

ESTIMATION OF (NEAR) LOW-RANK MATRICES WITH NOISE AND HIGH-DIMENSIONAL SCALING

BY SAHAND NEGAHBAN AND MARTIN J. WAINWRIGHT

University of California, Berkeley

We study an instance of high-dimensional inference in which the goal is to estimate a matrix $\Theta^* \in \mathbb{R}^{m_1 \times m_2}$ on the basis of N noisy observations. The unknown matrix Θ^* is assumed to be either exactly low rank, or “near” low-rank, meaning that it can be well-approximated by a matrix with low rank. We consider a standard M -estimator based on regularization by the nuclear or trace norm over matrices, and analyze its performance under high-dimensional scaling. We define the notion of restricted strong convexity (RSC) for the loss function, and use it to derive non-asymptotic bounds on the Frobenius norm error that hold for a general class of noisy observation models, and apply to both exactly low-rank and approximately low rank matrices. We then illustrate consequences of this general theory for a number of specific matrix models, including low-rank multivariate or multi-task regression, system identification in vector autoregressive processes, and recovery of low-rank matrices from random projections. These results involve non-asymptotic random matrix theory to establish that the RSC condition holds, and to determine an appropriate choice of regularization parameter. Simulation results show excellent agreement with the high-dimensional scaling of the error predicted by our theory.

1. Introduction. High-dimensional inference refers to instances of statistical estimation in which the ambient dimension of the data is comparable to (or possibly larger than) the sample size. Problems with a high-dimensional character arise in a variety of applications in science and engineering, including analysis of gene array data, medical imaging, remote sensing, and astronomical data analysis. In settings where the number of parameters may be large relative to the sample size, the utility of classical (fixed dimension) results is questionable, and accordingly, a line of on-going statistical research seeks to obtain results that hold under high-dimensional scaling, meaning that both the problem size and sample size (as well as other problem parameters) may tend to infinity simultaneously. It is usually impossible to obtain consistent procedures in such settings without imposing some sort of additional constraints. Accordingly, there are now various lines of work on high-dimensional inference based on imposing different types of structural constraints. A substantial body of past work has focused on mod-

els with sparsity constraints, including the problem of sparse linear regression [52, 16, 18, 40, 10], banded or sparse covariance matrices [7, 8, 19], sparse inverse covariance matrices [57, 24, 50, 46], sparse eigenstructure [30, 2, 44], and sparse regression matrices [43, 36, 56, 28]. A theme common to much of this work is the use of ℓ_1 -penalty as a surrogate function to enforce the sparsity constraint. A parallel line of work has focused on the use of concave penalties to achieve gains in model selection and sparsity recovery [20, 21].

In this paper, we focus on the problem of high-dimensional inference in the setting of matrix estimation. As mentioned above, there is already a substantial body of work on the problem of sparse matrix recovery. In contrast, our interest in this paper is the problem of estimating a matrix $\Theta^* \in \mathbb{R}^{m_1 \times m_2}$ that is either *exactly low rank*, meaning that it has at most $r \ll \min\{m_1, m_2\}$ non-zero singular values, or more generally is *near low-rank*, meaning that it can be well-approximated by a matrix of low rank. As we discuss at more length in the sequel, such exact or approximate low-rank conditions are appropriate for many applications, including multivariate or multi-task forms of regression, system identification for autoregressive processes, collaborative filtering, and matrix recovery from random projections. Analogous to the use of an ℓ_1 -regularizer for enforcing sparsity, we consider the use of the nuclear norm (also known as the trace norm) for enforcing a rank constraint in the matrix setting. By definition, the nuclear norm is the sum of the singular values of a matrix, and so encourages sparsity in the vector of singular values, or equivalently for the matrix to be low-rank. The problem of low-rank matrix approximation and the use of nuclear norm regularization have been studied by various researchers. In her Ph.D. thesis, Fazel [22] discusses the use of nuclear norm as a heuristic for restricting the rank of a matrix, showing that in practice it is often able to yield low-rank solutions. Other researchers have provided theoretical guarantees on the performance of nuclear norm and related methods for low-rank matrix approximation. Srebro et al. [51] proposed nuclear norm regularization for the collaborative filtering problem, and established risk consistency under certain settings. Recht et al. [47] provided sufficient conditions for exact recovery using the nuclear norm heuristic when observing random projections of a low-rank matrix, a set-up analogous to the compressed sensing model in sparse linear regression [18, 14]. Other researchers have studied a version of matrix completion in which a subset of entries are revealed, and the goal is to obtain perfect reconstruction either via the nuclear norm heuristic [15] or by other SVD-based methods [31]. For general observation models, Bach [6] has provided results on the consistency of nuclear norm minimization in noisy settings, but applicable to the classical “fixed p ” setting. In addition, Yuan

et al. [55] provide non-asymptotic bounds on the operator norm error of the estimate in the multi-task setting, provided that the design matrices are orthogonal. Under the assumption of RIP, Lee and Bresler [34] prove stability properties of least-squares under nuclear norm constraint when a form of restricted isometry property is imposed on the sampling operator. Liu and Vandenberghe [35] develop an efficient interior-point method for solving nuclear-norm constrained problems, and illustrate its usefulness for problems of system identification, an application also considered in this paper. Finally, in work posted shortly after our own, Rohde and Tsybakov [49] and Candes and Plan [13] have studied certain aspects of nuclear norm minimization under high-dimensional scaling. We discuss connections to this concurrent work at more length in Section 3.2 following the statement of our main results.

