

ROP: MATRIX RECOVERY VIA RANK-ONE PROJECTIONS¹

BY T. TONY CAI AND ANRU ZHANG

University of Pennsylvania

Estimation of low-rank matrices is of significant interest in a range of contemporary applications. In this paper, we introduce a rank-one projection model for low-rank matrix recovery and propose a constrained nuclear norm minimization method for stable recovery of low-rank matrices in the noisy case. The procedure is adaptive to the rank and robust against small perturbations. Both upper and lower bounds for the estimation accuracy under the Frobenius norm loss are obtained. The proposed estimator is shown to be rate-optimal under certain conditions. The estimator is easy to implement via convex programming and performs well numerically.

The techniques and main results developed in the paper also have implications to other related statistical problems. An application to estimation of spiked covariance matrices from one-dimensional random projections is considered. The results demonstrate that it is still possible to accurately estimate the covariance matrix of a high-dimensional distribution based only on one-dimensional projections.

1. Introduction. Accurate recovery of low-rank matrices has a wide range of applications, including quantum state tomography [1, 24], face recognition [3, 12], recommender systems [27] and linear system identification and control [36]. For example, a key step in reconstructing the quantum states in low-rank quantum tomography is the estimation of a low-rank matrix based on Pauli measurements [24, 42]. And phase retrieval, a problem which arises in a range of signal and image processing applications including X-ray crystallography, astronomical imaging and diffraction imaging, can be reformulated as a low-rank matrix recovery problem [12, 15]. See Recht et al. [36] and Candès and Plan [13] for further references and discussions.

Received March 2014; revised August 2014.

¹Supported in part by NSF FRG Grant DMS-08-54973, NSF Grant DMS-12-08982 and NIH Grant R01 CA127334-05.

AMS 2000 subject classifications. Primary 62H12; secondary 62H36, 62C20.

Key words and phrases. Constrained nuclear norm minimization, low-rank matrix recovery, optimal rate of convergence, rank-one projection, restricted uniform boundedness, spiked covariance matrix.

<p>This is an electronic reprint of the original article published by the Institute of Mathematical Statistics in <i>The Annals of Statistics</i>, 2015, Vol. 43, No. 1, 102–138. This reprint differs from the original in pagination and typographic detail.</p>
--

Motivated by these applications, low-rank matrix estimation based on a small number of measurements has drawn much recent attention in several fields, including statistics, electrical engineering, applied mathematics and computer science. For example, Candès and Recht [14], Candès and Tao [16] and Recht [35] considered the exact recovery of a low-rank matrix based on a subset of uniformly sampled entries. Negahban and Wainwright [30] investigated matrix completion under a row/column weighted random sampling scheme. Recht et al. [36], Candès and Plan [13] and Cai and Zhang [8–10] studied matrix recovery based on a small number of linear measurements in the framework of Restricted Isometry Property (RIP), and Koltchinskii et al. [26] proposed the penalized nuclear norm minimization method and derived a general sharp oracle inequality under the condition of restricted isometry in expectation.

The basic model for low-rank matrix recovery can be written as

$$(1.1) \quad y = \mathcal{X}(A) + z,$$

where $\mathcal{X}: \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$ is a linear map, $A \in \mathbb{R}^{p_1 \times p_2}$ is an unknown low-rank matrix and z is a noise vector. The goal is to recover the low-rank matrix A based on the measurements (\mathcal{X}, y) . The linear map \mathcal{X} can be equivalently specified by n $p_1 \times p_2$ measurement matrices X_1, \dots, X_n with

$$(1.2) \quad \mathcal{X}(A) = (\langle X_1, A \rangle, \langle X_2, A \rangle, \dots, \langle X_n, A \rangle)^\top,$$

where the inner product of two matrices of the same dimensions is defined as $\langle X, Y \rangle = \sum_{i,j} X_{ij} Y_{ij}$. Since $\langle X, Y \rangle = \text{trace}(X^\top Y)$, (1.1) is also known as trace regression.

A common approach to low-rank matrix recovery is the constrained nuclear norm minimization method which estimates A by

$$(1.3) \quad \hat{A} = \arg \min_M \{ \|M\|_* : y - \mathcal{X}(M) \in \mathcal{Z} \}.$$

Here, $\|X\|_*$ is the nuclear norm of the matrix X which is defined to be the sum of its singular values, and \mathcal{Z} is a bounded set determined by the noise structure. For example, $\mathcal{Z} = \{0\}$ in the noiseless case and \mathcal{Z} is the feasible set of the error vector z in the case of bounded noise. This constrained nuclear norm minimization method has been well studied. See, for example, [8–10, 13, 31, 36].

Two random design models for low-rank matrix recovery have been particularly well studied in the literature. One is the so-called ‘‘Gaussian ensemble’’ [13, 36], where the measurement matrices X_1, \dots, X_n are random matrices with i.i.d. Gaussian entries. By exploiting the low-dimensional structure, the number of linear measurements can be far smaller than the number of entries in the matrix to ensure stable recovery. It has been shown that a matrix A of rank r can be stably recovered by nuclear norm minimization with high

probability, provided that $n \gtrsim r(p_1 + p_2)$ [13]. One major disadvantage of the Gaussian ensemble design is that it requires $O(np_1p_2)$ bytes of storage space for \mathcal{X} , which can be excessively large for the recovery of large matrices. For example, at least 45 TB of space is need to store the measurement matrices M_i in order to ensure accurate reconstruction of $10,000 \times 10,000$ matrices of rank 10. (See more discussion in Section 5.) Another popular design is the “matrix completion” model [14, 16, 35], under which the individual entries of the matrix A are observed at randomly selected positions. In terms of the measurement matrices X_i in (1.2), this can be interpreted as

$$(1.4) \quad \mathcal{X}(A) = (\langle e_{i_1} e_{j_1}^\top, A \rangle, \langle e_{i_2} e_{j_2}^\top, A \rangle, \dots, \langle e_{i_n} e_{j_n}^\top, A \rangle)^\top,$$

where $e_i = (0, \dots, 0, \overbrace{1}^{\text{ith}}, 0, \dots, 0)$ is the i th standard basis vector, and i_1, \dots, i_n and j_1, \dots, j_n are randomly and uniformly drawn with replacement from $\{1, \dots, p_1\}$ and $\{1, \dots, p_2\}$, respectively. However, as pointed out in [14, 35], additional structural assumptions, which are not intuitive and difficult to check, on the unknown matrix A are needed in order to ensure stable recovery under the matrix completion model. For example, it is impossible to recover spiked matrices under the matrix completion model. This can be easily seen from a simple example where the matrix A has only one nonzero row. In this case, although the matrix is only of rank one, it is not recoverable under the matrix completion model unless all the elements on the nonzero row are observed.

In this paper, we introduce a “*Rank-One Projection*” (ROP) model for low-rank matrix recovery and propose a constrained nuclear norm minimization method for this model. Under the ROP model, we observe

$$(1.5) \quad y_i = (\beta^{(i)})^\top A \gamma^{(i)} + z_i, \quad i = 1, \dots, n,$$

where $\beta^{(i)}$ and $\gamma^{(i)}$ are random vectors with entries independently drawn from some distribution \mathcal{P} , and z_i are random errors. In terms of the linear map $\mathcal{X}: \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$ in (1.1), it can be defined as

$$(1.6) \quad [\mathcal{X}(A)]_i = (\beta^{(i)})^\top A \gamma^{(i)}, \quad i = 1, \dots, n.$$

Since the measurement matrices $X_i = \beta^{(i)}(\gamma^{(i)})^\top$ are of rank-one, we call the model (1.5) a “*Rank-One Projection*” (ROP) model. It is easy to see that the storage for the measurement vectors in the ROP model (1.5) is $O(n(p_1 + p_2))$ bytes which is significantly smaller than $O(np_1p_2)$ bytes required for the Gaussian ensemble.

We first establish a sufficient identifiability condition in Section 2 by considering the problem of exact recovery of low-rank matrices in the noiseless case. It is shown that, with high probability, ROP with $n \gtrsim r(p_1 + p_2)$

random projections is sufficient to ensure exact recovery of all rank- r matrices through the constrained nuclear norm minimization. The required number of measurements $O(r(p_1 + p_2))$ is rate optimal for any linear measurement model since a rank- r matrix $A \in \mathbb{R}^{p_1+p_2}$ has the degree of freedom $r(p_1 + p_2 - r)$. The Gaussian noise case is of particular interest in statistics. We propose a new constrained nuclear norm minimization estimator and investigate its theoretical and numerical properties in the Gaussian noise case. Both upper and lower bounds for the estimation accuracy under the Frobenius norm loss are obtained. The estimator is shown to be rate-optimal when the number of rank-one projections satisfies either $n \gtrsim (p_1 + p_2) \log(p_1 + p_2)$ or $n \sim r(p_1 + p_2)$. The lower bound also shows that if the number of measurements $n < r \max(p_1, p_2)$, then no estimator can recover rank- r matrices consistently. The general case where the matrix A is only approximately low-rank is also considered. The results show that the proposed estimator is adaptive to the rank r and robust against small perturbations. Extensions to the sub-Gaussian design and sub-Gaussian noise distribution are also considered.

The ROP model can be further simplified by taking $\beta^{(i)} = \gamma^{(i)}$ if the low-rank matrix A is known to be symmetric. This is the case in many applications, including low-dimensional Euclidean embedding [36, 38], phase retrieval [12, 15] and covariance matrix estimation [5, 6, 17]. In such a setting, the ROP design can be simplified to Symmetric Rank-One Projections (SROP)

$$[\mathcal{X}(A)]_i = (\beta^{(i)})^\top A \beta^{(i)}.$$

We will show that the results for the general ROP model continue to hold for the SROP model when A is known to be symmetric. Recovery of symmetric positive definite matrices in the noiseless and ℓ_1 -bounded noise settings has also been considered in a recent paper by Chen et al. [17] which was posted on arXiv at the time of the writing of the present paper. Their results and techniques for symmetric positive definite matrices are not applicable to the recovery of general low-rank matrices. See Section 6 for more discussions.

The techniques and main results developed in the paper also have implications to other related statistical problems. In particular, the results imply that it is possible to accurately estimate a spiked covariance matrix based only on one-dimensional projections. Spiked covariance matrix model has been well studied in the context of Principal Component Analysis (PCA) based on i.i.d. data where one observes p -dimensional vectors $X^{(1)}, \dots, X^{(n)} \stackrel{\text{i.i.d.}}{\sim} N(0, \Sigma)$ with $\Sigma = I_p + \Sigma_0$ and Σ_0 being low-rank [4–6, 25]. This covariance structure and its variations have been used in many applications including signal processing, financial econometrics, chemometrics and population genetics. See, for example, [21, 29, 33, 34, 43]. Suppose that

the random vectors $X^{(1)}, \dots, X^{(n)}$ are not directly observable. Instead, we observe only one-dimensional random projections of $X^{(i)}$,

$$\xi_i = \langle \beta^{(i)}, X^{(i)} \rangle, \quad i = 1, \dots, n,$$

where $\beta^{(i)} \stackrel{\text{i.i.d.}}{\sim} N(0, I_p)$. It is somewhat surprising that it is still possible to accurately estimate the spiked covariance matrix Σ based only on the one-dimensional projections $\{\xi_i : i = 1, \dots, n\}$. This covariance matrix recovery problem is also related to the recent literature on covariance sketching [18, 19], which aims to recover a symmetric matrix A (or a general rectangular matrix B) from low-dimensional projections of the form $X^\top AX$ (or $X^\top BY$). See Section 4 for further discussions.

The proposed methods can be efficiently implemented via convex programming. A simulation study is carried out to investigate the numerical performance of the proposed nuclear norm minimization estimators. The numerical results indicate that ROP with $n \geq 5r \max(p_1, p_2)$ random projections is sufficient to ensure the exact recovery of rank- r matrices through constrained nuclear norm minimization and show that the procedure is robust against small perturbations, which confirm the theoretical results developed in the paper. The proposed estimator outperforms two other alternative procedures numerically in the noisy case. In addition, the proposed method is illustrated through an image compression example.

The rest of the paper is organized as follows. In Section 2, after introducing basic notation and definitions, we consider exact recovery of low-rank matrices in the noiseless case and establish a sufficient identifiability condition. A constrained nuclear norm minimization estimator is introduced for the Gaussian noise case. Both upper and lower bounds are obtained for estimation under the Frobenius norm loss. Section 3 considers extensions to sub-Gaussian design and sub-Gaussian noise distributions. An application to estimation of spiked covariance matrices based on one-dimensional projections is discussed in detail in Section 4. Section 5 investigates the numerical performance of the proposed procedure through a simulation study and an image compression example. A brief discussion is given in Section 6. The main results are proved in Section 7 and the proofs of some technical lemmas are given in the supplementary material [11].

2. Matrix recovery under Gaussian noise. In this section, we first establish an identifiability condition for the ROP model by considering exact recovery in the noiseless case, and then focus on low-rank matrix recovery in the Gaussian noise case.

We begin with the basic notation and definitions. For a vector $\beta \in \mathbb{R}^n$, we use $\|\beta\|_q = \sqrt[q]{\sum_{i=1}^n |\beta_i|^q}$ to define its vector q -norm. For a matrix $X \in \mathbb{R}^{p_1 \times p_2}$, the Frobenius norm is $\|X\|_F = \sqrt{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} X_{ij}^2}$ and the spectral

norm $\|\cdot\|$ is $\|X\| = \sup_{\|\beta\|_2 \leq 1} \|X\beta\|_2$. For a linear map $\mathcal{X} = (X_1, \dots, X_n)$ from $\mathbb{R}^{p_1 \times p_2}$ to \mathbb{R}^n given by (1.2), its dual operator $\mathcal{X}^*: \mathbb{R}^n \rightarrow \mathbb{R}^{p_1 \times p_2}$ is defined as $\mathcal{X}^*(z) = \sum_{i=1}^n z_i X_i$. For a matrix $X \in \mathbb{R}^{p_1 \times p_2}$, let $X = \sum_i a_i u_i v_i^\top$ be the singular value decomposition of X with the singular values $a_1 \geq a_2 \geq \dots \geq 0$. We define $X_{\max(r)} = \sum_{i=1}^r a_i u_i v_i^\top$ and $X_{-\max(r)} = X - X_{\max(r)} = \sum_{i \geq r+1} a_i u_i v_i^\top$. For any two sequences $\{a_n\}$ and $\{b_n\}$ of positive numbers, denote by $a_n \gtrsim b_n$ when $a_n \geq C b_n$ for some uniform constant C and denote by $a_n \sim b_n$ if $a_n \gtrsim b_n$ and $b_n \gtrsim a_n$.

We use the phrase ‘‘rank- r matrices’’ to refer to matrices of rank at most r and denote by \mathbb{S}^p the set of all $p \times p$ symmetric matrices. A linear map $\mathcal{X}: \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$ is called ROP from distribution \mathcal{P} if \mathcal{X} is defined as in (1.6) with all the entries of $\beta^{(i)}$ and $\gamma^{(i)}$ independently drawn from the distribution \mathcal{P} .

