1/(\kappa^3 \mu r \log n)$ . With probability at least  $1 - O(n^{-3})$ , the iterates of Algorithm 1 satisfy*

$$(3.13a) \quad \|X^t Y^{t\top} - L^*\|_F \lesssim \kappa \delta_n \|L^*\|_F,$$

$$(3.13b) \quad \|X^t Y^{t\top} - L^*\|_\infty \lesssim \sqrt{\kappa^3 \mu r \log n \delta_n} \|L^*\|_\infty,$$

$$(3.13c) \quad \|X^t Y^{t\top} - L^*\| \lesssim \delta_n \|L^*\|$$

simultaneously for all  $t \leq t_0$ .

**PROOF.** See Chen et al. (2020a), Appendix D.12.  $\square$

### 3.4. Proof of Theorem 1.7. Define

$$(3.14) \quad t_* := \arg \min_{0 \leq t < t_0} \|\nabla f(X^t, Y^t; S^t)\|_F;$$

$$(3.15) \quad (X_{\text{ncvx}}, Y_{\text{ncvx}}, S_{\text{ncvx}}) := (X^{t_*} H^{t_*}, Y^{t_*} H^{t_*}, S^{t_*}).$$

Theorem 3.7 and Corollary 3.8 have established appealing statistical performance of the nonconvex solution  $(X_{\text{ncvx}}, Y_{\text{ncvx}}, S_{\text{ncvx}})$ . To transfer this desired statistical property to that of  $(L_{\text{cvx}}, S_{\text{cvx}})$ , it remains to show that the nonconvex estimator  $(X_{\text{ncvx}} Y_{\text{ncvx}}^\top, S_{\text{ncvx}})$  is extremely close to the convex estimator  $(L_{\text{cvx}}, S_{\text{cvx}})$ . Toward this end, we intend to invoke Theorem 3.6; therefore, it boils down to verifying the conditions therein.

1. The small gradient condition (cf. (3.7)) holds automatically under (3.12).
2. By virtue of the spectral norm bound (3.11b), one has

$$\|X_{\text{ncvx}} - X^*\| = \|X^{t_*} H^{t_*} - X^*\| \lesssim \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|L^*\| \leq \frac{\sqrt{\sigma_{\min}}}{10},$$

as long as  $\sigma \sqrt{\kappa n/p} \ll \sigma_{\min}$ . This together with the Weyl inequality verifies the constraints on the singular values of  $(X_{\text{ncvx}}, Y_{\text{ncvx}})$ .

3. The crude error bounds are valid in view of Theorem 3.1.

4. Regarding Condition 3.4 and Condition 3.5, Lemma 3.3 and standard inequalities about sub-Gaussian random variables imply that  $\|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{E})\| < \lambda/16$  and  $\|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{E})\|_{\infty} \leq \tau/4$ . In addition, the bounds (3.11d) and (3.13b) ensure the second assumption  $\|\mathbf{S}_{\text{ncvx}} - \mathbf{S}^*\| \leq \lambda/16$  and  $\|\mathbf{X}\mathbf{Y}^{\top} - \mathbf{L}^*\|_{\infty} \leq \tau/4$  in Condition 3.4. We are left with the last assumption in Condition 3.4 and Condition 3.5, which are guaranteed to hold in view of the following lemma (see Appendix 3 for the proof).

LEMMA 3.9. *Instate the notation and assumptions of Theorem 1.7. Then with probability exceeding  $1 - O(n^{-10})$ , we have*

$$(3.16a) \quad \|\mathcal{P}_{\Omega}(\mathbf{X}\mathbf{Y}^{\top} - \mathbf{M}^*) - p(\mathbf{X}\mathbf{Y}^{\top} - \mathbf{M}^*)\| < \lambda/8,$$

$$(3.16b) \quad \frac{1}{p} \|\mathcal{P}_{\Omega_{\text{obs}}}(\mathbf{H})\|_{\text{F}}^2 \geq \frac{1}{32\kappa} \|\mathbf{H}\|_{\text{F}}^2, \quad \forall \mathbf{H} \in T,$$