The goal of this paper is to analyze the nuclear norm relaxation for a general class of noisy observation models, and obtain non-asymptotic error bounds on the Frobenius norm that hold under high-dimensional scaling, and are applicable to both exactly and approximately low-rank matrices. We begin by presenting a generic observation model, and illustrating how it can be specialized to the several cases of interest, including low-rank multivariate regression, estimation of autoregressive processes, and random projection (compressed sensing) observations. In particular, this model is specified in terms of an operator \mathfrak{X} , which may be deterministic or random depending on the setting, that maps any matrix $\Theta^* \in \mathbb{R}^{m_1 \times m_2}$ to a vector of N noisy observations. We then present a single main theorem (Theorem 1) followed by two corollaries that cover the cases of exact low-rank constraints (Corollary 1) and near low-rank constraints (Corollary 2) respectively. These results demonstrate that high-dimensional error rates are controlled by two key quantities. First, the (random) observation operator \mathfrak{X} is required to satisfy a condition known as *restricted strong convexity* (RSC), introduced in a more general setting by Negahban et al. [41], which ensures that the loss function has sufficient curvature to guarantee consistent recovery of the unknown matrix Θ^* . As we show via various examples, this RSC condition is weaker than the RIP property, which requires that the sampling operator behave very much like an isometry on low-rank matrices. Second, our theory provides insight into the *choice of regularization parameter* that weights the nuclear norm, showing that an appropriate choice is to set it proportional to the spectral norm of a random matrix defined by the adjoint of observation operator \mathfrak{X} , and the observation noise in the problem.

This initial set of results, though appealing in terms of their simple statements and generality, are somewhat abstractly formulated. Our next con-

tribution is to show that by specializing our main result (Theorem 1) to three classes of models, we can obtain some concrete results based on readily interpretable conditions. In particular, Corollary 3 deals with the case of low-rank multivariate regression, relevant for applications in multitask learning. We show that the random operator \mathfrak{X} satisfies the RSC property for a broad class of observation models, and we use random matrix theory to provide an appropriate choice of the regularization parameter. Our next result, Corollary 4, deals with the case of estimating the matrix of parameters specifying a vector autoregressive (VAR) process [4, 37]. The usefulness of the nuclear norm in this context has been demonstrated by Liu and Vandenberghe [35]. Here we also establish that a suitable RSC property holds with high probability for the random operator \mathfrak{X} , and also specify a suitable choice of the regularization parameter. We note that the technical details here are considerably more subtle than the case of low-rank multivariate regression, due to dependencies introduced by the autoregressive sampling scheme. Accordingly, in addition to terms that involve the size, the matrix dimensions and rank, our bounds also depend on the mixing rate of the VAR process. Finally, we turn to the compressed sensing observation model for low-rank matrix recovery, as introduced by Recht and colleagues [47, 48]. In this setting, we again establish that the RSC property holds with high probability, specify a suitable choice of the regularization parameter, and thereby obtain a Frobenius error bound for noisy observations (Corollary 5). A technical result that we prove en route—namely, Proposition 1—is of possible independent interest, since it provides a bound on the constrained norm of a random Gaussian operator. In particular, this proposition allows us to obtain a sharp result (Corollary 6) for the problem of recovering a low-rank matrix from perfectly observed random Gaussian projections with a general dependency structure.

The remainder of this paper is organized as follows. Section 2 is devoted to background material, and the set-up of the problem. We present a generic observation model for low-rank matrices, and then illustrate how it captures various cases of interest. We then define the convex program based on nuclear norm regularization that we analyze in this paper. In Section 3, we state our main theoretical results and discuss their consequences for different model classes. Section 4 is devoted to the proofs of our results; in each case, we break down the key steps in a series of lemmas, with more technical details deferred to the appendices. In Section 5, we present the results of various simulations that illustrate excellent agreement between the theoretical bounds and empirical behavior.