2.1. *RUB, identifiability, and exact recovery in the noiseless case.* An important step toward understanding the constrained nuclear norm minimization is the study of exact recovery of low-rank matrices in the noiseless case which also leads to a sufficient identifiability condition. A widely used framework in the low-rank matrix recovery literature is the Restricted Isometry Property (RIP) in the matrix setting. See [8–10, 13, 36, 37]. However, the RIP framework is not well suited for the ROP model and would lead to suboptimal results. See Section 2.2 for more discussions on the RIP and other conditions used in the literature. See also [15]. In this section, we introduce a Restricted Uniform Boundedness (RUB) condition which will be shown to guarantee the exact recovery of low-rank matrices in the noiseless case and stable recovery in the noisy case through the constrained nuclear norm minimization. It will also be shown that the RUB condition are satisfied by a range of random linear maps with high probability.

DEFINITION 2.1 (Restricted Uniform Boundedness). For a linear map $\mathcal{X}: \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$, if there exist uniform constants C_1 and C_2 such that for all nonzero rank- r matrices $A \in \mathbb{R}^{p_1 \times p_2}$

$$C_1 \leq \frac{\|\mathcal{X}(A)\|_1/n}{\|A\|_F} \leq C_2,$$

where $\|\cdot\|_1$ means the vector ℓ_1 norm, then we say that \mathcal{X} satisfies the Restricted Uniform Boundedness (RUB) condition of order r and constants C_1 and C_2 .

In the noiseless case, we observe $y = \mathcal{X}(A)$ and estimate the matrix A through the constrained nuclear norm minimization

$$(2.1) \quad A_* = \arg \min_M \{\|M\|_* : \mathcal{X}(M) = y\}.$$

The following theorem shows that the RUB condition guarantees the exact recovery of all rank- r matrices.

THEOREM 2.1. *Let $k \geq 2$ be an integer. Suppose \mathcal{X} satisfies RUB of order kr with $C_2/C_1 < \sqrt{k}$, then the nuclear norm minimization method recovers all rank- r matrices. That is, for all rank- r matrices A and $y = \mathcal{X}(A)$, we have $A_* = A$, where A_* is given by (2.1).*

Theorem 2.1 shows that RUB of order kr with $C_2/C_1 < \sqrt{k}$ is a sufficient identifiability condition for the low-rank matrix recovery model (1.1) in the noisy case. The following result shows that the RUB condition is satisfied with high probability under the ROP model with a sufficient number of measurements.

THEOREM 2.2. *Suppose $\mathcal{X} : \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$ is ROP from the standard normal distribution. For integer $k \geq 2$, positive numbers $C_1 < \frac{1}{3}$ and $C_2 > 1$, there exist constants C and δ , not depending on p_1, p_2 and r , such that if*

$$(2.2) \quad n \geq Cr(p_1 + p_2),$$

then with probability at least $1 - e^{-n\delta}$, \mathcal{X} satisfies RUB of order kr and constants C_1 and C_2 .

REMARK 2.1. The condition $n \geq O(r(p_1 + p_2))$ on the number of measurements is indeed necessary for \mathcal{X} to satisfy nontrivial RUB with $C_1 > 0$. Note that the degree of freedom of all rank- r matrices of $\mathbb{R}^{p_1 \times p_2}$ is $r(p_1 + p_2 - r) \geq \frac{1}{2}r(p_1 + p_2)$. If $n < \frac{1}{2}r(p_1 + p_2)$, there must exist a nonzero rank- r matrix $A \in \mathbb{R}^{p_1 \times p_2}$ such that $\mathcal{X}(A) = 0$, which leads to the failure of any nontrivial RUB for \mathcal{X} .

As a direct consequence of Theorems 2.1 and 2.2, ROP with the number of measurements $n \geq Cr(p_1 + p_2)$ guarantees the exact recovery of all rank- r matrices with high probability.

COROLLARY 2.1. *Suppose $\mathcal{X} : \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$ is ROP from the standard normal distribution. There exist uniform constants C and δ such that, whenever $n \geq Cr(p_1 + p_2)$, the nuclear norm minimization estimator A_* given in (2.1) recovers all rank- r matrices $A \in \mathbb{R}^{p_1 \times p_2}$ exactly with probability at least $1 - e^{-n\delta}$.*

Note that the required number of measurements $O(r(p_1 + p_2))$ above is rate optimal, since the degree of freedom for a matrix $A \in \mathbb{R}^{p_1 + p_2}$ of rank r is $r(p_1 + p_2 - r)$, and thus at least $r(p_1 + p_2 - r)$ measurements are needed in order to recover A exactly using any method.

2.2. RUB, RIP and other conditions. We have shown that RUB implies exact recovery in the noiseless and proved that the random rank-one projections satisfy RUB with high probability whenever the number of measurements $n \geq Cr(p_1 + p_2)$. As mentioned earlier, other conditions, including the Restricted Isometry Property (RIP), RIP in expectation and Spherical Section Property (SSP), have been introduced for low-rank matrix recovery based on linear measurements. Among them, RIP is perhaps the most widely used. A linear map $\mathcal{X}: \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$ is said to satisfy RIP of order r with positive constants C_1 and C_2 if

$$C_1 \leq \frac{\|\mathcal{X}(A)\|_2/\sqrt{n}}{\|A\|_F} \leq C_2$$

for all rank- r matrices A . Many results have been given for low-rank matrices under the RIP framework. For example, Recht et al. [36] showed that Gaussian ensembles satisfy RIP with high probability under certain conditions on the dimensions. Candès and Plan [13] provided a lower bound and oracle inequality under the RIP condition. Cai and Zhang [8–10] established the sharp bounds for the RIP conditions that guarantee accurate recovery of low-rank matrices.

However, the RIP framework is not suitable for the ROP model considered in the present paper. The following lemma is proved in the supplementary material [11].

LEMMA 2.1. *Suppose $\mathcal{X}: \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$ is ROP from the standard normal distribution. Let*

$$C_1 = \min_{A: \text{rank}(A)=1} \frac{\|\mathcal{X}(A)\|_2/\sqrt{n}}{\|A\|_F} \quad \text{and} \quad C_2 = \max_{A: \text{rank}(A)=1} \frac{\|\mathcal{X}(A)\|_2/\sqrt{n}}{\|A\|_F}.$$

Then for all $t > 1$, $C_2/C_1 \geq \sqrt{p_1 p_2 / (4tn)}$ with probability at least $1 - e^{-p_1/4} - e^{-p_2/4} - \frac{8}{n(t-1)^2}$.

Lemma 2.1 implies that at least $O(p_1 p_2)$ number of measurements are needed in order to ensure that \mathcal{X} satisfies the RIP condition that guarantees the recovery of only rank-one matrices. Since $O(p_1 p_2)$ is the degree of freedom for all matrices $A \in \mathbb{R}^{p_1 \times p_2}$ and it is the number of measurements needed to recover all $p_1 \times p_2$ matrices (not just the low-rank matrices), Lemma 2.1 shows that the RIP framework is not suitable for the ROP model. In comparison, Theorem 2.2 shows that if $n \geq O(r(p_1 + p_2))$, then with high probability \mathcal{X} satisfies the RUB condition of order r with bounded C_2/C_1 , which ensures the exact recovery of all rank- r matrices.

The main technical reason for the failure of RIP under the ROP model is that RIP requires an upper bound for

$$(2.3) \quad \max_{A \in \mathcal{C}} \|\mathcal{X}(A)\|_2^2/n = \max_{A \in \mathcal{C}} \left(\sum_{j=1}^n ((\beta^{(j)})^\top A \gamma^{(j)})^2 \right) / n,$$

where \mathcal{C} is a set containing low-rank matrices. The right-hand side of (2.3) involves the 4th power of the Gaussian (or sub-Gaussian) variables $\beta^{(j)}$ and $\gamma^{(j)}$. A much larger n than the bound given in (2.2) is needed in order for the linear map \mathcal{X} to satisfy the required RIP condition, which would lead to suboptimal result.

Koltchinskii et al. [26] uses RIP in expectation, which is a weaker condition than RIP. A random linear map $\mathcal{X}: \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$ is said to satisfy RIP in expectation of order r with parameters $0 < \mu < \infty$ and $0 \leq \delta_r < 1$ if

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

for all rank- r matrices $A \in \mathbb{R}^{p_1 \times p_2}$. This condition was originally introduced by Koltchinskii et al. [26] to prove an oracle inequality for the estimator they proposed and a minimax lower bound. The condition is not sufficiently strong to guarantee the exact recovery of rank- r matrices in the noiseless case. To be more specific, the bounds in Theorems 1 and 2 in [26] depend on $\mathbf{M} = \|\frac{1}{n} \sum_{i=1}^n (y_i X_i - E(y_i X_i))\|$, which might be nonzero even in the noiseless case. In fact, in the ROP model considered in the present paper, we have

$$\begin{aligned} \frac{1}{n} E \|\mathcal{X}\|_2^2 &= \frac{1}{n} \sum_{i=1}^n E(\beta^{(i)T} A \gamma^{(i)})^2 = E(\beta^\top A \gamma \gamma^\top A^\top \beta) \\ &= E \operatorname{tr}(A \gamma \gamma^\top A^\top \beta \beta^\top) = \operatorname{tr}(A A^\top) = \|A\|_F^2 \end{aligned}$$

which means RIP in expectation is met for $\mu = 1$ and $\delta_r = 0$ for any number of measurements n . However, as we discussed earlier in this section that at least $O(r(p_1 + p_2))$ measurements are needed to guarantee the model identifiability for recovery of all rank- r matrices, we can see that RIP in expectation cannot ensure recovery.

Dvijotham and Fazel [20] and Oymak et al. [32] used a condition called the Spherical Section Property (SSP) which focuses on the null space of \mathcal{X} . $\operatorname{Null}(\mathcal{X})$ is said to satisfy Δ -SSP if for all $Z \in \operatorname{Null}(\mathcal{X}) \setminus \{0\}$, $\|Z\|_* / \|Z\|_F \geq \sqrt{\Delta}$. Dvijotham and Fazel [20] showed that if \mathcal{X} satisfies Δ -SSP, $p_1 \leq p_2$ and $\operatorname{rank}(A) < \min(3p_1/4 - \sqrt{9p_1^2/16 - p_1\Delta/4}, p_1/2)$, the nuclear norm minimization (2.1) recovers A exactly in the noiseless case. However, the SSP condition is difficult to utilize in the ROP framework since it is hard to characterize the matrices $Z \in \operatorname{Null}(\mathcal{X})$ when \mathcal{X} is rank-one projections.

2.3. Gaussian noise case. We now turn to the Gaussian noise case where $z_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$ in (1.5). We begin by introducing a constrained nuclear norm minimization estimator. Define two sets

$$(2.4) \quad \mathcal{Z}_1 = \{z : \|z\|_1/n \leq \sigma\} \quad \text{and} \quad \mathcal{Z}_2 = \{z : \|\mathcal{X}^*(z)\| \leq \eta\},$$

where $\eta = \sigma(12\sqrt{\log n}(p_1 + p_2) + 6\sqrt{2n(p_1 + p_2)})$, and let

$$(2.5) \quad \mathcal{Z}_G = \mathcal{Z}_1 \cap \mathcal{Z}_2.$$

Note that both \mathcal{Z}_1 and \mathcal{Z}_2 are convex sets and so is \mathcal{Z}_G . Our estimator of A is given by

$$(2.6) \quad \hat{A} = \arg \min_M \{\|M\|_* : y - \mathcal{X}(M) \in \mathcal{Z}_G\}.$$

The following theorem gives the rate of convergence for the estimator \hat{A} under the squared Frobenius norm loss.

THEOREM 2.3 (Upper bound). *Let \mathcal{X} be ROP from the standard normal distribution and let $z_1, \dots, z_n \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$. Then there exist uniform constants C, W and δ such that, whenever $n \geq Cr(p_1 + p_2)$, the estimator \hat{A} given in (2.6) satisfies*

$$(2.7) \quad \|\hat{A} - A\|_F^2 \leq W\sigma^2 \min\left(\frac{r \log n(p_1 + p_2)^2}{n^2} + \frac{r(p_1 + p_2)}{n}, 1\right)$$

for all rank- r matrices A , with probability at least $1 - 11/n - 3\exp(-\delta(p_1 + p_2))$.

Moreover, we have the following lower bound result for ROP.

THEOREM 2.4 (Lower bound). *Assume that \mathcal{X} is ROP from the standard normal distribution and that $z_1, \dots, z_n \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$. There exists a uniform constant C such that, when $n > Cr \max(p_1, p_2)$, with probability at least $1 - 26n^{-1}$,*

$$(2.8) \quad \inf_{\hat{A}} \sup_{A \in \mathbb{R}^{p_1 \times p_2} : \text{rank}(A)=r} P_z \left(\|\hat{A} - A\|_F^2 \geq \frac{\sigma^2 r(p_1 + p_2)}{32n} \right) \geq 1 - e^{-(p_1 + p_2)r/64},$$

$$(2.9) \quad \inf_{\hat{A}} \sup_{A \in \mathbb{R}^{p_1 \times p_2} : \text{rank}(A)=r} E_z \|\hat{A} - A\|_F^2 \geq \frac{\sigma^2 r(p_1 + p_2)}{4n},$$

where E_z , and P_z are the expectation and probability with respect to the distribution of z .

When $n < r \max(p_1, p_2)$, then

$$(2.10) \quad \inf_{\hat{A}} \sup_{A \in \mathbb{R}^{p_1 \times p_2} : \text{rank}(A)=r} E_z \|\hat{A} - A\|_F^2 = \infty.$$

Comparing Theorems 2.3 and 2.4, our proposed estimator is rate optimal in the Gaussian noise case when $n \gtrsim \log n(p_1 + p_2)$ [which is equivalent to $n \gtrsim (p_1 + p_2) \log(p_1 + p_2)$] or $n \sim r(p_1 + p_2)$. Since $n \gtrsim r(p_1 + p_2)$, this condition is also implied by $r \gtrsim \log(p_1 + p_2)$. Theorem 2.4 also shows that no method can recover matrices of rank r consistently if the number of measurements n is smaller than $r \max(p_1, p_2)$.

The result in Theorem 2.3 can also be extended to the more general case where the matrix of interest A is only approximately low-rank. Let $A = A_{\max(r)} + A_{-\max(r)}$.

PROPOSITION 2.1. *Under the assumptions of Theorem 2.3, there exist uniform constants C , W_1 , W_2 and δ such that, whenever $n \geq Cr(p_1 + p_2)$, the estimator \hat{A} given in (2.6) satisfies*

$$(2.11) \quad \begin{aligned} \|\hat{A} - A\|_F^2 &\leq W_1 \sigma^2 \min\left(\frac{r \log n(p_1 + p_2)^2}{n^2} + \frac{r(p_1 + p_2)}{n}, 1\right) \\ &\quad + W_2 \frac{\|A_{-\max(r)}\|_*^2}{r} \end{aligned}$$

for all matrices $A \in \mathbb{R}^{p_1 \times p_2}$, with probability at least $1 - 11/n - 3 \exp(-\delta(p_1 + p_2))$.