$$(3.16c) \quad p^{-1} \|\mathcal{P}_{\Omega^*}(\mathbf{H})\|_{\text{F}}^2 \leq \frac{1}{128\kappa} \|\mathbf{H}\|_{\text{F}}^2, \quad \forall \mathbf{H} \in T$$

simultaneously for all  $(\mathbf{X}, \mathbf{Y})$  obeying

$$(3.17a) \quad \|\mathbf{X} - \mathbf{X}^*\|_{2,\infty} \leq C_{\infty} \kappa \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \max\{\|\mathbf{X}^*\|_{2,\infty}, \|\mathbf{Y}^*\|_{2,\infty}\};$$

$$(3.17b) \quad \|\mathbf{Y} - \mathbf{Y}^*\|_{2,\infty} \leq C_{\infty} \kappa \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \max\{\|\mathbf{X}^*\|_{2,\infty}, \|\mathbf{Y}^*\|_{2,\infty}\}.$$

Here,  $T$  denotes the tangent space of the set of rank- $r$  matrices at the point  $\mathbf{X}\mathbf{Y}^{\top}$ , and  $C_{\infty} > 0$  is an absolute constant.

Armed with the above conditions, we can readily invoke Theorem 3.6 to reach

$$\|\mathbf{X}_{\text{ncvx}}\mathbf{Y}_{\text{ncvx}}^{\top} - \mathbf{L}_{\text{cvx}}\|_{\text{F}} \lesssim \frac{\sigma}{n^5} \quad \text{and} \quad \|\mathbf{S}_{\text{ncvx}} - \mathbf{S}_{\text{cvx}}\|_{\text{F}} \lesssim \frac{\sigma}{n^5}$$

with high probability. This taken collectively with Corollary 3.8 gives

$$\begin{aligned} \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_{\text{F}} &\leq \|\mathbf{X}_{\text{ncvx}}\mathbf{Y}_{\text{ncvx}}^{\top} - \mathbf{L}_{\text{cvx}}\|_{\text{F}} + \|\mathbf{X}_{\text{ncvx}}\mathbf{Y}_{\text{ncvx}}^{\top} - \mathbf{L}^*\|_{\text{F}} \\ &\lesssim \frac{\sigma}{n^5} + \kappa \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{L}^*\|_{\text{F}} \\ &\asymp \kappa \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{L}^*\|_{\text{F}}. \end{aligned}$$

Similar arguments lead to the advertised high-probability bounds

$$\begin{aligned} \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\|_{\infty} &\lesssim \sqrt{\kappa^3 \mu r} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|\mathbf{L}^*\|_{\infty}, \\ \|\mathbf{L}_{\text{cvx}} - \mathbf{L}^*\| &\lesssim \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{L}^*\|. \end{aligned}$$

Finally, given that  $\mathbf{X}_{\text{ncvx}}\mathbf{Y}_{\text{ncvx}}^{\top}$  is a rank- $r$  matrix, the rank- $r$  approximation  $\mathbf{L}_{\text{cvx},r} := \arg \min_{\mathbf{Z}: \text{rank}(\mathbf{Z}) \leq r} \|\mathbf{Z} - \mathbf{L}_{\text{cvx}}\|_{\text{F}}$  of  $\mathbf{L}_{\text{cvx}}$  necessarily satisfies

$$\|\mathbf{L}_{\text{cvx},r} - \mathbf{L}_{\text{cvx}}\|_{\text{F}} \leq \|\mathbf{X}_{\text{ncvx}}\mathbf{Y}_{\text{ncvx}}^{\top} - \mathbf{L}_{\text{cvx}}\|_{\text{F}} \lesssim \frac{\sigma}{n^5} \leq \frac{1}{n^5} \cdot \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{L}^*\|,$$

which establishes (1.22). In view of the triangle inequality, the properties (1.21) hold unchanged if  $\mathbf{L}_{\text{cvx}}$  is replaced by  $\mathbf{L}_{\text{cvx},r}$ .