Notation: For the convenience of the reader, we collect standard pieces

of notation here. For a pair of matrices Θ and Γ with commensurate dimensions, we let $\langle\langle \Theta, \Gamma \rangle\rangle = \text{trace}(\Theta^T \Gamma)$ denote the trace inner product on matrix space. For a matrix $\Theta \in \mathbb{R}^{m_1 \times m_2}$, we define $m = \min\{m_1, m_2\}$, and denote its (ordered) singular values by $\sigma_1(\Theta) \geq \sigma_2(\Theta) \geq \dots \geq \sigma_m(\Theta) \geq 0$. We also use the notation $\sigma_{\max}(\Theta) = \sigma_1(\Theta)$ and $\sigma_{\min}(\Theta) = \sigma_m(\Theta)$ to refer to the maximal and minimal singular values respectively. We use the notation $\|\cdot\|$ for various types of matrix norms based on these singular values, including the *nuclear norm* $\|\Theta\|_1 = \sum_{j=1}^m \sigma_j(\Theta)$, the *spectral or operator norm* $\|\Theta\|_{\text{op}} = \sigma_1(\Theta)$, and the *Frobenius norm* $\|\Theta\|_F = \sqrt{\text{trace}(\Theta^T \Theta)} = \sqrt{\sum_{j=1}^m \sigma_j^2(\Theta)}$. We refer the reader to Horn and Johnson [26, 27] for more background on these matrix norms and their properties.

2. Background and problem set-up. We begin with some background on problems and applications in which rank constraints arise, before describing a generic observation model. We then introduce the semidefinite program (SDP) based on nuclear norm regularization that we study in this paper.

2.1. Models with rank constraints. Imposing a rank r constraint on a matrix $\Theta^* \in \mathbb{R}^{m_1 \times m_2}$ is equivalent to requiring the rows (or columns) of Θ^* lie in some r -dimensional subspace of \mathbb{R}^{m_2} (or \mathbb{R}^{m_1} respectively). Such types of rank constraints (or approximate forms thereof) arise in a variety of applications, as we discuss here. In some sense, rank constraints are a generalization of sparsity constraints; rather than assuming that the data is sparse in a known basis, a rank constraint implicitly imposes sparsity but without assuming the basis.

We first consider the problem of multivariate regression, also referred to as multi-task learning in statistical machine learning. The goal of *multivariate regression* is to estimate a prediction function that maps covariates $Z_j \in \mathbb{R}^m$ to multi-dimensional output vectors $Y_j \in \mathbb{R}^{m_1}$. More specifically, let us consider the linear model, specified by a matrix $\Theta^* \in \mathbb{R}^{m_1 \times m_2}$, of the form

$$(1) \quad Y_a = \Theta^* Z_a + W_a, \quad \text{for } a = 1, \dots, n,$$

where $\{W_a\}_{a=1}^n$ is an i.i.d. sequence of m_1 -dimensional zero-mean noise vectors. Given a collection of observations $\{Z_a, Y_a\}_{a=1}^n$ of covariate-output pairs, our goal is to estimate the unknown matrix Θ^* . This type of model has been used in many applications, including analysis of fMRI image data [25], analysis of EEG data decoding [3], neural response modeling [12] and analysis of financial data. This model and closely related ones also arise in the problem of collaborative filtering [51], in which the goal is to predict users' preferences

for items (such as movies or music) based on their and other users' ratings of related items. The papers [1, 5] discuss additional instances of low-rank decompositions. In all of these settings, the low-rank condition translates into the existence of a smaller set of “features” that are actually controlling the prediction.

As a second (not unrelated) example, we now consider the problem of system identification in vector autoregressive processes (see the book [37] for detailed background). A *vector autoregressive* (VAR) process in m -dimensions is a stochastic process $\{Z_t\}_{t=1}^{\infty}$ specified by an initialization $Z_1 \in \mathbb{R}^m$, followed by the recursion

$$(2) \quad Z_{t+1} = \Theta^* Z_t + W_t, \quad \text{for } t = 1, 2, 3, \dots$$

In this recursion, the sequence $\{W_t\}_{t=1}^{\infty}$ consists of i.i.d. samples of innovations noise. We assume that each vector $W_t \in \mathbb{R}^m$ is zero-mean with covariance matrix $C \succ 0$, so that the process $\{Z_t\}_{t=1}^{\infty}$ is zero-mean, and has a covariance matrix Σ given by the solution of the discrete-time Ricatti equation

$$(3) \quad \Sigma = \Theta^* \Sigma (\Theta^*)^T + C.$$

The goal of system identification in a VAR process is to estimate the unknown matrix $\Theta^* \in \mathbb{R}^{m \times m}$ on the basis of a sequence of samples $\{Z_t\}_{t=1}^n$. In many application domains, it is natural to expect that the system is controlled primarily by a low-dimensional subset of variables. For instance, models of financial data might have an ambient dimension m of thousands (including stocks, bonds, and other financial instruments), but the behavior of the market might be governed by a much smaller set of macro-variables (combinations of these financial instruments). Similar statements apply to other types of time series data, including neural data [12, 23], subspace tracking models in signal processing, and motion models in computer vision. While the form of system identification formulated here assumes direct observation of the state variables $\{Z_t\}_{t=1}^n$, it is also possible to tackle the more general problem when only noisy versions are observed (e.g., see Liu and Vandenberghe [35]). An interesting feature of the system identification problem is that the matrix Θ^* , in addition to having low rank, might also be required to satisfy some type of structural constraint (e.g., having a Hankel-type structure), and the estimator that we consider here allows for this possibility.