If the matrix A is approximately of rank r , then $\|A_{-\max(r)}\|_*$ is small, and the estimator \hat{A} continues to perform well. This result shows that the constrained nuclear norm minimization estimator is adaptive to the rank r and robust against perturbations of small amplitude.

REMARK 2.2. All the results remain true if the Gaussian design is replaced by the Rademacher design where entries of $\beta^{(i)}$ and $\gamma^{(i)}$ are i.i.d. ± 1 with probability $\frac{1}{2}$. More general sub-Gaussian design case will be discussed in Section 3.

REMARK 2.3. The estimator \hat{A} we propose here is the minimizer of the nuclear norm under the constraint of the intersection of two convex sets \mathcal{Z}_1 and \mathcal{Z}_2 . Nuclear norm minimization under either one of the two constraints, called “ ℓ_1 constraint nuclear norm minimization” ($\mathcal{Z} = \mathcal{Z}_1$) and “matrix Dantzig Selector” ($\mathcal{Z} = \mathcal{Z}_2$), has been studied before in various settings [8–10, 13, 17, 36]. Our analysis indicates the following:

1. The ℓ_1 constraint minimization performs better than the matrix Dantzig Selector for small n ($n \sim r(p_1 + p_2)$) when $r \ll \log n$.
2. The matrix Dantzig Selector outperforms the ℓ_1 constraint minimization for large n as the loss of the matrix Dantzig Selector decays at the rate $O(n^{-1})$.
3. The proposed estimator \hat{A} combines the advantages of the two estimators.

See Section 5 for a comparison of numerical performances of the three methods.

2.4. Recovery of symmetric matrices. For applications such as low-dimensional Euclidean embedding [36, 38], phase retrieval [12, 15] and covariance matrix estimation [5, 6, 17], the low-rank matrix A of interest is known to be symmetric. Examples of such matrices include distance matrices, Gram matrices, and covariance matrices. When the matrix A is known to be symmetric, the ROP design can be further simplified by taking $\beta^{(i)} = \gamma^{(i)}$.

Denote by \mathbb{S}^p the set of all $p \times p$ symmetric matrices in $\mathbb{R}^{p \times p}$. Let $\beta^{(1)}, \beta^{(2)}, \dots, \beta^{(n)}$ be independent p -dimensional random vectors with i.i.d. entries generated from some distribution \mathcal{P} . Define a linear map $\mathcal{X}: \mathbb{S}^p \rightarrow \mathbb{R}^n$ by

$$[\mathcal{X}(A)]_i = (\beta^{(i)})^\top A \beta^{(i)}, \quad i = 1, \dots, n.$$

We call such a linear map \mathcal{X} ‘‘Symmetric Rank-One Projections’’ (SROP) from the distribution \mathcal{P} .

Suppose we observe

$$(2.12) \quad y_i = (\beta^{(i)})^\top A \beta^{(i)} + z_i, \quad i = 1, \dots, n$$

and wish to recover the symmetric matrix A . As for the ROP model, in the noiseless case we estimate A under the SROP model by

$$(2.13) \quad A_* = \arg \min_{M \in \mathbb{S}^p} \{\|M\|_* : y = \mathcal{X}(M)\}.$$

PROPOSITION 2.2. *Let \mathcal{X} be SROP from the standard normal distribution. Similar to Corollary 2.1, there exist uniform constants C and δ such that, whenever $n \geq Crp$, the nuclear norm minimization estimator A_* given by (2.13) recovers exactly all rank- r symmetric matrices $A \in \mathbb{S}^p$ with probability at least $1 - e^{-n\delta}$.*

For the noisy case, we propose a constraint nuclear norm minimization estimator similar to (2.6). Define the linear map $\tilde{\mathcal{X}}: \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^{\lfloor n/2 \rfloor}$ by

$$(2.14) \quad [\tilde{\mathcal{X}}(A)]_i = [\mathcal{X}(A)]_{2i-1} - [\mathcal{X}(A)]_{2i}, \quad i = 1, \dots, \left\lfloor \frac{n}{2} \right\rfloor$$

and define $\tilde{y} \in \mathbb{R}^{\lfloor n/2 \rfloor}$ by

$$(2.15) \quad \tilde{y}_i = y_{2i-1} - y_{2i}, \quad i = 1, \dots, \left\lfloor \frac{n}{2} \right\rfloor.$$

Based on the definition of $\tilde{\mathcal{X}}$, the dual map $\tilde{\mathcal{X}}^* : \mathbb{R}^{\lfloor n/2 \rfloor} \rightarrow \mathbb{S}^p$ is

$$(2.16) \quad \tilde{\mathcal{X}}^*(z) = \sum_{i=1}^{\lfloor n/2 \rfloor} z_i (\beta^{(2i-1)} \beta^{(2i-1)\top} - \beta^{(2i)} \beta^{(2i)\top}).$$

Let $\eta = 24\sigma(\sqrt{pn} + 2p\sqrt{2\log n})$. The estimator \hat{A} of the matrix A is given by

$$(2.17) \quad \hat{A} = \arg \min_{M \in \mathbb{S}^p} \{ \|M\|_* : \|y - \mathcal{X}(M)\|_1/n \leq \sigma, \|\tilde{\mathcal{X}}^*(\tilde{y} - \tilde{\mathcal{X}}(M))\| \leq \eta \}.$$

REMARK 2.4. An important property in the ROP model considered in Section 2.3 is that $E\mathcal{X} = 0$, that is, $EX_i = 0$ for all the measurement matrices X_i . However, under the SROP model $X_i = \beta^{(i)}(\beta^{(i)})^\top$ and so $E\mathcal{X} \neq 0$. The step of taking the pairwise differences in (2.14) and (2.15) is to ensure that $E\tilde{\mathcal{X}} = 0$.

The following result is similar to the upper bound given in Proposition 2.1 for ROP.

PROPOSITION 2.3. *Let \mathcal{X} be SROP from the standard normal distribution and let $z_1, \dots, z_n \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$. There exist constants C, W_1, W_2 and δ such that, whenever $n \geq Crp$, the estimator \hat{A} given in (2.17) satisfies*

$$(2.18) \quad \|\hat{A} - A\|_F^2 \leq W_1 \sigma^2 \min\left(\frac{rp^2 \log n}{n^2} + \frac{rp}{n}, 1\right) + W_2 \frac{\|A_{-\max(r)}\|_*^2}{r}$$

for all matrices $A \in \mathbb{S}^p$, with probability at least $1 - 15/n - 5\exp(-p\delta)$.

In addition, we also have lower bounds for SROP, which show that the proposed estimator is rate-optimal when $n \gtrsim p \log n$ or $n \sim rp$, and no estimator can recover a rank- r matrix consistently if the number of measurements $n < \lfloor \frac{r}{2} \rfloor \cdot \lfloor \frac{p}{2} \rfloor$.

PROPOSITION 2.4 (Lower bound). *Assume that \mathcal{X} is SROP from the standard normal distribution and that $z_1, \dots, z_n \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$. Then there exists a uniform constant C such that, when $n > Crp$ and $p, r \geq 2$, with*

probability at least $1 - 26n^{-1}$,

$$\inf_{\hat{A}} \sup_{A \in \mathbb{S}^p : \text{rank}(A)=r} P_z \left(\|\hat{A} - A\|_F^2 \geq \frac{\sigma^2 r p}{192n} \right) \geq 1 - e^{-pr/192},$$

$$\inf_{\hat{A}} \sup_{A \in \mathbb{S}^p : \text{rank}(A)=r} E_z \|\hat{A} - A\|_F^2 \geq \frac{\sigma^2 r p}{24n},$$

where \hat{A} is any estimator of A , E_z, P_z are the expectation and probability with respect to z .

When $n < \lfloor \frac{r}{2} \rfloor \cdot \lfloor \frac{p}{2} \rfloor$ and $p, r \geq 2$, then

$$\inf_{\hat{A}} \sup_{A \in \mathbb{S}^p : \text{rank}(A)=r} E_z \|\hat{A} - A\|_F^2 = \infty.$$

3. Sub-Gaussian design and sub-Gaussian noise. We have focused on the Gaussian design and Gaussian noise distribution in Section 2. These results can be further extended to more general distributions. In this section, we consider the case where the ROP design is from a symmetric sub-Gaussian distribution \mathcal{P} and the errors z_i are also from a sub-Gaussian distribution. We say the distribution of a random variable Z is sub-Gaussian with parameter τ if

$$(3.1) \quad P(|Z| \geq t) \leq 2 \exp(-t^2/(2\tau^2)) \quad \text{for all } t > 0.$$

The following lemma provides a necessary and sufficient condition for symmetric sub-Gaussian distributions.

LEMMA 3.1. *Let \mathcal{P} be a symmetric distribution and let the random variable $X \sim \mathcal{P}$. Define*

$$(3.2) \quad \alpha_{\mathcal{P}} = \sup_{k \geq 1} \left(\frac{EX^{2k}}{(2k-1)!!} \right)^{1/2k}.$$

Then the distribution \mathcal{P} is sub-Gaussian if and only if $\alpha_{\mathcal{P}}$ is finite.

For the sub-Gaussian ROP design and sub-Gaussian noise, we estimate the low-rank matrix A by the estimator \hat{A} given in (1.3) with

$$(3.3) \quad \begin{aligned} \mathcal{Z}_G &= \{z : \|z\|_1/n \leq 6\tau\} \\ &\cap \{z : \|\mathcal{X}^*(z)\| \leq 6\alpha_{\mathcal{P}}^2 \tau (\sqrt{6n(p_1 + p_2)} + 2\sqrt{\log n(p_1 + p_2)})\}, \end{aligned}$$

where $\alpha_{\mathcal{P}}$ is given in (3.2).

THEOREM 3.1. *Suppose $\mathcal{X} : \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$ is ROP from a symmetric and variance 1 sub-Gaussian distribution \mathcal{P} . Assume that z_i are i.i.d. sub-Gaussian with parameter τ and \hat{A} is given by (1.3) with $\mathcal{Z} = \mathcal{Z}_G$ defined in (3.3). Then there exist constants C, W_1, W_2, δ which only depend on \mathcal{P} , such that if $n \geq Cr(p_1 + p_2)$, we have*

$$(3.4) \quad \begin{aligned} \|\hat{A} - A\|_F^2 &\leq W_1 \tau^2 \min\left(\frac{r \log n (p_1 + p_2)^2}{n^2} + \frac{r(p_1 + p_2)}{n}, 1\right) \\ &\quad + W_2 \frac{\|A_{-\max(r)}\|_*^2}{r} \end{aligned}$$

with probability at least $1 - 2/n - 5e^{-\delta(p_1 + p_2)}$.

An exact recovery result in the noiseless case for the sub-Gaussian design follows directly from Theorem 3.1. If $z = 0$, then, with high probability, all rank- r matrices A can be recovered exactly via the constrained nuclear minimization (2.1) whenever $n \geq C_{\mathcal{P}} r (p_1 + p_2)$ for some constant $C_{\mathcal{P}} > 0$.

REMARK 3.1. For the SROP model considered in Section 2.4, we can similarly extend the results to the case of sub-Gaussian design and sub-Gaussian noise. Suppose \mathcal{X} is SROP from a symmetric variance 1 sub-Gaussian distribution \mathcal{P} (other than the Rademacher ± 1 distribution) and z satisfies (3.1). Define the estimator of the low-rank matrix A by

$$(3.5) \quad \hat{A} = \arg \min_{M \in \mathbb{S}^p} \{\|M\|_* : \|y - \mathcal{X}(M)\|_1/n \leq 6\tau, \|\tilde{\mathcal{X}}^*(\tilde{y} - \tilde{\mathcal{X}}(M))\| \leq \eta\},$$

where $\eta = C_{\mathcal{P}}(\sqrt{np} + \sqrt{\log np})$ with $C_{\mathcal{P}}$ some constant depending on \mathcal{P} .

PROPOSITION 3.1. *Suppose $\mathcal{X} : \mathbb{R}^{p \times p} \rightarrow \mathbb{R}^n$ is SROP from a symmetric sub-Gaussian distribution \mathcal{P} with variance 1. Also, assume that $\text{Var}(\mathcal{P}^2) > 0$ [i.e., $\text{Var}(w^2) > 0$ where $w \sim \mathcal{P}$]. Let \hat{A} be given by (3.5). Then there exist constants $C, C_{\mathcal{P}}, W_1, W_2$ and δ which only depend on \mathcal{P} , such that for $n \geq Crp$,*

$$(3.6) \quad \|\hat{A} - A\|_F^2 \leq W_1 \tau^2 \min\left(\frac{rp^2 \log n}{n^2} + \frac{rp}{n}, 1\right) + W_2 \frac{\|A_{-\max(r)}\|_*^2}{r}$$

with probability at least $1 - 2/n - 5e^{-\delta p}$.

By restricting $\text{Var}(\mathcal{P}^2) > 0$, Rademacher ± 1 is the only symmetric and variance 1 distribution that has been excluded. The reason why the Rademacher ± 1 distribution is an exception for the SROP design is as follows. If $\beta^{(i)}$ are i.i.d. Rademacher ± 1 distributed, then

$$[\mathcal{X}(A)]_i = (\beta^{(i)})^\top A \beta^{(i)} = \sum_{j=1}^p a_{jj} + \sum_{j \neq k} \beta_j^{(i)} \beta_k^{(i)} a_{jk}, \quad i = 1, \dots, n.$$

So the only information contained in $\mathcal{X}(A)$ about $\text{diag}(A)$ is $\text{trace}(A)$, which makes it impossible to recover the whole matrix A .

4. Application to estimation of spiked covariance matrix. In this section, we consider an interesting application of the methods and results developed in the previous sections to estimation of a spiked covariance matrix based on one-dimensional projections. As mentioned in the [Introduction](#), spiked covariance matrix model has been used in a wide range of applications and it has been well studied in the context of PCA based on i.i.d. data where one observes i.i.d. p -dimensional random vectors $X^{(1)}, \dots, X^{(n)}$ with mean 0 and covariance matrix Σ , where $\Sigma = I_p + \Sigma_0$ and Σ_0 being low-rank. See, for example, [\[4–6, 25\]](#). Here, we consider estimation of Σ_0 (or equivalently Σ) based only on one-dimensional random projections of $X^{(i)}$. More specifically, suppose that the random vectors $X^{(1)}, \dots, X^{(n)}$ are not directly observable and instead we observe

$$(4.1) \quad \xi_i = \langle \beta^{(i)}, X^{(i)} \rangle = \sum_{j=1}^p \beta_j^{(i)} X_j^{(i)}, \quad i = 1, \dots, n,$$

where $\beta^{(i)} \stackrel{\text{i.i.d.}}{\sim} N(0, I_p)$. The goal is to recover Σ_0 from the projections $\{\xi_i, i = 1, \dots, n\}$.