**4. Discussion.** This paper investigates the unreasonable effectiveness of convex programming in estimating an unknown low-rank matrix from grossly corrupted data. We develop an improved theory that confirms the optimality of convex relaxation in the presence of random noise, gross sparse outliers and missing data. In particular, our results significantly improve upon the prior statistical guarantees (Zhou et al. (2010)) under random noise, while further allowing for missing data. Our theoretical analysis is built upon an appealing connection between convex and nonconvex optimization, which has not been established previously.

Having said this, our current work leaves open several important issues that call for further investigation. To begin with, the conditions (1.20) stated in the main theorem are likely suboptimal in terms of the dependency on both the rank  $r$  and the condition number  $\kappa$ . For example, we shall keep in mind that in the noise-free setting, the sample size can be as low as  $O(nr \text{poly} \log n)$  and the tolerable outlier fraction can be as large as a constant (Chen et al. (2013), Li (2013)), both of which exhibit more favorable scalings w.r.t.  $r$  and  $\kappa$  compared to our current condition (1.20). Moving forward, our analysis ideas suggest a possible route for analyzing convex relaxation for other structured estimation problems under both random noise and outliers, including but not limited to sparse PCA (the case with a simultaneously low-rank and sparse matrix) (Cai, Ma and Wu (2013)), low-rank Hankel matrix estimation (the case involving a low-rank Hankel matrix) (Chen and Chi (2014)), and blind deconvolution<sup>9</sup> (the case that aims to recover a low-rank matrix from structured Fourier measurements) (Ahmed, Recht and Romberg (2014)). Last but not least, we would like to point out that it is possible to design a similar debiasing procedure as in Chen et al. (2019b) for correcting the bias in the convex estimator, which further allows uncertainty quantification and statistical inference on the unknown low-rank matrix of interest.

**Acknowledgments.** Author names are sorted alphabetically. Y. Chen is the corresponding author.

**Funding.** Y. Chen is supported in part by the AFOSR YIP award FA9550-19-1-0030, by the ONR Grant N00014-19-1-2120, by the ARO Grants W911NF-20-1-0097 and W911NF-18-1-0303, by NSF Grants CCF-1907661, IIS-1900140, IIS-2100158, and DMS-2014279 and by the Princeton SEAS innovation award. J. Fan is supported in part by the NSF Grants DMS-1662139 and DMS-1712591, the ONR Grant N00014-19-1-2120 and the NIH Grant 2R01-GM072611-14.

## SUPPLEMENTARY MATERIAL

**Additional proofs** (DOI: [10.1214/21-AOS2066SUPP](https://doi.org/10.1214/21-AOS2066SUPP); .pdf). Additional proofs of the results in the paper can be found in the Supplementary Material (Chen et al. (2021)).

## REFERENCES

- AGARWAL, A., NEGAHBAN, S. and WAINWRIGHT, M. J. (2012). Noisy matrix decomposition via convex relaxation: Optimal rates in high dimensions. *Ann. Statist.* **40** 1171–1197. MR2985947 <https://doi.org/10.1214/12-AOS1000>
- AHMED, A., RECHT, B. and ROMBERG, J. (2014). Blind deconvolution using convex programming. *IEEE Trans. Inf. Theory* **60** 1711–1732. MR3168432 <https://doi.org/10.1109/TIT.2013.2294644>
- CAI, H., CAI, J.-F. and WEI, K. (2019). Accelerated alternating projections for robust principal component analysis. *J. Mach. Learn. Res.* **20** Paper No. 20, 33. MR3911427
- CAI, T. T., MA, Z. and WU, Y. (2013). Sparse PCA: Optimal rates and adaptive estimation. *Ann. Statist.* **41** 3074–3110. MR3161458 <https://doi.org/10.1214/13-AOS1178>

<sup>9</sup>Our ongoing work Chen et al. (2020c) is pursuing this direction.