A third example that we consider in this paper is a *compressed sensing* observation model, in which one observes random projections of the unknown

matrix Θ^* . This observation model has been studied extensively in the context of estimating sparse vectors [18, 14], and Recht and colleagues [47, 48] suggested and studied its extension to low-rank matrices. In their set-up, one observes trace inner products of the form $\langle\langle X_i, \Theta^* \rangle\rangle = \text{trace}(X_i^T \Theta^*)$, where $X_i \in \mathbb{R}^{m_1 \times m_2}$ is a random matrix (for instance, filled with standard normal $N(0, 1)$ entries), so that $\langle\langle X_i, \Theta^* \rangle\rangle$ is a standard random projection. In the sequel, we consider this model with a more general family of random projections involving matrices with dependent entries. Like compressed sensing for sparse vectors, applications of this model include computationally efficient updating in large databases (where the matrix Θ^* measures the difference between the data base at two different time instants), and matrix denoising.

2.2. A generic observation model. We now introduce a generic observation model that will allow us to deal with these different observation models in an unified manner. For pairs of matrices $A, B \in \mathbb{R}^{m_1 \times m_2}$, recall the Frobenius or trace inner product $\langle\langle A, B \rangle\rangle := \text{trace}(BA^T)$. We then consider a linear observation model of the form

$$(4) \quad y_i = \langle\langle X_i, \Theta^* \rangle\rangle + \varepsilon_i, \quad \text{for } i = 1, 2, \dots, N,$$

which is specified by the sequence of observation matrices $\{X_i\}_{i=1}^N$ and observation noise $\{\varepsilon_i\}_{i=1}^N$. This observation model can be written in a more compact manner using operator-theoretic notation. In particular, let us define the observation vector

$$\vec{y} = [y_1 \quad \dots \quad y_N]^T \in \mathbb{R}^N,$$

with a similar definition for $\vec{\varepsilon} \in \mathbb{R}^N$ in terms of $\{\varepsilon_i\}_{i=1}^N$. We then use the observation matrices $\{X_i\}_{i=1}^N$ to define an operator $\mathfrak{X} : \mathbb{R}^{m_1 \times m_2} \rightarrow \mathbb{R}^N$ via $[\mathfrak{X}(\Theta)]_i = \langle\langle X_i, \Theta \rangle\rangle$. With this notation, the observation model (4) can be re-written as

$$(5) \quad \vec{y} = \mathfrak{X}(\Theta^*) + \vec{\varepsilon}.$$

Let us illustrate the form of the observation model (5) for some of the applications that we considered earlier.

EXAMPLE 1 (Multivariate regression). Recall the observation model (1) for multivariate regression. In this case, we make n observations of vector pairs $(Y_a, Z_a) \in \mathbb{R}^{m_1} \times \mathbb{R}^{m_2}$. Accounting for the m_1 -dimensional nature of the output, after the model is scalarized, we receive a total of $N = m_1 n$ observations. Let us introduce the quantity $b = 1, \dots, m_1$ to index the different

elements of the output, so that we can write

$$(6) \quad Y_{ab} = \langle\langle e_b Z_a^T, \Theta^* \rangle\rangle + W_{ab}, \quad \text{for } b = 1, 2, \dots, m_1.$$

By re-indexing this collection of $N = nm_1$ observations via the mapping

$$(a, b) \mapsto i = a + (b - 1) m_1,$$

we recognize multivariate regression as an instance of the observation model (4) with observation matrix $X_i = e_b Z_a^T$ and scalar observation $y_i = Y_{ab}$.

EXAMPLE 2 (Vector autoregressive processes). Recall that a vector autoregressive (VAR) process is defined by the recursion (2), and suppose that we observe an n -sequence $\{Z_t\}_{t=1}^n$ produced by this recursion. Since each $Z_t = [Z_{t1} \ \dots \ Z_{tm}]^T$ is m -variate, the scalarized sample size is $N = nm$. Letting $b = 1, 2, \dots, m$ index the dimension, we have

$$(7) \quad Z_{(t+1)b} = \langle\langle e_b Z_t^T, \Theta^* \rangle\rangle + W_{tb}.$$

In this case, we re-index the collection of $N = nm$ observations via the mapping

$$(t, b) \mapsto i = t + (b - 1) m.$$

After doing so, we see that the autoregressive problem can be written in the form (4) with $y_i = Z_{(t+1)b}$ and observation matrix $X_i = e_b Z_t^T$.

EXAMPLE 3 (Compressed sensing). As mentioned earlier, this is a natural extension of the compressed sensing observation model for sparse vectors to the case of low-rank matrices [47, 48]. In a typical form of compressed sensing, the observation matrix $X_i \in \mathbb{R}^{m_1 \times m_2}$ has i.i.d. standard normal $N(0, 1)$ entries, so that one makes observations of the form

$$(8) \quad y_i = \langle\langle X_i, \Theta^* \rangle\rangle + \varepsilon_i, \quad \text{for } i = 1, 2, \dots, N.$$

By construction, these observations are an instance of the model (4). In the sequel, we study a more general observation model, in which the entries of X_i are allowed to have general Gaussian dependencies. For this problem, the more compact form (5) involves a random Gaussian operator mapping $\mathbb{R}^{m_1 \times m_2}$ to \mathbb{R}^N , and we study some of its properties in the sequel.