Let $y = (y_1, \dots, y_n)^\top$ with $y_i = \xi_i^2 - \beta^{(i)\top} \beta^{(i)}$. Note that

$$E(\xi^2 | \beta) = E\left(\sum_{i,j} \beta_i \beta_j X_i X_j \mid \beta\right) = \sum_{i,j} \beta_i \beta_j \sigma_{i,j} = \beta^\top \Sigma \beta$$

and so $E(\xi^2 - \beta^\top \Sigma \beta) = \beta^\top \Sigma_0 \beta$. Define a linear map $\mathcal{X} : \mathbb{S}^p \rightarrow \mathbb{R}^n$ by

$$(4.2) \quad [\mathcal{X}(A)]_i = \beta^{(i)\top} A \beta^{(i)}.$$

Then y can be formally written as

$$(4.3) \quad y = \mathcal{X}(\Sigma_0) + z,$$

where $z = y - \mathcal{X}(\Sigma_0)$. We define the corresponding $\tilde{\mathcal{X}}$ and \tilde{y} as in [\(2.14\)](#) and [\(2.15\)](#), respectively, and apply the constraint nuclear norm minimization to recover the low-rank matrix Σ_0 by

$$(4.4) \quad \hat{\Sigma}_0 = \arg \min_M \{\|M\|_* : \|y - \mathcal{X}(M)\| \leq \eta_1, \|\tilde{\mathcal{X}}^*(\tilde{y} - \tilde{\mathcal{X}}(M))\| \leq \eta_2\}.$$

The tuning parameters η_1 and η_2 are chosen as

$$(4.5) \quad \eta_1 = c_1 \sum_{i=1}^n \xi_i^2 \quad \text{and} \quad \eta_2 = 24c_2 \sqrt{p \sum_{i=1}^n \xi_i^4 + 48c_3 p \log n \max_{1 \leq i \leq n} \xi_i^2},$$

where $c_1 > \sqrt{2}$, $c_2, c_3 > 1$ are constants.

We have the following result on the estimator (4.4) for spiked covariance matrix estimation.

THEOREM 4.1. *Suppose $n \geq 3$, we observe $\xi_i, i = 1, \dots, n$, as in (4.1), where $\beta^{(i)} \stackrel{i.i.d.}{\sim} N(0, I_p)$ and $X^{(1)}, \dots, X^{(n)} \stackrel{i.i.d.}{\sim} N(0, \Sigma)$ with $\Sigma = I_p + \Sigma_0$ and Σ_0 positive semidefinite and $\text{rank}(\Sigma_0) \leq r$. Let $\hat{\Sigma}_0$ be given by (4.4). Then there exist uniform constants C, D, δ such that when $n \geq Drp$,*

$$(4.6) \quad \begin{aligned} & \|\hat{\Sigma}_0 - \Sigma_0\|_F^2 \\ & \leq C \min \left(\frac{rp}{n} \|\Sigma\|_*^2 + \frac{rp^2 \log^4 n}{n^2} (\|\Sigma\|_*^2 + \log^2 n \|\Sigma\|^2), \|\Sigma\|_*^2 \right) \end{aligned}$$

with probability at least $1 - O(1/n) - 4 \exp(-p\delta) - \frac{2}{\sqrt{2\pi \log n}}$.

REMARK 4.1. We have focused estimation of spiked covariance matrices on the setting where the random vectors $X^{(i)}$ are Gaussian. Similar to the discussion in Section 3, the results given here can be extended to more general distributions under certain moment conditions.

REMARK 4.2. The problem considered in this section is related to the so-called covariance sketching problem considered in Dasarthy et al. [18]. In covariance sketching, the goal is to estimate the covariance matrix of high-dimensional random vectors $X^{(1)}, \dots, X^{(n)}$ based on the low-dimensional projections

$$y^{(i)} = QX^{(i)}, \quad i = 1, \dots, n,$$

where Q is a fixed $m \times p$ projection matrix with $m < p$. The main differences between the two settings are that the projection matrix in covariance sketch is the same for all $X^{(i)}$ and the dimension m is still relatively large with $m \geq C\sqrt{p} \log^3 p$ for some $C > 0$. In our setting, $m = 1$ and Q is random and varies with i . The techniques for solving the two problems are very different. Comparing to [18], the results in this section indicate that there is a significant advantage to have different random projections for different random vectors $X^{(i)}$ as opposed to having the same projection for all $X^{(i)}$.

5. Simulation results. The constrained nuclear norm minimization methods can be efficiently implemented. The estimator \hat{A} proposed in Section 2.3 can be implemented by the following convex programming:

$$(5.1) \quad \begin{aligned} & \text{minimize} && \text{Tr}(B_1) + \text{Tr}(B_2) \\ & \text{subject to} && \begin{bmatrix} B_1 & A \\ A^T & B_2 \end{bmatrix} \succeq 0, \quad \|y - \mathcal{X}(A)\|_1 \leq \lambda_1, \\ & && \|\mathcal{X}^*(y - \mathcal{X}(A))\| \leq \lambda_2, \end{aligned}$$

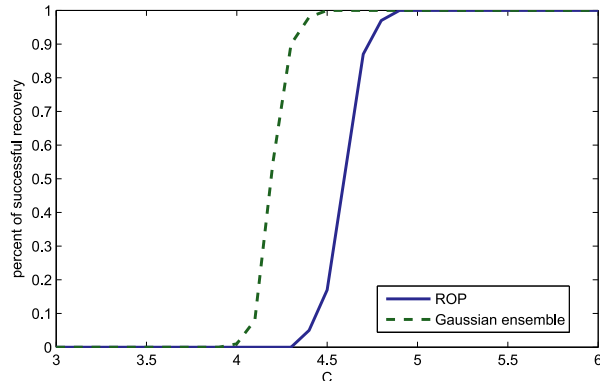


FIG. 1. Rates of successful recovery for the ROP and Gaussian ensemble with $p_1 = p_2 = 100$, $r = 5$, and $n = Cr \max(p_1, p_2)$ for C ranging from 3 to 6.

with optimization variables $B_1 \in \mathbb{S}^{p_1}$, $B_2 \in \mathbb{S}^{p_2}$, $A \in \mathbb{R}^{p_1 \times p_2}$. We use the CVX package [22, 23] to implement the proposed procedures. In this section, a simulation study is carried out to investigate the numerical performance of the proposed procedures for low-rank matrix recovery in various settings.

We begin with the noiseless case. In this setting, Theorem 2.2 and Corollary 2.1 show that the nuclear norm minimization recovers a rank r matrix exactly whenever

$$(5.2) \quad n \geq Cr \max(p_1, p_2).$$

A similar result holds for the Gaussian ensemble [13]. However, the minimum constant C that guarantees the exact recovery with high probability is not specified in either case. It is of practical interest to find the minimum constant C . For this purpose, we randomly generate $p_1 \times p_2$ rank- r matrices A as $A = X^\top Y$, where $X \in \mathbb{R}^{r \times p_1}$, $Y \in \mathbb{R}^{r \times p_2}$ are i.i.d. Gaussian matrices. We compare ROP from the standard Gaussian distribution and the Gaussian ensemble, with the number of measurements $n = Cr \max(p_1, p_2)$ from a range of values of C using the constrained nuclear norm minimization (2.1). A recovery is considered successful if $\|\hat{A} - A\|_F / \|A\|_F \leq 10^{-4}$. Figure 1 shows the rate of successful recovery when $p_1 = p_2 = 100$ and $r = 5$.

The numerical results show that for ROP from the Gaussian distribution, the minimum constant C to ensure exact recovery with high probability is slightly less than 5 in the small scale problems ($p_1, p_2 \leq 100$) we tested. The corresponding minimum constant C for the Gaussian ensemble is about 4.5. Matrix completion requires much larger number of measurements. Based on the theoretical analyses given in [14, 35], the required number of measurements for matrix completion is $O(\mu r(p_1 + p_2) \log^2(p_1 + p_2))$, where $\mu \geq 1$ is some coherence constant describing the “spikedness” of the matrix A .

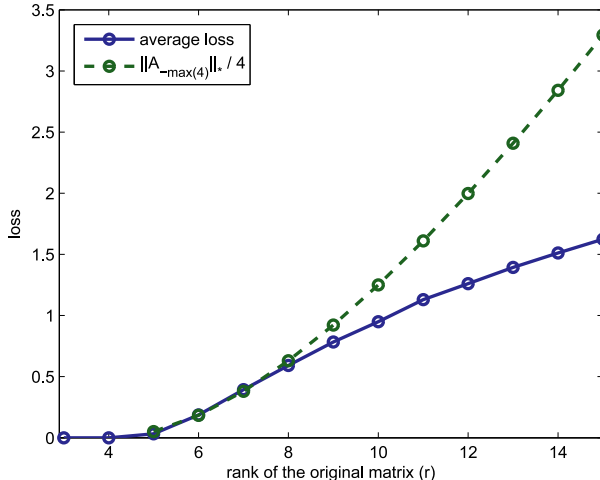


FIG. 2. Recovery accuracy (solid line) for approximately low-rank matrices with different values of r , where $p_1 = p_2 = 100$, $n = 2000$, $\sigma(A) = (1, 1/\sqrt{2}, \dots, 1/\sqrt{r})$. The dashed line is the theoretical upper bound.

Hence, for matrix completion, the factor C in (5.2) needs to grow with the dimensions p_1 and p_2 and it requires $C \gtrsim \mu \log^2(p_1 + p_2)$, which is much larger than what is needed for the ROP or Gaussian ensemble. The required storage space for the Gaussian ensemble is much greater than that for the ROP. In order to ensure accurate recovery of $p \times p$ matrices of rank r , one needs at least $4.5rp^3$ bytes of space to store the measurement matrices, which could be prohibitively large for the recovery of high-dimensional matrices. In contrast, the storage space for the projection vectors in ROP is only $10rp^2$ bytes, which is far smaller than what is required by the Gaussian ensemble in the high-dimensional case.

We then consider the recovery of approximately low-rank matrices to investigate the robustness of the method against small perturbations. To this end, we randomly draw 100×100 matrix A as $A = U \cdot \text{diag}(1, 2^{-1/2}, \dots, r^{-1/2}) \cdot V^\top$, where $U \in \mathbb{R}^{100 \times r}$ and $V \in \mathbb{R}^{100 \times r}$ are random matrices with orthonormal columns. We then observe $n = 2000$ random rank-one projections with the measurement vectors being i.i.d. Gaussian. Based on the observations, the nuclear minimization procedure (2.1) is applied to estimate A . The results for different values of r are shown in Figure 2. It can be seen from the plot that in this setting one can exactly recover a matrix of rank at most 4 with 2000 measurements. However, when the rank r of the true matrix A exceeds 4, the estimate is still stable. The theoretical result in Proposition 2.1 bounds the loss (solid line) at $O(\|A_{-\max(4)}\|_*^2/4)$ (shown in the dashed line) with high probability, which corresponds to Figure 2.

We now turn to the noisy case. The low-rank matrices A are generated by $A = X^T Y$, where $X \in \mathbb{R}^{r \times p_1}$ and $Y \in \mathbb{R}^{r \times p_2}$ are i.i.d. Gaussian matrices. The ROP \mathcal{X} is from the standard Gaussian distribution and the noise vector $z \sim N_n(0, \sigma^2)$. Based on (\mathcal{X}, y) with $y = \mathcal{X}(A) + z$, we compare our proposed estimator \hat{A} with the ℓ_1 constraint minimization estimator \hat{A}^{ℓ_1} [17] and the matrix Dantzig Selector \hat{A}^{DS} [13], where

$$\hat{A} = \arg \min_M \{\|M\|_* : y - \mathcal{X}(M) \in \mathcal{Z}_1 \cap \mathcal{Z}_2\},$$

$$\hat{A}^{\ell_1} = \arg \min_M \{\|M\|_* : y - \mathcal{X}(M) \in \mathcal{Z}_1\},$$

$$\hat{A}^{\text{DS}} = \arg \min_M \{\|M\|_* : y - \mathcal{X}(M) \in \mathcal{Z}_2\},$$

with $\mathcal{Z}_1 = \{z : \|z\|_1/n \leq \sigma\}$ and $\mathcal{Z}_2 = \{z : \|\mathcal{X}(z)\| \leq \sigma(\sqrt{\log n}(p_1 + p_2) + \sqrt{n(p_1 + p_2)})\}$. Note that \hat{A}^{ℓ_1} is similar to the estimator proposed in Chen et al. [17], except their estimator is for symmetric matrices under the SROP but ours is for general low-rank matrices under the ROP. Figure 3 compares the performance of the three estimators. It can be seen from the left panel that for small n , ℓ_1 constrained minimization outperforms the matrix Dantzig Selector, while our estimator outperforms both \hat{A}^{ℓ_1} and \hat{A}^{DS} . When n is large, our estimator and \hat{A}^{DS} are essentially the same and both outperforms \hat{A}^{ℓ_1} . The right panel of Figure 3 plots the ratio of the squared Frobenius norm loss of \hat{A}^{ℓ_1} to that of our estimator. The ratio increases with n . These numerical results are consistent with the observations made in Remark 2.3.

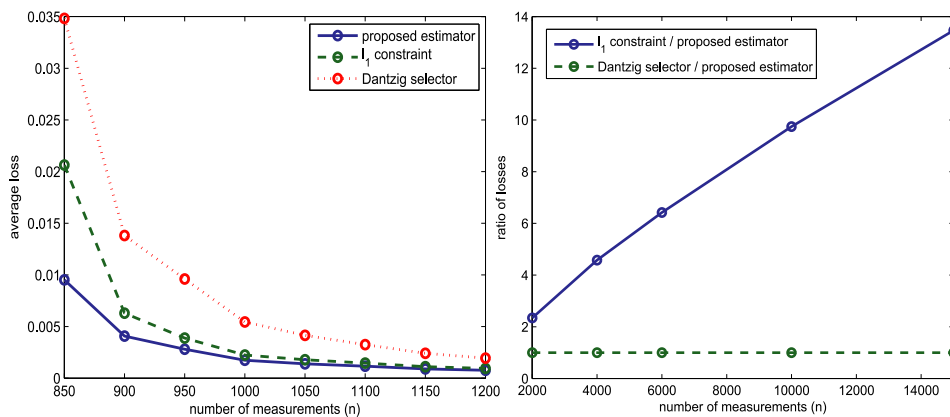


FIG. 3. *Left panel: Comparison of the proposed estimator with \hat{A}^{ℓ_1} and \hat{A}^{DS} for $p_1 = p_2 = 50$, $r = 4$, $\sigma = 0.01$, and n ranging from 850 to 1200. Right panel: Ratio of the squared Frobenius norm loss of \hat{A}^{ℓ_1} to that of the proposed estimator for $p_1 = p_2 = 50$, $r = 4$, and n varying from 2000 to 15,000.*

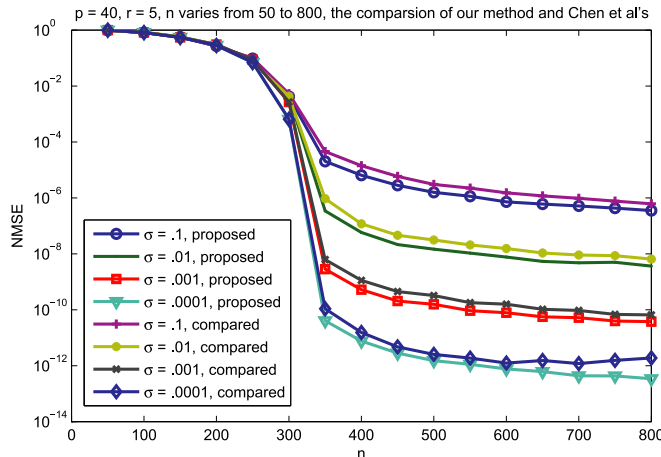


FIG. 4. Comparison of the proposed estimator \hat{A} with the \hat{A}^{ℓ_1} . Here $p = 40$, $r = 5$, $\sigma = 0.1, 0.01, 0.001, 0.0001$ and n ranges from 50 to 800.