- CAI, J.-F. and WEI, K. (2018). Solving systems of phaseless equations via Riemannian optimization with optimal sampling complexity. Preprint. [arXiv:1809.02773](https://arxiv.org/abs/1809.02773).
- CANDÈS, E. J., LI, X. and SOLTANOLKOTABI, M. (2015). Phase retrieval via Wirtinger flow: Theory and algorithms. *IEEE Trans. Inf. Theory* **61** 1985–2007. MR3332993 <https://doi.org/10.1109/TIT.2015.2399924>
- CANDÈS, E. and PLAN, Y. (2010). Matrix completion with noise. *Proc. IEEE* **98** 925–936.
- 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>
- CANDÈS, E. J. and TAO, T. (2010). The power of convex relaxation: Near-optimal matrix completion. *IEEE Trans. Inf. Theory* **56** 2053–2080. MR2723472 <https://doi.org/10.1109/TIT.2010.2044061>
- CANDÈS, E. J., LI, X., MA, Y. and WRIGHT, J. (2011). Robust principal component analysis? *J. ACM* **58** Art. 11, 37. MR2811000 <https://doi.org/10.1145/1970392.1970395>
- CHANDRASEKARAN, V., PARRILO, P. A. and WILLSKY, A. S. (2012). Latent variable graphical model selection via convex optimization. *Ann. Statist.* **40** 1935–1967. MR3059067 <https://doi.org/10.1214/11-AOS949>
- CHANDRASEKARAN, V., SANGHAVI, S., PARRILO, P. A. and WILLSKY, A. S. (2011). Rank-sparsity incoherence for matrix decomposition. *SIAM J. Optim.* **21** 572–596. MR2817479 <https://doi.org/10.1137/090761793>
- CHARISOPOULOS, V., DAVIS, D., DÍAZ, M. and DRUSVYATSKIY, D. (2019). Composite optimization for robust blind deconvolution. Preprint. [arXiv:1901.01624](https://arxiv.org/abs/1901.01624).
- CHEN, Y. (2015). Incoherence-optimal matrix completion. *IEEE Trans. Inf. Theory* **61** 2909–2923. MR3342311 <https://doi.org/10.1109/TIT.2015.2415195>
- CHEN, Y. and CANDÈS, E. J. (2017). Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Comm. Pure Appl. Math.* **70** 822–883. MR3628877 <https://doi.org/10.1002/cpa.21638>
- CHEN, Y. and CANDÈS, E. J. (2018). The projected power method: An efficient algorithm for joint alignment from pairwise differences. *Comm. Pure Appl. Math.* **71** 1648–1714. MR3847751 <https://doi.org/10.1002/cpa.21760>
- CHEN, Y. and CHI, Y. (2014). Robust spectral compressed sensing via structured matrix completion. *IEEE Trans. Inf. Theory* **60** 6576–6601. MR3265040 <https://doi.org/10.1109/TIT.2014.2343623>
- CHEN, Y., GUIBAS, L. J. and HUANG, Q. (2014). Near-optimal joint optimal matching via convex relaxation. In *International Conference on Machine Learning (ICML)* 100–108.
- CHEN, J., LIU, D. and LI, X. (2020). Nonconvex rectangular matrix completion via gradient descent without  $\ell_{2,\infty}$  regularization. *IEEE Trans. Inform. Theory* **66** 5806–5841. MR4158648 <https://doi.org/10.1109/TIT.2020.2992234>