2.3. *Regression with nuclear norm regularization.* We now consider an estimator that is naturally suited to the problems described in the previous section. Recall that the *nuclear or trace norm* of a matrix $\Theta \in \mathbb{R}^{m_1 \times m_2}$ is given by $\|\Theta\|_1 = \sum_{j=1}^m \sigma_j(\Theta)$, corresponding to the sum of its singular values. Given a collection of observations $(y_i, X_i) \in \mathbb{R} \times \mathbb{R}^{m_1 \times m_2}$, for $i = 1, \dots, N$ from the observation model (4), we consider estimating the unknown $\Theta^* \in \mathcal{S}$ by solving the following optimization problem

$$(9) \quad \hat{\Theta} \in \arg \min_{\Theta \in \mathcal{S}} \left\{ \frac{1}{2N} \|\vec{y} - \mathfrak{X}(\Theta)\|_2^2 + \lambda_N \|\Theta\|_1 \right\},$$

where \mathcal{S} is a convex subset of $\mathbb{R}^{m_1 \times m_2}$, and $\lambda_N > 0$ is a regularization parameter. When $\mathcal{S} = \mathbb{R}^{m_1 \times m_2}$, the optimization problem (9) can be viewed as the analog of the Lasso estimator [52], tailored to low-rank matrices as opposed to sparse vectors. We include the possibility of a more general convex set \mathcal{S} since they arise naturally in certain applications (e.g., Hankel-type constraints in system identification [35]). When \mathcal{S} is a polytope (with $\mathcal{S} = \mathbb{R}^{m_1 \times m_2}$ as a special case), then the optimization problem (9) can be solved in time polynomial in the sample size N and the matrix dimensions m_1 and m_2 . Indeed, the optimization problem (9) is an instance of a *semidefinite program* [53], a class of convex optimization problems that can be solved efficiently by various polynomial-time algorithms [11]. For instance, Liu and Vandenberghe [35] develop an efficient interior point method for solving constrained versions of nuclear norm programs. Moreover, as we discuss in Section 5, there are a variety of first-order methods for solving the semidefinite program (SDP) defining our M -estimator [42, 29]. These first-order methods are well-suited to the high-dimensional problems arising in statistical settings, and we make use of one in performing our simulations.

Like in any typical M -estimator for statistical inference, the regularization parameter λ_N is specified by the statistician. As part of the theoretical results in the next section, we provide suitable choices of this parameter so that the estimate $\hat{\Theta}$ is close in Frobenius norm to the unknown matrix Θ^* . The setting of the regularizer depends on the knowledge of the noise variance. While in general one might need to estimate this parameter through cross validation [20, 9], we assume knowledge of the noise variance in order to most succinctly demonstrate the empirical behavior of our results through the experiments.

3. Main results and some consequences. In this section, we state our main results and discuss some of their consequences. Section 3.1 is devoted to results that apply to generic instances of low-rank problems, whereas Section 3.3 is devoted to the consequences of these results for more

specific problem classes, including low-rank multivariate regression, estimation of vector autoregressive processes, and recovery of low-rank matrices from random projections.

3.1. Results for general model classes. We begin by introducing the key technical condition that allows us to control the error $\widehat{\Theta} - \Theta^*$ between an SDP solution $\widehat{\Theta}$ and the unknown matrix Θ^* . We refer to it as the *restricted strong convexity* condition [41], since it amounts to guaranteeing that the quadratic loss function in the SDP (9) is strictly convex over a restricted set of directions. Letting $\mathcal{C} \subseteq \mathbb{R}^{m_1 \times m_2}$ denote the restricted set of directions, we say that the operator \mathfrak{X} satisfies restricted strong convexity (RSC) over the set \mathcal{C} if there exists some $\kappa(\mathfrak{X}) > 0$ such that

$$(10) \quad \frac{1}{2N} \|\mathfrak{X}(\Delta)\|_2^2 \geq \kappa(\mathfrak{X}) \|\Delta\|_F^2 \quad \text{for all } \Delta \in \mathcal{C}.$$

We note that analogous conditions have been used to establish error bounds in the context of sparse linear regression [10, 17], in which case the set \mathcal{C} corresponded to certain subsets of sparse vectors. These types of conditions are weaker than restricted isometry properties, since they involve only lower bounds on the operator \mathfrak{X} , and the constant $\kappa(\mathfrak{X})$ can be arbitrarily small.