We now turn to the recovery of symmetric low-rank matrices under the SROP model (2.12). Let \mathcal{X} be SROP from the standard normal distribution. We consider the setting where $p = 40$, n varies from 50 to 600, $z_i \sim \sigma \cdot \mathcal{U}[-1, 1]$ with $\sigma = 0.1, 0.01, 0.001$ or 0.0001 , and A is randomly generated as rank-5 matrix by the same procedure discussed above. The setting is identical to the one considered in Section 5.1 of [17]. Although we cannot exactly repeat the simulation study in [17] as they did not specify the choice of the tuning parameter, we can implement both our procedure

$$\hat{A} = \arg \min_M \left\{ \|M\|_* : \|y - \mathcal{X}(M)\|_1 \leq \frac{n\sigma}{2}, \right. \\ \left. \|\tilde{\mathcal{X}}^*(\tilde{y} - \tilde{\mathcal{X}}(M))\| \leq \frac{\sigma(\sqrt{\log np} + \sqrt{np})}{3} \right\}$$

and the estimator \hat{A}^{ℓ_1} with only the ℓ_1 constraint which was proposed by Chen et al. [17]

$$\hat{A}^{\ell_1} = \arg \min_M \left\{ \|M\|_* : \|y - \mathcal{X}(M)\|_1 \leq \frac{n\sigma}{2} \right\}.$$

The results are given in Figure 4. It can be seen that our estimator \hat{A} outperforms the estimator \hat{A}^{ℓ_1} .

5.1. *Data driven selection of tuning parameters.* We have so far considered the estimators

$$(5.3) \quad \hat{A} = \arg \min_B \{ \|B\|_* : \|y - \mathcal{X}(B)\|_1/n \leq \lambda, \|\mathcal{X}^*(y - \mathcal{X}(B))\| \leq \eta \},$$

$$(5.4) \quad \hat{A} = \arg \min_M \{ \|M\|_* : \|y - \mathcal{X}(M)\|_1/n \leq \lambda, \|\tilde{\mathcal{X}}^*(\tilde{y} - \tilde{\mathcal{X}}(M))\| \leq \eta \}$$

for the ROP and SROP, respectively. The theoretical choice of the tuning parameters λ and η depends on the knowledge of the error distribution such as the variance. In real applications, such information may not be available and/or the theoretical choice may not be the best. It is thus desirable to have a data driven choice of the tuning parameters. We now introduce a practical method for selecting the tuning parameters using K -fold cross-validation.

Let $(\mathcal{X}, y) = \{(X_i, y_i), i = 1, \dots, n\}$ be the observed sample and let T be a grid of positive real values. For each $t \in T$, set

$$(5.5) \quad (\lambda, \eta) = (\lambda(t), \eta(t)) = \begin{cases} (t, t(\sqrt{\log n}(p_1 + p_2) + \sqrt{n(p_1 + p_2)})), & \text{for ROP;} \\ (t, t(\sqrt{\log np} + \sqrt{np})), & \text{for SROP.} \end{cases}$$

Randomly split the n samples $(X_i, y_i), i = 1, \dots, n$ into two groups of sizes $n_1 \sim \frac{(K-1)n}{K}$ and $n_2 \sim \frac{n}{K}$ for I times. Denote by $J_1^i, J_2^i \subseteq \{1, \dots, n\}$ the index sets for Groups 1 and 2, respectively, for the i th split. Apply our procedure [(5.3) for ROP and (5.4) for SROP, resp.] to the sub-samples in Group 1 with the tuning parameters $(\lambda(t), \eta(t))$ and denote the estimators by $\hat{A}^i(t)$, $i = 1, \dots, I$. Evaluate the prediction error of $\hat{A}^i(t)$ over the subsample in Group 2 and set

$$\hat{R}(t) = \sum_{i=1}^I \sum_{j \in J_2^i} |y_j - \langle \hat{A}^i(t), X_j \rangle|^2, \quad t \in T.$$

We select

$$t_* = \arg \min_T \hat{R}(t)$$

and choose the tuning parameters $(\lambda(t_*), \eta(t_*))$ as in (5.5) with $t = t_*$ and the final estimator \hat{A} based on (5.3) or (5.4) with the chosen tuning parameters.

We compare the numerical result by 5-fold cross-validation with the result based on the known σ by simulation in Figure 5. Both the ROP and SROP are considered. It can be seen that the estimator with the tuning parameters chosen through 5-fold cross-validation has the same performance as or outperforms the one with the theoretical choice of the tuning parameters.

5.2. Image compression. Since a two-dimensional image can be considered as a matrix, one approach to image compression is by using low-rank matrix approximation via the singular value decomposition. See, for example, [2, 36, 40]. Here, we use an image recovery example to further illustrate the nuclear norm minimization method under the ROP model.

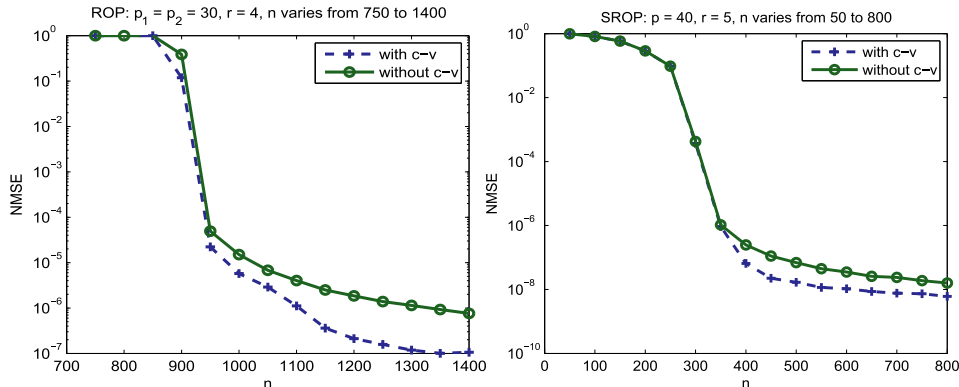


FIG. 5. Comparison of the performance with cross validation and without cross-validation in both ROP and SROP. Left panel: ROP, $p_1 = p_2 = 30$, $r = 4$, n varies from 750 to 1400. Right panel: SROP, $p = 40$, $r = 5$, n varies from 50 to 800.

For a grayscale image, let $A = (a_{i,j}) \in \mathbb{R}^{m \times n}$ be the intensity matrix associated with the image, where a_{ij} is the grayscale intensity of the (i, j) pixel. When the matrix A is approximately low-rank, the ROP model and nuclear norm minimization method can be used for image compression and recovery. To illustrate this point, let us consider the following grayscale MIT Logo image (Figure 6).

The matrix associated with MIT logo is of the size 50×80 and of rank 6. We take rank-one random projections $\mathcal{X}(A)$ as the observed sample, with various sample sizes. Then the constrained nuclear norm minimization method is applied to reconstruct the original low-rank matrix. The recovery results are shown in Figure 7. The results show that the original image can be compressed and recovered well via the ROP model and the nuclear norm minimization.

6. Discussions. This paper introduces the ROP model for the recovery of general low-rank matrices. A constrained nuclear norm minimization method is proposed and its theoretical and numerical properties are studied. The proposed estimator is shown to be rate-optimal when the number of rank-one projections $n \gtrsim \log n(p_1 + p_2)$ or $n \sim r(p_1 + p_2)$. It is also shown that the



FIG. 6. Original grayscale MIT logo.



FIG. 7. Recovery of MIT logo based on different number of measurements. Left: 900; Middle: 1000; Right: 1080.

procedure is adaptive to the rank and robust against small perturbations. The method and results are applied to estimation of a spiked covariance matrix. It is somewhat unexpected that it is possible to accurately recover a spiked covariance matrix from only one-dimensional projections. An interesting open problem is to estimate the principal components/subspace based on the one-dimensional random projections. We leave this as future work.

In a recent paper, Chen et al. [17] considered quadratic measurements for the recovery of symmetric positive definite matrices, which is similar to the special case of SROP that we studied here. The paper was posted on arXiv as we finish writing the present paper. They considered the noiseless and ℓ_1 bounded noise cases and introduced the so-called “RIP- ℓ_2/ℓ_1 ” condition. The “RIP- ℓ_2/ℓ_1 ” condition is similar to RUB in our work. But these two conditions are not identical as the RIP- ℓ_2/ℓ_1 condition can only be applied to symmetric low-rank matrices as only symmetric operators are considered in the paper. In contrast, RUB applies to all low-rank matrices.

Chen et al. ([17] version 4) considered ℓ_1 -bounded noise case under the SROP model and gave an upper bound in their Theorem 3 (after a slight change of notation)

$$(6.1) \quad \|\hat{\Sigma} - \Sigma\|_F \leq C_1 \frac{\|\Sigma - \Sigma_\Omega\|_*}{\sqrt{r}} + C_2 \frac{\varepsilon}{n}.$$

This result for ℓ_1 bounded noise case is not applicable to the i.i.d. random noise setting. When the entries of the noise term $\eta \in \mathbb{R}^n$ are of constant order, which is the typical case for i.i.d. noise with constant variance, one has $\|\eta\|_1 \sim Cn$ with high probability. In such a case, the term $C_2 \frac{\varepsilon}{n}$ on the right-hand side of (6.1) does not even converge to 0 as the sample size $n \rightarrow \infty$.

In comparison, the bound (3.6) in Proposition 3.1 can be equivalently rewritten as

$$(6.2) \quad \|\hat{A} - A\|_F \leq W_2 \frac{\|A_{-\max(r)}\|_*}{\sqrt{r}} + W_1 \tau \min\left(\frac{\sqrt{r \log np}}{n} + \sqrt{\frac{rp}{n}}, 1\right),$$

where the first term $W_2 \frac{\|A_{-\max(r)}\|_*}{\sqrt{r}}$ is of the same order as $C_1 \frac{\|\Sigma - \Sigma_\Omega\|_*}{\sqrt{r}}$ in (6.1) while the second term decays to 0 as $n \rightarrow \infty$. Hence, for the recovery

of rank- r matrices, as the sample size n increases our bound decays to 0 but the bound (6.1) given in Chen et al. [17] does not. The main reason of this phenomenon lies in the difference in the two methods: we use nuclear norm minimization under two convex constraints (see Remark 2.3), but Chen et al. [17] used only the ℓ_1 constraint. Both theoretical results (see Remark 2.3) and numerical results (Figure 3 in Section 5) show that the additional constraint \mathcal{Z}_2 improves the performance of the estimator.

Moreover, the results and techniques in [17] for symmetric positive definite matrices are not applicable to the recovery of general nonsymmetric matrices. This is due to the fact that for a nonsymmetric square matrix $A = (a_{ij})$, the quadratic measurements $(\beta^{(i)})^\top A \beta^{(i)}$ satisfy

$$(\beta^{(i)})^\top A \beta^{(i)} = (\beta^{(i)})^\top A^s \beta^{(i)},$$

where $A^s = \frac{1}{2}(A + A^\top)$. Hence, for a nonsymmetric matrix A , only its symmetrized version A^s can be possibly identified and estimated based on the quadratic measurements, the matrix A itself is neither identifiable nor estimable.

7. Proofs. We prove the main results in this section. We begin by collecting a few important technical lemmas that will be used in the proofs of the main results. The proofs of some of these technical lemmas are involved and are postponed to the supplementary material [11].

7.1. *Technical tools.* Lemmas 7.1 and 7.2 below are used for deriving the RUB condition (see Definition 2.1) from the ROP design.

LEMMA 7.1. *Suppose $A \in \mathbb{R}^{p_1 \times p_2}$ is a fixed matrix and \mathcal{X} is ROP from a symmetric sub-Gaussian distribution \mathcal{P} , that is,*

$$[\mathcal{X}(A)]_j = \beta^{(j)T} A \gamma^{(j)}, \quad j = 1, \dots, n,$$

where $\beta^{(j)} = (\beta_1^{(j)}, \dots, \beta_{p_1}^{(j)})^T$, $\gamma^{(j)} = (\gamma_1^{(j)}, \dots, \gamma_{p_2}^{(j)})^T$ are random vectors with entries i.i.d. generated from \mathcal{P} . Then for $\delta > 0$, we have

$$\left(\frac{1}{3\alpha_{\mathcal{P}}^4} - 2\alpha_{\mathcal{P}}^2\delta - \alpha_{\mathcal{P}}^2\delta^2 \right) \|A\|_F \leq \|\mathcal{X}(A)\|_1/n \leq (1 + 2\alpha_{\mathcal{P}}^2\delta + \alpha_{\mathcal{P}}^2\delta^2) \|A\|_F$$

with probability at least $1 - 2\exp(-\delta^2 n)$. Here, $\alpha_{\mathcal{P}}$ is defined by (3.2).

LEMMA 7.2. *Suppose $A \in \mathbb{R}^{p_1 \times p_2}$ is a fixed matrix. $\beta = (\beta_1, \dots, \beta_{p_1})^T$, $\gamma = (\gamma_1, \dots, \gamma_{p_2})^T$ are random vectors such that $\beta_1, \dots, \beta_{p_1}, \gamma_1, \dots, \gamma_{p_2} \stackrel{i.i.d.}{\sim} \mathcal{P}$, where \mathcal{P} is some symmetric variance 1 sub-Gaussian distribution, then we have*

$$\frac{\|A\|_F}{3\alpha_{\mathcal{P}}^4} \leq E|\beta^T A \gamma| \leq \|A\|_F,$$

where $\alpha_{\mathcal{P}}$ is given by (3.2).

Let $z \in \mathbb{R}^n$ be i.i.d. sub-Gaussian distributed. By measure concentration theory, $\|z\|_p^p/n$, $1 \leq p \leq \infty$, are essentially bounded; specifically, we have the following lemma.