- CHEN, Y. and WAINWRIGHT, M. J. (2015). Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. [arXiv:1509.03025](https://arxiv.org/abs/1509.03025).
- CHEN, Y., JALALI, A., SANGHAVI, S. and CARAMANIS, C. (2013). Low-rank matrix recovery from errors and erasures. *IEEE Transactions on Information Theory* **59** 4324–4337.
- CHEN, Y., JALALI, A., SANGHAVI, S. and XU, H. (2014). Clustering partially observed graphs via convex optimization. *J. Mach. Learn. Res.* **15** 2213–2238. MR3231602
- CHEN, Y., CHI, Y., FAN, J. and MA, C. (2019a). Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval. *Math. Program.* **176** 5–37. MR3960803 <https://doi.org/10.1007/s10107-019-01363-6>
- CHEN, Y., FAN, J., MA, C. and YAN, Y. (2019b). Inference and uncertainty quantification for noisy matrix completion. *Proc. Natl. Acad. Sci. USA* **116** 22931–22937. MR4036123 <https://doi.org/10.1073/pnas.1910053116>
- CHEN, Y., CHI, Y., FAN, J., MA, C. and YAN, Y. (2020a). Noisy matrix completion: Understanding statistical guarantees for convex relaxation via nonconvex optimization. *SIAM J. Optim.* **30** 3098–3121. MR4167625 <https://doi.org/10.1137/19M1290000>
- CHEN, Y., CHI, Y., FAN, J. and MA, C. (2020b). Spectral methods for data science: A statistical perspective. Preprint. [arXiv:2012.08496](https://arxiv.org/abs/2012.08496).
- CHEN, Y., FAN, J., WANG, B. and YAN, Y. (2020c). Convex and nonconvex optimization are both minimax-optimal for noisy blind deconvolution. Preprint. [arXiv:2008.01724](https://arxiv.org/abs/2008.01724).
- CHEN, Y., FAN, J., MA, C. and YAN, Y. (2021). Supplement to “Bridging convex and nonconvex optimization in robust PCA: Noise, outliers and missing data.” <https://doi.org/10.1214/21-AOS2066SUPP>
- CHERAPANAMJERI, Y., GUPTA, K. and JAIN, P. (2017). Nearly optimal robust matrix completion. In *Proceedings of the 34th International Conference on Machine Learning* **70** 797–805. JMLR.org.
- CHI, Y., LU, Y. M. and CHEN, Y. (2019). Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Trans. Signal Process.* **67** 5239–5269. MR4016283 <https://doi.org/10.1109/TSP.2019.2937282>
- DAVENPORT, M. A. and ROMBERG, J. (2016). An overview of low-rank matrix recovery from incomplete observations. *IEEE Journal of Selected Topics in Signal Processing* **10** 608–622.
- DONOHO, D. and GAVISH, M. (2014). Minimax risk of matrix denoising by singular value thresholding. *Ann. Statist.* **42** 2413–2440. MR3269984 <https://doi.org/10.1214/14-AOS1257>