Of course, the definition (10) hinges on the choice of the restricted set \mathcal{C} . In order to specify some appropriate sets for the case of (near) low-rank matrices, we require some additional notation. Any matrix $\Theta^* \in \mathbb{R}^{m_1 \times m_2}$ has a singular value decomposition of the form $\Theta^* = UDV^T$, where $U \in \mathbb{R}^{m_1 \times m_1}$ and $V \in \mathbb{R}^{m_2 \times m_2}$ are orthonormal matrices. For each integer $r \in \{1, 2, \dots, m\}$, we let $U^r \in \mathbb{R}^{m_1 \times r}$ and $V^r \in \mathbb{R}^{m_2 \times r}$ be the sub-matrices of singular vectors associated with the top r singular values of Θ^* . We then define the following two subspaces of $\mathbb{R}^{m_1 \times m_2}$

$$(11a) \quad \mathcal{A}(U^r, V^r) := \{\Delta \in \mathbb{R}^{m_1 \times m_2} \mid \text{row}(\Delta) \subseteq V^r \text{ and } \text{col}(\Delta) \subseteq U^r\}, \quad \text{and}$$

$$(11b) \quad \mathcal{B}(U^r, V^r) := \{\Delta \in \mathbb{R}^{m_1 \times m_2} \mid \text{row}(\Delta) \perp V^r \text{ and } \text{col}(\Delta) \perp U^r\},$$

where $\text{row}(\Delta) \subseteq \mathbb{R}^{m_2}$ and $\text{col}(\Delta) \subseteq \mathbb{R}^{m_1}$ denote the row space and column space, respectively, of the matrix Δ . When (U^r, V^r) are clear from the context, we adopt the shorthand notation \mathcal{A}^r and \mathcal{B}^r .

We can now define the subsets of interest. Let $\Pi_{\mathcal{B}^r}$ denote the projection operator onto the subspace \mathcal{B}^r , and define $\Delta'' = \Pi_{\mathcal{B}^r}(\Delta)$ and $\Delta' = \Delta - \Delta''$. For a positive integer $r \leq m = \min\{m_1, m_2\}$ and a tolerance parameter $\delta \geq$

0, consider the following subset of matrices

$$(12) \quad \mathcal{C}(r; \delta) := \left\{ \Delta \in \mathbb{R}^{m_1 \times m_2} \mid \|\Delta\|_F \geq \delta, \|\Delta''\|_1 \leq 3\|\Delta'\|_1 + 4 \sum_{j=r+1}^m \sigma_j(\Theta^*) \right\}.$$

Note that this set corresponds to matrices Δ for which the quantity $\|\Delta''\|_1$ is relatively small compared to $\Delta - \Delta''$ and the remaining $m - r$ singular values of Θ^* .

The next ingredient is the choice of the regularization parameter λ_N used in solving the SDP (9). Our theory specifies a choice for this quantity in terms of the adjoint of the operator \mathfrak{X} —namely, the operator $\mathfrak{X}^* : \mathbb{R}^N \rightarrow \mathbb{R}^{m_1 \times m_2}$ defined by

$$(13) \quad \mathfrak{X}^*(\vec{\varepsilon}) := \sum_{i=1}^N \varepsilon_i X_i.$$

With this notation, we come to the first result of our paper. It is a deterministic result, which specifies two conditions—namely, an RSC condition and a choice of the regularizer—that suffice to guarantee that any solution of the SDP (9) falls within a certain radius.

THEOREM 1. *Suppose $\Theta^* \in \mathcal{S}$ and that the operator \mathfrak{X} satisfies restricted strong convexity with parameter $\kappa(\mathfrak{X}) > 0$ over the set $\mathcal{C}(r; \delta)$, and that the regularization parameter λ_N is chosen such that $\lambda_N \geq 2\|\mathfrak{X}^*(\vec{\varepsilon})\|_{\text{op}}/N$. Then any solution $\hat{\Theta}$ to the semidefinite program (9) satisfies*

$$(14) \quad \|\hat{\Theta} - \Theta^*\|_F \leq \max \left\{ \delta, \frac{32\lambda_N \sqrt{r}}{\kappa(\mathfrak{X})}, \left[\frac{16 \lambda_N \sum_{j=r+1}^m \sigma_j(\Theta^*)}{\kappa(\mathfrak{X})} \right]^{1/2} \right\}.$$

Apart from the tolerance parameter δ , the two main terms in the bound (14) have a natural interpretation. The first term (involving \sqrt{r}) corresponds to *estimation error*, capturing the difficulty of estimating a rank r matrix. The second is an *approximation error* that describes the gap between the true matrix Θ^* and the best rank r approximation. Understanding the magnitude of the tolerance parameter δ is a bit more subtle, and it depends on the geometry of the set $\mathcal{C}(r; \delta)$, and more specifically, the inequality

$$(15) \quad \|\Delta''\|_1 \leq 3\|\Delta'\|_1 + 4 \sum_{j=r+1}^m \sigma_j(\Theta^*).$$