LEMMA 7.3. *Suppose $z \in \mathbb{R}^n$ and $z_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$, we have*

$$\begin{aligned} P(\|z\|_1 \geq \sigma n) &\leq \frac{9}{n}, \\ P(\|z\|_2 \geq \sigma \sqrt{n + 2\sqrt{n \log n}}) &\leq \frac{1}{n}, \\ P(\|z\|_\infty \geq 2\sigma \sqrt{\log n}) &\leq \frac{1}{n\sqrt{2\pi \log n}}. \end{aligned}$$

More general, when z_i are i.i.d. sub-Gaussian distributed such that (3.1) holds, then

$$\begin{aligned} P(\|z\|_1 \geq Cn) &\leq \exp\left(-\frac{n(C - 2\sqrt{2\pi}\gamma)^2}{2\gamma^2}\right) \quad \forall C > 2\sqrt{2\pi}\gamma, \\ P(\|z\|_2 \geq \sqrt{Cn}) &\leq \exp\left(-\frac{n(C - 4\gamma^2)^2}{8\gamma^2 C}\right) \quad \forall C > 4\gamma^2, \\ P(\|z\|_\infty \geq C\gamma\sqrt{\log n}) &\leq 2n^{-C^2/2-1} \quad \forall C > 0. \end{aligned}$$

Lemma 7.4 below presents an upper bound for the spectral norm of $\mathcal{X}(z)$ for a fixed vector z .

LEMMA 7.4. *Suppose \mathcal{X} is ROP from some symmetric sub-Gaussian distribution \mathcal{P} and $z \in \mathbb{R}^n$ is some fixed vector, then for $C > \log 7$, we have*

$$\|\mathcal{X}^*(z)\| \leq 3\alpha_{\mathcal{P}}^2(C(p_1 + p_2)\|z\|_\infty + \sqrt{2C(p_1 + p_2)}\|z\|_2)$$

with probability at least $1 - 2\exp(-(C - \log 7)(p_1 + p_2))$. Here, $\alpha_{\mathcal{P}}$ is defined by (3.2).

We are now ready to prove the main results of the paper.

7.2. *Proof of Theorem 2.1.* We introduce the following two technical lemmas that will be used in the proof of theorem.

The *null space property* below is a well-known result in affine rank minimization problem (see [32]). It provides a necessary, sufficient and easier-to-check condition for exact recovery in the noiseless setting.

LEMMA 7.5 (Null space property). *Using (2.1), one can recover all matrices A of rank at most r if and only if for all $R \in \mathcal{N}(\mathcal{X}) \setminus \{0\}$,*

$$\|R_{\max(r)}\|_* < \|R_{-\max(r)}\|_*.$$

The following lemma is given in [10], which provides a way to decompose the general vectors to sparse ones.

LEMMA 7.6 (Sparse representation of a polytope). *Suppose s is a non-negative integer, $v \in \mathbb{R}^p$ and $\theta \geq 0$. Then $\|v\|_\infty \leq \theta$, $\|v\|_1 \leq s\theta$, if and only if v can be expressed as a weighted mean,*

$$v = \sum_{i=1}^N \lambda_i u_i, \quad 0 \leq \lambda_i \leq 1, \quad \sum_{i=1}^N \lambda_i = 1,$$

where u_i satisfies

$$(7.1) \quad \begin{aligned} u_i \text{ is } s\text{-sparse}, \quad \text{supp}(u_i) \subseteq \text{supp}(v), \\ \|u_i\|_1 = \|v\|_1, \quad \|u_i\|_\infty \leq \theta. \end{aligned}$$

For the proof of Theorem 2.1, by null space property (Lemma 7.5), we only need to show for all nonzero R with $\mathcal{X}(R) = 0$, we must have $\|R_{\max(r)}\|_* < \|R_{-\max(r)}\|_*$.

If this does not hold, suppose there exists nonzero R with $\mathcal{X}(R) = 0$ and $\|R_{\max(r)}\|_* \geq \|R_{-\max(r)}\|_*$. We denote $p = \min(p_1, p_2)$ and assume the singular value decomposition of R is

$$R = \sum_{i=1}^p \sigma_i u_i v_i^\top = U \text{diag}(\vec{\sigma}) V^\top,$$

where u_i, v_i are orthogonal basis in $\mathbb{R}^{p_1}, \mathbb{R}^{p_2}$, respectively, and $\vec{\sigma}$ is the singular value vector such that $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$. Without loss of generality, we can assume $p \geq kr$, otherwise we can set the undefined entries of σ as 0.

Consider the singular value vector $\vec{\sigma} = (\sigma_1, \sigma_2, \dots, \sigma_p)$, we note that $\vec{\sigma}_{-\max(kr)}$ satisfies

$$\begin{aligned} \|\vec{\sigma}_{-\max(kr)}\|_\infty &\leq \sigma_{kr}, \\ \|\vec{\sigma}_{-\max(kr)}\|_1 &= \|\vec{\sigma}_{-\max(r)}\|_1 - (\sigma_{r+1} + \dots + \sigma_{kr}) \\ &\leq \|\vec{\sigma}_{-\max(r)}\|_1 - (k-1)r\sigma_{kr} \\ &\leq \|\vec{\sigma}_{\max(r)}\|_1 - (k-1)r\sigma_{kr}. \end{aligned}$$

Denote $\theta = \max\{\sigma_{kr}, (\|\vec{\sigma}_{\max(r)}\|_1 - r(k-1)\sigma_{kr})/(kr)\}$, by the two inequalities above we have $\|\vec{\sigma}_{-\max(kr)}\|_\infty \leq \theta$ and $\|\vec{\sigma}_{-\max(kr)}\|_1 \leq kr\theta$. Now apply Lemma 7.6, we can get $b^{(i)} \in \mathbb{R}^p$, $\lambda_i \geq 0, i = 1, \dots, N$ such that $\sum_{i=1}^N \lambda_i = 1$, $\vec{\sigma}_{-\max(kr)} = \sum_{i=1}^N \lambda_i b^{(i)}$ and

$$(7.2) \quad \begin{aligned} \text{supp}(b^{(i)}) &\subseteq \text{supp}(\vec{\sigma}_{-\max(kr)}), \quad \|b^{(i)}\|_0 \leq kr, \\ \|b^{(i)}\|_1 &= \|\vec{\sigma}_{-\max(kr)}\|_1, \quad \|b^{(i)}\|_\infty \leq \theta, \end{aligned}$$

which leads to

$$\|b^{(i)}\|_2 \leq \sqrt{\|b^{(i)}\|_1 \cdot \|b^{(i)}\|_\infty} \leq \sqrt{(\|\vec{\sigma}_{\max(r)}\|_1 - r(k-1)\sigma_{kr}) \cdot \theta}.$$

If $\theta = \sigma_{kr}$, we have

$$\begin{aligned} \|b^{(i)}\|_2 &\leq \sqrt{(\|\vec{\sigma}_{\max(r)}\|_1 - r(k-1)\sigma_{kr})\sigma_{kr}} \\ &\leq \sqrt{\left(\|\vec{\sigma}_{\max(r)}\|_1 - r(k-1)\frac{\|\vec{\sigma}_{\max(r)}\|_1}{2r(k-1)}\right)\frac{\|\vec{\sigma}_{\max(r)}\|_1}{2r(k-1)}} \\ &\leq \frac{\|\vec{\sigma}_{\max(r)}\|_1}{\sqrt{4r(k-1)}} \leq \frac{\|\vec{\sigma}_{\max(r)}\|_2}{\sqrt{4(k-1)}}. \end{aligned}$$

If $\theta = (\|\vec{\sigma}_{\max(r)}\|_1 - r(k-1)\sigma_{kr})/(kr)$, we have

$$\|b^{(i)}\|_2 \leq \sqrt{\frac{1}{kr}(\|\vec{\sigma}_{\max(r)}\|_1 - r(k-1)\sigma_{kr})} \leq \sqrt{\frac{1}{kr}}\|\vec{\sigma}_{\max(r)}\|_1 \leq \frac{\|\vec{\sigma}_{\max(r)}\|_2}{\sqrt{k}}.$$

Since $k \geq 2$, we always have $\|b^{(i)}\|_2 \leq \|\vec{\sigma}_{\max(r)}\|_2/\sqrt{k}$. Finally, we define $B_i = U \text{diag}(b^{(i)})V^\top$, then the rank of B_i are all at most kr and $\sum_{i=1}^N \lambda_i B_i = R_{-\max(kr)}$ and

$$\|B_i\|_F = \|b^{(i)}\|_2 \leq \|\vec{\sigma}_{\max(r)}\|_2/\sqrt{k} = \|R_{\max(r)}\|_F/\sqrt{k}.$$

Hence,

$$\begin{aligned} 0 &= \|\mathcal{X}(R)\|_1 \geq \|\mathcal{X}(R_{\max(kr)})\|_1 - \|\mathcal{X}(R_{-\max(kr)})\|_1 \\ &\geq C_1 \|R_{\max(kr)}\|_F - \sum_{i=1}^N \|\mathcal{X}(\lambda_i B_i)\|_1 \\ &\geq C_1 \|R_{\max(r)}\|_F - \sum_{i=1}^N \lambda_i C_2 \|B_i\|_F \\ &\geq C_1 \|R_{\max(r)}\|_F - C_2 \|R_{\max(r)}\|_F/\sqrt{k} > 0. \end{aligned}$$

Here, we used the RUB condition. The last inequality is due to $C_2/C_1 < \sqrt{k}$ and $R \neq 0$ (so $R_{\max(r)} \neq 0$). This is a contradiction, which completes the proof of the theorem.

7.3. Proof of Theorem 2.2. Notice that for \mathcal{P} as standard Gaussian distribution, the constant $\alpha_{\mathcal{P}}$ [defined as (3.2)] equals 1. We will prove the following more general result than Theorem 2.2 instead. The proof is provided in the supplementary material [11].

PROPOSITION 7.1. *Suppose $\mathcal{X} : \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$ is ROP from some variance 1 symmetric sub-Gaussian distribution \mathcal{P} . For integer $k \geq 2$, positive $C_1 < \frac{1}{3\alpha_{\mathcal{P}}}$ [$\alpha_{\mathcal{P}}$ is defined as (3.2)] and $C_2 > 1$, there exists constants C and δ , only depending on \mathcal{P}, C_1, C_2 but not on p_1, p_2, r , such that if $n \geq Cr(p_1 + p_2)$, then with probability at least $1 - e^{-n\delta}$, \mathcal{X} satisfies RUB of order kr and constants C_1 and C_2 .*

7.4. *Proof of Theorems 2.3 and 3.1, Proposition 2.1.* In order to prove the result, we introduce the following technical lemma as an extension of null space property (Lemma 7.5) from exact low-rank into the approximate low-rank setting.

LEMMA 7.7. *Suppose $A_*, A \in \mathbb{R}^{p_1 \times p_2}$, $R = A_* - A$. If $\|A_*\|_* \leq \|A\|_*$, we have*

$$(7.3) \quad \|R_{-\max(r)}\|_* \leq \|R_{\max(r)}\|_* + 2\|A_{-\max(r)}\|_*.$$

The following two lemmas described the separate effect of constraint $\mathcal{Z}_1 = \{z : \|z\|_1/n \leq \lambda_1\}$ and $\mathcal{Z}_2 = \{z : \|\mathcal{X}^*(z)\| \leq \lambda_2\}$ on the estimator.

LEMMA 7.8. *Suppose \mathcal{X} satisfies RUB condition of order kr with constants C_1, C_2 such that $C_1 > C_2/\sqrt{k}$. Assume that $A_*, A \in \mathbb{R}^{p_1 \times p_2}$ satisfy $\|A_*\|_* \leq \|A\|_*$, $\|\mathcal{X}(A_* - A)\|_1/n \leq \lambda_1$. Then we have*

$$\|A_* - A\|_F \leq \frac{2}{C_1 - C_2/\sqrt{k}} \lambda_1 + \left(\frac{3}{\sqrt{k}C_1/C_2 - 1} + \frac{1}{\sqrt{k} - 1} \right) \frac{\|A_{-\max(r)}\|_*}{\sqrt{r}}.$$

LEMMA 7.9. *Suppose \mathcal{X} satisfies RUB condition of order kr with constants C_1, C_2 such that $C_1 > C_2/\sqrt{k}$. Assume that \hat{A}^{DS} satisfies $\|\mathcal{X}^*\mathcal{X}(A_* - A)\| \leq \lambda_2$. Then we have*

$$\begin{aligned} \|A_* - A\|_F &\leq \frac{4}{(C_1 - C_2/\sqrt{k})^2} \cdot \frac{\sqrt{r}\lambda_2}{n} \\ &\quad + \left(\frac{5}{\sqrt{k}C_1/C_2 - 1} + \frac{1}{\sqrt{k} - 1} + 1 \right) \frac{\|A_{-\max(r)}\|_*}{\sqrt{r}}. \end{aligned}$$

The proof of Lemmas 7.7, 7.8 and 7.9 are listed in the supplementary material [11]. Now we prove Theorem 2.3 and Proposition 2.1. We only need to prove Proposition 2.1 since Theorem 2.3 is a special case of Proposition 2.1. By Lemmas 7.3 and 7.4, we have

$$P_z(\|z\|_1 \leq \sigma n) \leq \frac{9}{n},$$

$$\begin{aligned}
P_{\mathcal{X},z}(\|\mathcal{X}^*(z)\| \geq \sigma(12(p_1 + p_2)\sqrt{\log n} + 6\sqrt{2(p_1 + p_2)n})) \\
\leq P_{\mathcal{X}}(\|\mathcal{X}^*(z)\| \geq (6(p_1 + p_2)\|z\|_{\infty} + 6\sqrt{p_1 + p_2}\|z\|_2)) \\
+ P_z(\|z\|_{\infty} \geq 2\sigma\sqrt{\log n}) + P_z(\|z\|_2 \geq \sigma\sqrt{2n}) \\
\leq 2\exp(-(2 - \log 7)(p_1 + p_2)) + \frac{1}{n\sqrt{2\pi\log n}} + \frac{1}{n}.
\end{aligned}$$

Here, $P_{\mathcal{X}}$ (P_z or $P_{\mathcal{X},z}$) means the probability with respect to \mathcal{X} [z or (\mathcal{X}, z)]. Hence, we have

$$P(z \in \mathcal{Z}_1 \cap \mathcal{Z}_2) \geq 1 - 2\exp(-(2 - \log 7)(p_1 + p_2)) - \frac{11}{n}.$$

Under the event that $z \in \mathcal{Z}_1 \cap \mathcal{Z}_2$, A is in the feasible set of the programming (2.6), which implies $\|\hat{A}\|_* \leq \|A\|_*$ by the definition of \hat{A} . Moreover, we have

$$\begin{aligned}
\|\mathcal{X}(\hat{A} - A)\|_1/n &\leq \|y - \mathcal{X}(A)\|_1/n + \|y - \mathcal{X}(\hat{A})\|_1/n \\
&\leq \|z\|_1/n + \|y - \mathcal{X}(\hat{A})\|_1/n \leq 2\sigma, \\
\|\mathcal{X}^*\mathcal{X}(\hat{A} - A)\| &\leq \|\mathcal{X}^*(y - \mathcal{X}(\hat{A}))\| + \|\mathcal{X}^*(y - \mathcal{X}(A))\| \\
&\leq \|\mathcal{X}^*(y - \mathcal{X}(\hat{A}))\| + \|\mathcal{X}^*(z)\| \leq 2\eta.
\end{aligned}$$

On the other hand, suppose $k = 10$, by Theorem 2.2, we can have find a uniform constant C and δ such that if $n \geq Crk(p_1 + p_2)$, \mathcal{X} satisfies RUB of order $10r$ and constants $C_1 = 0.32, C_2 = 1.02$ with probability at least $1 - e^{-n\delta}$. Hence, we have $D(= Ck)$ and δ' such that if $n \geq Dr(p_1 + p_2)$, \mathcal{X} satisfies RUB of order $10r$ and constants C_1, C_2 satisfying $C_2/C_1 < \sqrt{10}$ with probability at least $1 - e^{-n\delta'}$.