- DRUSVYATSKIY, D., IOFFE, A. D. and LEWIS, A. S. (2021). Nonsmooth optimization using Taylor-like models: Error bounds, convergence, and termination criteria. *Math. Program.* **185** 357–383. MR4201717 <https://doi.org/10.1007/s10107-019-01432-w>
- FAN, J., FAN, Y. and LV, J. (2008). High dimensional covariance matrix estimation using a factor model. *J. Econometrics* **147** 186–197. MR2472991 <https://doi.org/10.1016/j.jeconom.2008.09.017>
- FAN, J., LIAO, Y. and MINCHEVA, M. (2013). Large covariance estimation by thresholding principal orthogonal complements. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **75** 603–680. MR3091653 <https://doi.org/10.1111/rssb.12016>
- FAN, J., WANG, W. and ZHONG, Y. (2017). An  $\ell_\infty$  eigenvector perturbation bound and its application to robust covariance estimation. *J. Mach. Learn. Res.* **18** Paper No. 207, 42. MR3827095
- FAN, J., WANG, W. and ZHONG, Y. (2019). Robust covariance estimation for approximate factor models. *J. Econometrics* **208** 5–22. MR3906959 <https://doi.org/10.1016/j.jeconom.2018.09.003>
- FAN, J., SUN, Q., ZHOU, W.-X. and ZHU, Z. (2018). Principal component analysis for big data. Preprint. arXiv::1801.01602.
- FENG, J., XU, H. and YAN, S. (2013). Online robust PCA via stochastic optimization. In *Advances in Neural Information Processing Systems* 404–412.
- GANESH, A., WRIGHT, J., LI, X., CANDES, E. J. and MA, Y. (2010). Dense error correction for low-rank matrices via principal component pursuit. In *2010 IEEE International Symposium on Information Theory* 1513–1517. IEEE.
- GOLDFARB, D., MA, S. and SCHEINBERG, K. (2013). Fast alternating linearization methods for minimizing the sum of two convex functions. *Math. Program.* **141** 349–382. MR3097290 <https://doi.org/10.1007/s10107-012-0530-2>
- GROSS, D. (2011). Recovering low-rank matrices from few coefficients in any basis. *IEEE Trans. Inform. Theory* **57** 1548–1566. MR2815834 <https://doi.org/10.1109/TIT.2011.2104999>
- GU, Q., WANG, Z. W. and LIU, H. (2016). Low-rank and sparse structure pursuit via alternating minimization. In *Artificial Intelligence and Statistics* 600–609.
- GUO, H., QIU, C. and VASWANI, N. (2014). An online algorithm for separating sparse and low-dimensional signal sequences from their sum. *IEEE Trans. Signal Process.* **62** 4284–4297. MR3260427 <https://doi.org/10.1109/TSP.2014.2331612>
- HSU, D., KAKADE, S. M. and ZHANG, T. (2011). Robust matrix decomposition with sparse corruptions. *IEEE Trans. Inform. Theory* **57** 7221–7234. MR2883652 <https://doi.org/10.1109/TIT.2011.2158250>
- HUANG, Q. and GUIBAS, L. (2013). Consistent shape maps via semidefinite programming. *Computer Graphics Forum* **32** 177–186.
- JAIN, P., NETRAPALLI, P. and SANGHAVI, S. (2013). Low-rank matrix completion using alternating minimization (extended abstract). In *STOC'13—Proceedings of the 2013 ACM Symposium on Theory of Computing* 665–674. ACM, New York. MR3210828 <https://doi.org/10.1145/2488608.2488693>
- JOLLIFFE, I. T. (1986). *Principal Component Analysis*. Springer Series in Statistics. Springer, New York. MR0841268 <https://doi.org/10.1007/978-1-4757-1904-8>
- KESHAVAN, R. H., MONTANARI, A. and OH, S. (2010). Matrix completion from a few entries. *IEEE Trans. Inform. Theory* **56** 2980–2998. MR2683452 <https://doi.org/10.1109/TIT.2010.2046205>
- KLOPP, O. (2014). Noisy low-rank matrix completion with general sampling distribution. *Bernoulli* **20** 282–303. MR3160583 <https://doi.org/10.3150/12-BEJ486>
- 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>
- 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>
- KRAHMER, F. and STÖGER, D. (2021). On the convex geometry of blind deconvolution and matrix completion. *Comm. Pure Appl. Math.* **74** 790–832. MR4221934
- 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>
- LI, Y., MA, C., CHEN, Y. and CHI, Y. (2019). Nonconvex matrix factorization from rank-one measurements. In *The 22nd International Conference on Artificial Intelligence and Statistics* 1496–1505.
- MA, S. and AYBAT, N. S. (2018). Efficient optimization algorithms for robust principal component analysis and its variants. *Proceedings of the IEEE* **106** 1411–1426.
- 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>
- MAZUMDER, R., HASTIE, T. and TIBSHIRANI, R. (2010). Spectral regularization algorithms for learning large incomplete matrices. *J. Mach. Learn. Res.* **11** 2287–2322. MR2719857