In the simplest case, when Θ^* is at most rank r , then we have $\sum_{j=r+1}^m \sigma_j(\Theta^*) = 0$, so the constraint (15) defines a cone. This cone completely excludes certain directions, and thus it is possible that the operator \mathfrak{X} , while failing RSC in a global sense, can satisfy it over the cone. Therefore, there is no need for a non-zero tolerance parameter δ in the exact low-rank case. In contrast, when Θ^* is only approximately low-rank, then the constraint (15) no longer defines a cone; rather, it includes an open ball around the origin. Thus, if \mathfrak{X} fails RSC in a global sense, then it will also fail it under the constraint (15). The purpose of the additional constraint $\|\Delta\|_F \geq \delta$ is to eliminate the open ball centered at the origin, so that it is possible that \mathfrak{X} satisfies RSC over $\mathcal{C}(r, \delta)$.

Let us now illustrate the consequences of Theorem 1 when the true matrix Θ^* has exactly rank r , in which case the approximation error term is zero. For the technical reasons mentioned above, it suffices to set $\delta = 0$ in the case of exact rank constraints, and we thus obtain the following result:

COROLLARY 1 (Exact low-rank recovery). *Suppose that $\Theta^* \in \mathcal{S}$ has rank r , and \mathfrak{X} satisfies RSC with respect to $\mathcal{C}(r; 0)$. Then as long as $\lambda_N \geq 2\|\mathfrak{X}^*(\varepsilon)\|_{\text{op}}/N$, any optimal solution $\hat{\Theta}$ to the SDP (9) satisfies the bound*

$$(16) \quad \|\hat{\Theta} - \Theta^*\|_F \leq \frac{32\sqrt{r} \lambda_N}{\kappa(\mathfrak{X})}.$$

Like Theorem 1, Corollary 1 is a deterministic statement on the SDP error. It takes a much simpler form since when Θ^* is exactly low rank, then neither tolerance parameter δ nor the approximation term are required.

As a more delicate example, suppose instead that Θ^* is *nearly low-rank*, an assumption that we can formalize by requiring that its singular value sequence $\{\sigma_i(\Theta^*)\}_{i=1}^m$ decays quickly enough. In particular, for a parameter $q \in [0, 1]$ and a positive radius R_q , we define the set

$$(17) \quad \mathbb{B}_q(R_q) := \left\{ \Theta \in \mathbb{R}^{m_1 \times m_2} \mid \sum_{i=1}^m |\sigma_i(\Theta)|^q \leq R_q \right\},$$

where $m = \min\{m_1, m_2\}$. Note that when $q = 0$, the set $\mathbb{B}_0(R_0)$ corresponds to the set of matrices with rank at most R_0 .

COROLLARY 2 (Near low-rank recovery). *Suppose that $\Theta^* \in \mathbb{B}_q(R_q) \cap \mathcal{S}$, the regularization parameter is lower bounded as $\lambda_N \geq 2\|\mathfrak{X}^*(\varepsilon)\|_{\text{op}}/N$,*

and the operator \mathfrak{X} satisfies RSC with parameter $\kappa(\mathfrak{X}) \in (0, 1]$ over the set $\mathcal{C}(R_q \lambda_N^{-q}; \delta)$. Then any solution $\hat{\Theta}$ to the SDP (9) satisfies

$$(18) \quad \|\hat{\Theta} - \Theta^*\|_F \leq \max \left\{ \delta, 32 \sqrt{R_q} \left(\frac{\lambda_N}{\kappa(\mathfrak{X})} \right)^{1-q/2} \right\}.$$

Note that the error bound (18) reduces to the exact low rank case (16) when $q = 0$, and $\delta = 0$. The quantity $\lambda_N^{-q} R_q$ acts as the “effective rank” in this setting; as clarified by our proof in Section 4.2. This particular choice is designed to provide an optimal trade-off between the approximation and estimation error terms in Theorem 1. Since λ_N is chosen to decay to zero as the sample size N increases, this effective rank will increase, reflecting the fact that as we obtain more samples, we can afford to estimate more of the smaller singular values of the matrix Θ^* .

3.2. Comparison to related work. Past work by Lee and Bresler [34] provides stability results on minimizing the nuclear norm with a quadratic constraint, or equivalently, performing least-squares with nuclear norm constraints. Their results are based on the restricted isometry property (RIP), which is more restrictive than than the RSC condition given here; see Example 4 and Example 5 for concrete examples of operators \mathfrak{X} that satisfy RSC but fail RIP. In our notation, their stability results guarantee that the error $\|\hat{\Theta} - \Theta^*\|_F$ is bounded by a quantity proportional $t := \|y - \mathfrak{X}(\Theta^*)\|_2 / \sqrt{N}$. Given the observation model (5) with a noise vector $\vec{\varepsilon}$ in which each entry is i.i.d., zero mean with variance ν^2 , note that we have $t \approx \nu$ with high probability. Thus, although such a result guarantees stability, it does not guarantee consistency, since for any fixed noise variance $\nu^2 > 0$, the error bound does not tend to zero as the sample size N increases. In contrast, our bounds all depend on the noise and sample size via the regularization parameter, whose optimal choice is $\lambda_N^* = 2 \|\mathfrak{X}^*(\vec{\varepsilon})\|_{\text{op}} / N$. As will be clarified in Corollaries 3 through 5 to follow, for noise $\vec{\varepsilon}$ with variance ν and various choices of \mathfrak{X} , this regularization parameter satisfies the scaling $\lambda_N^* \asymp \nu \sqrt{\frac{m_1 + m_2}{N}}$. Thus, our results guarantee consistency of the estimator, meaning that the error tends to zero as the sample size increases.