Now under the event that:

1. \mathcal{X} satisfies RUB of order $10r$ and constants C_1, C_2 satisfying $C_2/C_1 < \sqrt{10}$,
2. $z \in \mathcal{Z}_1 \cap \mathcal{Z}_2$,

apply Lemmas 7.8 and 7.9 with $A_* = \hat{A}$, we can get (2.11). The probability that these two events both happen is at least $1 - 2\exp(-(2 - \log 7)(p_1 + p_2)) - \frac{11}{n} - \exp(-\delta'n)$. Set $\delta = \min(2 - \log 7, \delta')$, we finished the proof of Proposition 2.1.

For Theorem 3.1, the proof is similar. We apply the latter part of Lemmas 7.3 and 7.4 and get

$$\begin{aligned}
P(z \notin \mathcal{Z}_1 \cap \mathcal{Z}_2) \\
\leq P(\|z\|_1/n > 6\tau) \\
+ P(\|\mathcal{X}(z)\| > \tau\alpha_p^2(6\sqrt{6n(p_1 + p_2)} + 12\sqrt{\log n}(p_1 + p_2)))
\end{aligned}$$

$$\begin{aligned}
&\leq P(\|z\|/n > 6\tau) + P(\|z\|_2 > \sqrt{6n\tau}) + P(\|z\|_\infty > 2\sqrt{\log n\tau}) \\
&\quad + P_{\mathcal{X}}(\|\mathcal{X}(z)\| > \alpha_{\mathcal{P}}^2(6(p_1 + p_2)\|z\|_\infty + 6\sqrt{p_1 + p_2}\|z\|_2)) \\
&\leq \exp(-n(6 - 2\sqrt{2\pi})^2/2) + \exp(-n/12) \\
&\quad + \frac{2}{n} + 2\exp(-(2 - \log 7)(p_1 + p_2)).
\end{aligned}$$

Besides, we choose $k > (3\alpha_{\mathcal{P}}^4)^2$, then we can find $C_1 < 1/(3\alpha_{\mathcal{P}}^4)$ and $C_2 > 1$ such that $C_2/C_1 < \sqrt{k}$. Apply Proposition 7.1, there exists C, δ' only depending on \mathcal{P}, C_1, C_2 such that if $n \geq Ckr(p_1 + p_2)$, \mathcal{X} satisfies RUB of order kr with constants C_1 and C_2 with probability at least $1 - \exp(-\delta'(p_1 + p_2))$. Note that C_1, C_2 only depends on \mathcal{P} , we can conclude that there exist constants $D(= Ck), \delta'$ only depending on \mathcal{P} such that if $n \geq Dr(p_1 + p_2)$, \mathcal{X} satisfies RUB of order kr with constants C_1, C_2 satisfying $C_2/C_1 \leq \sqrt{k}$.

Similarly, to the proof of Proposition 2.1, under the event that:

1. \mathcal{X} satisfies RUB of order kr and constants C_1, C_2 satisfying $C_2/C_1 < \sqrt{k}$,
2. $z \in \mathcal{Z}_1 \cap \mathcal{Z}_2$,

we can get (3.4) (we shall note that W_1 depends on \mathcal{P} , so its value can also depend on $\alpha_{\mathcal{P}}$). The probability that those events happen is at least $1 - 2/n - 5\exp(-\delta(p_1 + p_2))$ for $\delta \leq \min((6 - 2\sqrt{2\pi})^2/2, 1/12, 2 - \log 7, \delta')$.

7.5. *Proof of Theorem 2.4.* Without loss of generality, we assume that $p_1 \leq p_2$. We consider the class of rank- r matrices

$$\mathcal{F}_c = \{A \in \mathbb{R}^{p_1 \times p_2} : A_{ij} = 0, \text{ whenever } i \geq r + 1\}$$

namely the matrices with all nonzero entries in the first r rows. The model (1.1) become

$$y_i = \beta_{1:r}^{(i)T} A_r \gamma^{(i)} + z_i, \quad i = 1, \dots, n,$$

where $\beta_{1:r}^{(i)}$ is the vector of the first to the r th entries of $\beta^{(i)}$. Note that this is a linear regression model with variable $A_r \in \mathbb{R}^{r \times p_2}$, by Lemma 3.11 in [13], we have

$$(7.4) \quad \inf_{\hat{A}} \sup_{A \in \mathcal{F}_c} E \|\hat{A}(y) - A\|_F^2 = \sigma^2 \text{trace}[(\mathcal{X}_r^* \mathcal{X}_r)^{-1}],$$

$$(7.5) \quad \inf_{\hat{A}} \sup_{A \in \mathcal{F}_c} E \|\hat{A}(y) - A\|_F^2 = \infty \quad \text{when } \mathcal{X}_r^* \mathcal{X}_r \text{ is singular,}$$

where $\mathcal{X}_r: \mathbb{R}^{r \times p_2} \rightarrow \mathbb{R}^n$ is the \mathcal{X} constrained on \mathcal{F}_c , Then \mathcal{X}_r sends A_r to $(\beta_{1:r}^{(1)} A_r \gamma^{(1)}, \dots, \beta_{1:r}^{(n)} A_r \gamma^{(n)})^\top$. When $n < p_2 r$, \mathcal{X}_r is singular, hence we have (2.10).

When $n \geq p_2 r$, we can see in order to show (2.9), we only need to show $\text{trace}(\mathcal{X}_r^* \mathcal{X}_r) \geq \frac{p_2 r}{2n}$ with probability at least $1 - 26n^{-1}$. Suppose the singular value of \mathcal{X}_r are $\sigma_i(\mathcal{X}_r)$, $i = 1, \dots, rp_2$, then $\text{trace}(\mathcal{X}_r^* \mathcal{X}_r) = \sum_{i=1}^{p_2 r} \sigma_i^2(\mathcal{X}_r)$.

Suppose \mathcal{X} is ROP while $B \in \mathbb{R}^{r \times p_2}$ is i.i.d. standard Gaussian random matrix (both \mathcal{X} and B_r are random). Then by some calculation, we can see

$$E_{B, \mathcal{X}_r} \|\mathcal{X}_r(B)\|_2^2 = n E_{B, \beta, \gamma} (\beta_{1:r}^T B \gamma)^2 = n \sum_{j=1}^r \sum_{k=1}^{p_2} E(\beta_j B_{jk} \gamma_k)^2 = np_2 r.$$

Note (0.20) in the proof of Lemma 7.1 in the supplementary material [11], we know $E(\beta_{1:r}^{(i)T} B \gamma^{(i)} \|B\|_F^4) \leq 9 \|B\|_F^4$. Hence,

$$\begin{aligned} E \|\mathcal{X}_r(B)\|_2^4 &= \sum_{i=1}^n E(\beta_{1:r}^{(i)T} B \gamma^{(i)})^4 \\ &\quad + 2 \sum_{1 \leq i < l \leq n} E \sum_{j=1}^n (\beta_{1:r}^{(i)T} B \gamma^{(i)})^2 \cdot E \sum_{j=1}^n (\beta_{1:r}^{(l)T} B \gamma^{(l)})^2 \\ &= n \cdot 9 E \|B\|_F^4 + n(n-1)(p_2 r)^2 \\ &= 9n E(\chi^2(p_2 r))^2 + n(n-1)p_2^2 r^2 \\ &= 9n(p_2^2 r^2 + 2p_2 r) + n(n-1)p_2^2 r^2 \\ &= n^2 p_2^2 r^2 + 2np_2 r(4p_2 r + 9) \leq n^2 p_2^2 r^2 + 26np_2^2 r^2. \end{aligned}$$

Besides,

$$\begin{aligned} E \|\mathcal{X}_r(B_r)\|_2^2 &= E(E(\|\mathcal{X}_r(B_r)\|_2^2 | \mathcal{X}_r)) = E\left(\sum_{i=1}^{rp_2} \sigma_i^2(\mathcal{X}_r)\right), \\ E \|\mathcal{X}_r(B_r)\|_2^4 &= E(E(\|\mathcal{X}_r(B_r)\|_2^4 | \mathcal{X}_r)) \\ &= E\left(\sum_{i=1}^{rp_2} 3\sigma_i^4(\mathcal{X}_r) + 2 \sum_{1 \leq i < j \leq rp_2} \sigma_i^2(\mathcal{X}_r) \sigma_j^2(\mathcal{X}_r)\right) \\ &\geq E\left(\sum_{i=1}^{rp_2} \sigma_i^2(\mathcal{X}_r)^2\right)^2. \end{aligned}$$

Hence,

$$\begin{aligned} E\left(\sum_{i=1}^{rp_2} \sigma_i^2(\mathcal{X}_r)^2\right) &= np_2 r, \\ \text{Var}\left(\sum_{i=1}^{rp_2} \sigma_i^2(\mathcal{X}_r)^2\right) &= E\left(\sum_{i=1}^{rp_2} \sigma_i^2(\mathcal{X}_r)^2\right)^2 - \left(E \sum_{i=1}^{rp_2} \sigma_i^2(\mathcal{X}_r)^2\right)^2 \leq 26np_2^2 r^2. \end{aligned}$$

Then by Chebyshev's inequality, we have

$$(7.6) \quad \sum_{i=1}^{rp_2} \sigma_i^2(\mathcal{X}_r) \leq 2np_2r$$

with probability at least $1 - \frac{26np_2^2r^2}{(npr)^2} = 1 - \frac{26}{n}$. By Cauchy-Schwarz's inequality, we have

$$\text{trace}((\mathcal{X}_r^* \mathcal{X}_r)^{-1}) = \sum_{i=1}^{rp_2} \sigma_i^{-2}(\mathcal{X}_r) \geq \frac{(p_2r)^2}{\sum_{i=1}^{rp_2} \sigma_i^2(\mathcal{X}_r)}.$$

Therefore, we have

$$\text{trace}((\mathcal{X}_r^* \mathcal{X}_r)^{-1}) \geq \frac{p_2r}{2n}$$

with probability at least $1 - 26/n$, which shows (2.9).

Finally, we consider (2.8). Suppose inequality (7.6) holds, then

$$(7.7) \quad \begin{aligned} |\{i : \sigma_i^2(X_r) \geq 4n\}| &\leq \frac{p_2r}{2} \\ \Rightarrow \left| \left\{ i : \sigma_i^{-2}(X_r) \leq \frac{1}{4n} \right\} \right| &\geq \frac{p_2r}{2} \\ \Rightarrow \left| \left\{ i : \sigma_i^{-2}(X_r) \geq \frac{1}{4n} \right\} \right| &\geq \frac{p_2r}{2}. \end{aligned}$$

By Lemma 3.12 in [13], we know

$$\begin{aligned} &\inf_{\hat{A}} \sup_{A \in \mathcal{F}_c} P_z \left(\|\hat{A} - A\|_F^2 \geq \frac{p_2r\sigma^2}{16n} \right) \\ &= \inf_{\hat{A}} \sup_{A \in \mathcal{F}_c} E_z 1_{\{x \geq p_2r\sigma^2/16n\}} (\|\hat{A} - A\|_F^2) \\ &= E_z 1_{\{x \geq p_2r\sigma^2/16n\}} (\|(\mathcal{X}_r^* \mathcal{X}_r)^{-1} \mathcal{X}_r^*(z)\|_F^2) \\ &= P_z \left(\|(\mathcal{X}_r^* \mathcal{X}_r)^{-1} \mathcal{X}_r^*(z)\|_F^2 \geq \frac{p_2r\sigma^2}{16n} \right), \end{aligned}$$

where $1_{\{x \geq p_2r\sigma^2/16n\}}(\cdot)$ is the indicator function. Note that when $z \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$, $\|(\mathcal{X}_r^* \mathcal{X}_r)^{-1} \mathcal{X}_r^*(z)\|_F^2$ is identical distributed as $\sum_{i=1}^{rp_2} \frac{y_i^2}{\sigma_i^2(\mathcal{X}_r)}$, where $y_1, \dots, y_{rp_2} \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$, hence,

$$P \left(\|(\mathcal{X}_r^* \mathcal{X}_r)^{-1} \mathcal{X}_r^*(z)\|_F^2 \leq \frac{p_2r\sigma^2}{16n} \right)$$

$$\begin{aligned}
&= P\left(\sum_{i=1}^{rp_2} \frac{y_i^2}{\sigma_i^2(\mathcal{X}_r)} \leq \frac{p_2 r \sigma^2}{16n}\right) \\
&\leq P\left(\sum_{i: \sigma_i^{-2}(\mathcal{X}_r) \geq 1/(4n)} y_i^2 \sigma_i^{-2}(\mathcal{X}_r) \leq \frac{p_2 r \sigma^2}{16n}\right) \\
&\leq P\left(\sum_{i: \sigma_i^{-2}(\mathcal{X}_r) \geq 1/(4n)} \frac{y_i^2}{4n} \leq \frac{p_2 r \sigma^2}{16n}\right) \leq P\left(\chi^2\left(\left\lceil \frac{rp_2}{2} \right\rceil\right) \leq \frac{p_2 r}{4}\right) \\
&\leq \exp\left(-\frac{rp_2}{32}\right).
\end{aligned}$$

The last inequality is due to the tail bound of χ^2 distribution given by Lemma 1 in [28]; the second last inequality is due to (7.7). In summary, when (7.6) holds, we have

$$\inf_{\hat{A}} \sup_{A \in \mathcal{F}_c} P_z\left(\|\hat{A} - A\|_F^2 \geq \frac{p_2 r \sigma^2}{16n}\right) \leq \exp\left(-\frac{rp_2}{32}\right).$$

Finally, since $p_2 \geq (p_1 + p_2)/2$, we showed that with probability at least $1 - 26n^{-1}$, \mathcal{X} satisfies (2.8).

7.6. Proof of Theorem 4.1. We first introduce the following lemma about the upper bound of $\|z\|_1, \|z\|_2, \|z\|_\infty$.