- 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](#)
- NETRAPALLI, P., JAIN, P. and SANGHAVI, S. (2015). Phase retrieval using alternating minimization. *IEEE Trans. Signal Process.* **63** 4814–4826. [MR3385838](#) <https://doi.org/10.1109/TSP.2015.2448516>
- NETRAPALLI, P., NIRANJAN, U., SANGHAVI, S., ANANDKUMAR, A. and JAIN, P. (2014). Non-convex robust PCA. In *Advances in Neural Information Processing Systems* 1107–1115.
- PEARSON, K. (1901). LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **2** 559–572.
- QIU, C. and VASWANI, N. (2010). Real-time robust principal components' pursuit. In *2010 48th Annual Allerton Conference on Communication, Control, and Computing (Allerton)* 591–598. IEEE.
- QIU, C., VASWANI, N., LOIS, B. and HOGBEN, L. (2014). Recursive robust PCA or recursive sparse recovery in large but structured noise. *IEEE Trans. Inform. Theory* **60** 5007–5039. [MR3245369](#) <https://doi.org/10.1109/TIT.2014.2331344>
- SHEN, Y., WEN, Z. and ZHANG, Y. (2014). Augmented Lagrangian alternating direction method for matrix separation based on low-rank factorization. *Optim. Methods Softw.* **29** 239–263. [MR3175484](#) <https://doi.org/10.1080/10556788.2012.700713>
- SINGER, A. (2011). Angular synchronization by eigenvectors and semidefinite programming. *Appl. Comput. Harmon. Anal.* **30** 20–36. [MR2737931](#) <https://doi.org/10.1016/j.acha.2010.02.001>
- SREBRO, N. and SHRAIBMAN, A. (2005). Rank, trace-norm and max-norm. In *Learning Theory. Lecture Notes in Computer Science* **3559** 545–560. Springer, Berlin. [MR2203286](#) [https://doi.org/10.1007/11503415\\_37](https://doi.org/10.1007/11503415_37)
- SUN, R. and LUO, Z.-Q. (2016). Guaranteed matrix completion via non-convex factorization. *IEEE Trans. Inform. Theory* **62** 6535–6579. [MR3565131](#) <https://doi.org/10.1109/TIT.2016.2598574>
- TAO, M. and YUAN, X. (2011). Recovering low-rank and sparse components of matrices from incomplete and noisy observations. *SIAM J. Optim.* **21** 57–81. [MR2765489](#) <https://doi.org/10.1137/100781894>
- VASWANI, N. and NARAYANAMURTHY, P. (2018). Static and dynamic robust PCA and matrix completion: A review. *Proceedings of the IEEE* **106** 1359–1379.
- VERSHYNIN, R. (2012). Introduction to the non-asymptotic analysis of random matrices. In *Compressed Sensing* 210–268. Cambridge Univ. Press, Cambridge. [MR2963170](#)
- WANG, G., GIANNAKIS, G. B. and EL DAR, Y. C. (2018). Solving systems of random quadratic equations via truncated amplitude flow. *IEEE Trans. Inform. Theory* **64** 773–794. [MR3762591](#) <https://doi.org/10.1109/TIT.2017.2756858>
- 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>
- WONG, R. K. W. and LEE, T. C. M. (2017). Matrix completion with noisy entries and outliers. *J. Mach. Learn. Res.* **18** Paper No. 147, 25. [MR3763781](#)
- YI, X., PARK, D., CHEN, Y. and CARAMANIS, C. (2016). Fast algorithms for robust PCA via gradient descent. In *NIPS* 4152–4160.
- ZHAN, J., LOIS, B., GUO, H. and VASWANI, N. (2016). Online (and offline) robust PCA: Novel algorithms and performance guarantees. In *Artificial Intelligence and Statistics* 1488–1496.
- ZHANG, H., CHI, Y. and LIANG, Y. (2016). Provable non-convex phase retrieval with outliers: Median truncated Wirtinger flow. In *International Conference on Machine Learning* 1022–1031.
- ZHANG, X., WANG, L. W. and GU, Q. (2018). A unified framework for nonconvex low-rank plus sparse matrix recovery. In *International Conference on Artificial Intelligence and Statistics*.
- ZHENG, Q. and LAFFERTY, J. (2016). Convergence analysis for rectangular matrix completion using Burer–Monteiro factorization and gradient descent. [arXiv:1605.07051](#).
- ZHOU, Z., LI, X., WRIGHT, J., CANDÈS, E. and MA, Y. (2010). Stable principal component pursuit. In *International Symposium on Information Theory* 1518–1522.