LEMMA 7.10. *Suppose z is defined as (4.3), then for constants $C_1 > \sqrt{2}$, $M_1 > 1$, we have*

$$(7.8) \quad P\left(\|z\|_1/n \leq \frac{C_1}{n} \sum_{i=1}^n \xi_i^2\right) \geq 1 - \frac{9C_1^2 + 6}{n(C_1 - \sqrt{2})^2},$$

$$P\left(\frac{C_1}{n} \sum_{i=1}^n \xi_i^2 \leq M_1 C_1 \|\Sigma\|_*\right) \geq 1 - \frac{9}{n(M_1 - 1)^2};$$

for constants $C_2 > 1$, $M_2 > 9$,

$$(7.9) \quad P\left(\|z\|_2^2/n \leq \frac{C_2^2 \sum_{i=1}^n \xi_i^4}{n}\right) \geq 1 - \frac{105(105C_2^4 + 60)}{n(3C_2^2 - 2)^2},$$

$$P\left(\frac{C_2^2 \sum_{i=1}^n \xi_i^4}{n} \leq M_2 C_2^2 \|\Sigma\|_*^2\right) \geq 1 - \frac{105^2}{n(M_2 - 9)^2};$$

for constants $C_3 > 1$, $M_3 > 1$,

$$P\left(\|z\|_\infty \leq C_3 \log n \max_{1 \leq i \leq n} \xi_i^2\right) \geq 1 - \frac{2}{\sqrt{2\pi C_3 \log n}},$$

$$(7.10) \quad P\left(C_3 \log n \max_{1 \leq i \leq n} \xi_i^2 \leq 2C_3 M_3 \log^2 n (\sqrt{\|\Sigma\|_*} + \sqrt{2M_3 \log n \|\Sigma\|})^2\right) \\ \geq 1 - 2n^{-M_3+1}.$$

The proof of Lemma 7.10 is listed in the supplementary material [11]. The rest of the proof is basically the same as Proposition 2.3. Suppose $\mathcal{X}_1, \mathcal{X}_2$ and \tilde{z} are given by (0.36), (0.37) and (0.39) in the supplementary material [11], then $\mathcal{X}_1, \mathcal{X}_2$ are ROP. By Lemma 7.4,

$$(7.11) \quad \|\mathcal{X}_1^*(\tilde{z})\| \leq 6(2p\|\tilde{z}\|_\infty + \sqrt{2p}\|\tilde{z}\|_2),$$

$$(7.12) \quad \|\mathcal{X}_2^*(\tilde{z})\| \leq 6(2p\|\tilde{z}\|_\infty + \sqrt{2p}\|\tilde{z}\|_2)$$

with probability at least $1 - 4 \exp(-2(2 - \log 7)p)$. Hence, there exists $\delta > 0$ such that

$$\begin{aligned} & P(\Sigma_0 \text{ is NOT in the feasible set of (4.4)}) \\ &= P(\|z\|_1/n > \eta_1 \text{ or } \|\tilde{\mathcal{X}}^*(\tilde{z})\| > \eta_2) \\ &\leq P\left(\|z\|_1/n > \frac{c_1}{n} \sum_{i=1}^n \xi_i^2\right) + P\left(\|\tilde{z}\|_\infty > 2c_3 \log n \max_{1 \leq i \leq n} \xi_i^2\right) \\ &\quad + P\left(\|\tilde{z}\|_2 > c_2 \sqrt{2 \sum_{i=1}^n \xi_i^4}\right) \\ &\quad + P(\|\tilde{\mathcal{X}}^*(z)\| > 24p\|\tilde{z}\|_\infty + 12\sqrt{2p}\|\tilde{z}\|_2) \\ &\leq P\left(\|z\|_1/n > \frac{c_1}{n} \sum_{i=1}^n \xi_i^2\right) + P\left(\|z\|_\infty > c_3 \log n \max_{1 \leq i \leq n} \xi_i^2\right) \\ &\quad + P\left(\|z\|_2 > c_2 \sqrt{\sum_{i=1}^n \xi_i^4}\right) \\ &\quad + P(\|\mathcal{X}_1^*(\tilde{z})\| > 12p\|\tilde{z}\|_\infty + 6\sqrt{2p}\|\tilde{z}\|_2) \\ &\quad + P(\|\mathcal{X}_2^*(z)\| > 12p\|\tilde{z}\|_\infty + 6\sqrt{2p}\|\tilde{z}\|_2) \\ &\leq O(1/n) + 4 \exp(-2(2 - \log 7)p) + \frac{2}{\sqrt{2\pi c_3 \log n}}. \end{aligned}$$

Here, we used the fact that $\tilde{\mathcal{X}}^* = \mathcal{X}_1^* + \mathcal{X}_2^*$,

$$\|\tilde{z}\|_2 = \sqrt{\sum_{i=1}^{\lfloor n/2 \rfloor} (z_{2i-1} - z_{2i})^2} \leq \sqrt{\sum_{i=1}^{\lfloor n/2 \rfloor} 2(z_{2i-1}^2 + z_{2i}^2)} \leq \sqrt{2}\|z\|_2,$$

$$\|\tilde{z}\|_\infty = \max_i |z_{2i-1} - z_{2i}| \leq 2 \max_i |z_i| \leq 2\|z\|_\infty.$$

Similarly to the proof of Proposition 2.3, since \mathcal{X}_1 is ROP, there exists constants D and δ' such that if $n \geq Drp$, \mathcal{X}_1 satisfies RUB of order $10k$ with constants C_1, C_2 satisfying $C_2/C_1 < \sqrt{10}$ with probability at least $1 - e^{-n\delta'}$.

Now under the event that:

1. A is feasible in (4.4),
2. \mathcal{X}_1 satisfies RUB of order $10k$ with constants C_1, C_2 satisfying $C_2/C_1 < \sqrt{10}$,
3. the latter part of (7.8), (7.9) and (7.10) hold for some $M_1 > 1$, $M_2 > 9$, $M_3 > 2$,

we can prove (4.6) similarly as the proof of Proposition 2.3, which we omit the proof here.

Acknowledgments. We thank the Associate Editor and the referees for their thorough and useful comments which have helped to improve the presentation of the paper.

SUPPLEMENTARY MATERIAL

Supplement to “ROP: Matrix recovery via rank-one projections”: (DOI: [10.1214/14-AOS1267SUPP](https://doi.org/10.1214/14-AOS1267SUPP); .pdf). We prove the technical lemmas used in the proofs of the main results in this supplement. The proofs rely on results in [7, 13, 28, 36, 39, 41] and [31].

REFERENCES

- [1] ALQUIER, P., BUTUCEA, C., HEBIRI, M. and MEZIANI, K. (2013). Rank penalized estimation of a quantum system. *Phys. Rev. A* **88** 032133.
- [2] ANDREWS, H. C. and PATTERSON, C. L. III (1976). Singular value decomposition (SVD) image coding. *IEEE Trans. Commun.* **24** 425–432.
- [3] BASRI, R. and JACOBS, D. W. (2003). Lambertian reflectance and linear sub-spaces. *IEEE Trans. Pattern Anal. Mach. Intell.* **25** 218–233.
- [4] BIRNBAUM, A., JOHNSTONE, I. M., NADLER, B. and PAUL, D. (2013). Minimax bounds for sparse PCA with noisy high-dimensional data. *Ann. Statist.* **41** 1055–1084. [MR3113803](#)
- [5] CAI, T. T., MA, Z. and WU, Y. (2013). Sparse PCA: Optimal rates and adaptive estimation. *Ann. Statist.* **41** 3074–3110. [MR3161458](#)
- [6] CAI, T. T., MA, Z. and WU, Y. (2014). Optimal estimation and rank detection for sparse spiked covariance matrices. *Probab. Theory Related Fields*. To appear.
- [7] CAI, T. T., XU, G. and ZHANG, J. (2009). On recovery of sparse signals via ℓ_1 minimization. *IEEE Trans. Inform. Theory* **55** 3388–3397. [MR2598028](#)
- [8] CAI, T. T. and ZHANG, A. (2013). Sharp RIP bound for sparse signal and low-rank matrix recovery. *Appl. Comput. Harmon. Anal.* **35** 74–93. [MR3053747](#)
- [9] CAI, T. T. and ZHANG, A. (2013). Compressed sensing and affine rank minimization under restricted isometry. *IEEE Trans. Signal Process.* **61** 3279–3290. [MR3070321](#)

- [10] CAI, T. T. and ZHANG, A. (2014). Sparse representation of a polytope and recovery in sparse signals and low-rank matrices. *IEEE Trans. Inform. Theory* **60** 122–132. [MR3150915](#)
- [11] CAI, T. and ZHANG, A. (2014). Supplement to “ROP: Matrix recovery via rank-one projections.” DOI:[10.1214/14-AOS1267SUPP](#).
- [12] CANDÈS, E. J., LI, X., MA, Y. and WRIGHT, J. (2011). Robust principal component analysis? *J. ACM* **58** Art. 11, 37. [MR2811000](#)
- [13] CANDÈS, E. J. and PLAN, Y. (2011). Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Trans. Inform. Theory* **57** 2342–2359. [MR2809094](#)
- [14] CANDÈS, E. J. and RECHT, B. (2009). Exact matrix completion via convex optimization. *Found. Comput. Math.* **9** 717–772. [MR2565240](#)
- [15] CANDÈS, E. J., STROHMER, T. and VORONINSKI, V. (2013). PhaseLift: Exact and stable signal recovery from magnitude measurements via convex programming. *Comm. Pure Appl. Math.* **66** 1241–1274. [MR3069958](#)
- [16] CANDÈS, E. J. and TAO, T. (2010). The power of convex relaxation: Near-optimal matrix completion. *IEEE Trans. Inform. Theory* **56** 2053–2080. [MR2723472](#)
- [17] CHEN, Y., CHI, Y. and GOLDSMITH, A. (2013). Exact and stable covariance estimation from quadratic sampling via convex programming. Preprint. Available at [arXiv:1310.0807](#).
- [18] DASARATHY, G., SHAH, P., BHASKAR, B. N. and NOWAK, R. (2012). Covariance sketching. In *50th Annual Allerton Conference on Communication, Control, and Computing* 1026–1033.
- [19] DASARATHY, G., SHAH, P., BHASKAR, B. N. and NOWAK, R. (2013). Sketching sparse matrices. Preprint. Available at [arXiv:1303.6544](#).
- [20] DVIJOTHAM, K. and FAZEL, M. (2010). A nullspace analysis of the nuclear norm heuristic for rank minimization. In *2010 IEEE International Conference on Acoustics Speech and Signal Processing (ICASSP)* 3586–3589.
- [21] FAN, J., FAN, Y. and LV, J. (2008). High dimensional covariance matrix estimation using a factor model. *J. Econometrics* **147** 186–197. [MR2472991](#)
- [22] GRANT, M. and BOYD, S. (2012). CVX: Matlab software for disciplined convex programming, version 2.0 beta. Available at <http://cvxr.com/cvx>.
- [23] GRANT, M. C. and BOYD, S. P. (2008). Graph implementations for nonsmooth convex programs. In *Recent Advances in Learning and Control (a tribute to M. Vidyasagar)* (V. BLONDEL ET AL., eds.). *Lecture Notes in Control and Inform. Sci.* **371** 95–110. Springer, London. [MR2409077](#)
- [24] GROSS, D., LIU, Y. K., FLAMMIA, S. T., BECKER, S. and EISERT, J. (2010). Quantum state tomography via compressed sensing. *Phys. Rev. Lett.* **105** 150401–150404.
- [25] JOHNSTONE, I. M. (2001). On the distribution of the largest eigenvalue in principal components analysis. *Ann. Statist.* **29** 295–327. [MR1863961](#)
- [26] 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](#)
- [27] KOREN, Y., BELL, R. and VOLINSKY, C. (2009). Matrix factorization techniques for recommender systems. *Computer* **42** 30–37.
- [28] LAURENT, B. and MASSART, P. (2000). Adaptive estimation of a quadratic functional by model selection. *Ann. Statist.* **28** 1302–1338. [MR1805785](#)

- [29] NADLER, B. (2010). Nonparametric detection of signals by information theoretic criteria: Performance analysis and an improved estimator. *IEEE Trans. Signal Process.* **58** 2746–2756. [MR2789420](#)
- [30] NEGAHBAN, S. and WAINWRIGHT, M. J. (2011). Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *Ann. Statist.* **39** 1069–1097. [MR2816348](#)
- [31] OYMAK, S. and HASSIBI, B. (2010). New null space results and recovery thresholds for matrix rank minimization. Preprint. Available at [arXiv:1011.6326](#).
- [32] binproceedings OYMAK, S., MOHAN, K., FAZEL, M. and HASSIBI, B. (2011). A simplified approach to recovery conditions for low-rank matrices. In *Proc. Intl. Sympo. Information Theory (ISIT)* 2318–2322. IEEE, Piscataway, NJ.
- [33] PATTERSON, N., PRICE, A. L. and REICH, D. (2006). Population structure and eigenanalysis. *PLoS Genet.* **2** e190.
- [34] PRICE, A. L., PATTERSON, N. J., PLENGE, R. M., WEINBLATT, M. E., SHADICK, N. A. and REICH, D. (2006). Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* **38** 904–909.
- [35] RECHT, B. (2011). A simpler approach to matrix completion. *J. Mach. Learn. Res.* **12** 3413–3430. [MR2877360](#)
- [36] RECHT, B., FAZEL, M. and PARRILO, P. A. (2010). Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Rev.* **52** 471–501. [MR2680543](#)
- [37] ROHDE, A. and TSYBAKOV, A. B. (2011). Estimation of high-dimensional low-rank matrices. *Ann. Statist.* **39** 887–930. [MR2816342](#)
- [38] TROSSET, M. W. (2000). Distance matrix completion by numerical optimization. *Comput. Optim. Appl.* **17** 11–22. [MR1791595](#)
- [39] VERSHYNIN, R. (2011). Spectral norm of products of random and deterministic matrices. *Probab. Theory Related Fields* **150** 471–509. [MR2824864](#)
- [40] WAKIN, M., LASKA, J., DUARTE, M., BARON, D., SARVOTHAM, S., TAKHAR, D., KELLY, K. and BARANIUK, R. (2006). An architecture for compressive imaging. In *Proceedings of the International Conference on Image Processing (ICIP 2006)* 1273–1276.
- [41] WANG, H. and LI, S. (2013). The bounds of restricted isometry constants for low rank matrices recovery. *Sci. China Ser. A* **56** 1117–1127.
- [42] WANG, Y. (2013). Asymptotic equivalence of quantum state tomography and noisy matrix completion. *Ann. Statist.* **41** 2462–2504. [MR3127872](#)
- [43] WAX, M. and KAILATH, T. (1985). Detection of signals by information theoretic criteria. *IEEE Trans. Acoust. Speech Signal Process.* **33** 387–392. [MR0788604](#)

DEPARTMENT OF STATISTICS
 THE WHARTON SCHOOL
 UNIVERSITY OF PENNSYLVANIA
 PHILADELPHIA, PENNSYLVANIA 19104
 USA
 E-MAIL: tcai@wharton.upenn.edu
anzhang@wharton.upenn.edu
 URL: <http://www-stat.wharton.upenn.edu/~tcai/>
<http://www-stat.wharton.upenn.edu/~anzhang/>