As previously noted, some concurrent work [13, 49] has also provided results on estimation of high-dimensional matrices in the noisy and statistical setting. Rohde and Tsybakov [49] derive results for estimating low-rank matrices based on a quadratic loss term regularized by the Schatten- q norm for $0 < q \leq 1$. Note that the the nuclear norm ($q = 1$) is a convex program, whereas the values $q \in (0, 1)$ provide analogs on concave regularized least squares [20] in the linear regression setting. They provide results on both

multivariate regression and matrix completion; most closely related to our work are the results on multivariate regression, which we discuss at more length following Corollary 3 below. On the other hand, Candes and Plan [13] present error rates in the Frobenius norm for estimating approximately low-rank matrices under the compressed sensing model, and we discuss below the connection to our Corollary 5 for this particular observation model. A major difference between our work and this body of work lies in the assumptions imposed on the observation operator \mathfrak{X} . All of the papers [34, 13, 49] impose the restricted isometry property (RIP), which requires that all restricted singular values of \mathfrak{X} very close to 1 (so that it is a near-isometry). In contrast, we require only the restricted strong convexity (RSC) condition, which imposes only an arbitrarily small but positive lower bound on the operator. It is straightforward to construct operators \mathfrak{X} that satisfy RSC while failing RIP, as we discuss in Examples 4 and Example 5 to follow.

3.3. Results for specific model classes. As stated, Corollaries 1 and 2 are fairly abstract in nature. More importantly, it is not immediately clear how the key underlying assumption—namely, the RSC condition—can be verified, since it is specified via subspaces that depend on Θ^* , which is itself the unknown quantity that we are trying to estimate. Nonetheless, we now show how, when specialized to more concrete models, these results yield concrete and readily interpretable results. As will be clear in the proofs of these results, each corollary requires overcoming two main technical obstacles: establishing that the appropriate form of the RSC property holds in a uniform sense (so that a priori knowledge of Θ^* is not required), and specifying an appropriate choice of the regularization parameter λ_N . Each of these two steps is non-trivial, requiring some random matrix theory, but the end results are simply stated upper bounds that hold with high probability.

We begin with the case of rank-constrained multivariate regression. As discussed earlier in Example 1, recall that we observe pairs $(Y_i, Z_i) \in \mathbb{R}^{m_1} \times \mathbb{R}^{m_2}$ linked by the linear model $Y_i = \Theta^* Z_i + W_i$, where $W_i \sim N(0, \nu^2 I_{m_1 \times m_1})$ is observation noise. Here we treat the case of *random design regression*, meaning that the covariates Z_i are modeled as random. In particular, in the following result, we assume that $Z_i \sim N(0, \Sigma)$, i.i.d. for some m_2 -dimensional covariance matrix $\Sigma \succ 0$. Recalling that $\sigma_{\max}(\Sigma)$ and $\sigma_{\min}(\Sigma)$ denote the maximum and minimum eigenvalues respectively, we have:

COROLLARY 3 (Low-rank multivariate regression). *Consider the random design multivariate regression model where $\Theta^* \in \mathbb{B}_q(R_q) \cap \mathcal{S}$. There are universal constants $\{c_i, i = 1, 2, 3\}$ such that if we solve the SDP (9) with*

regularization parameter $\lambda_N = 10 \frac{\nu}{m_1} \sqrt{\sigma_{\max}(\Sigma)} \sqrt{\frac{(m_1+m_2)}{n}}$, we have

$$(19) \quad \|\widehat{\Theta} - \Theta^*\|_F^2 \leq c_1 \left(\frac{\nu^2 \sigma_{\max}(\Sigma)}{\sigma_{\min}^2(\Sigma)} \right)^{1-q/2} R_q \left(\frac{m_1 + m_2}{n} \right)^{1-q/2}$$

with probability greater than $1 - c_2 \exp(-c_3(m_1 + m_2))$.

Remarks: Corollary 3 takes a particularly simple form when $\Sigma = I_{m_2 \times m_2}$: then there exists a constant c'_1 such that $\|\widehat{\Theta} - \Theta^*\|_F^2 \leq c'_1 \nu^{2-q} R_q \left(\frac{m_1+m_2}{n} \right)^{1-q/2}$. When Θ^* is exactly low rank—that is, $q = 0$, and Θ^* has rank $r = R_0$ —this simplifies even further to

$$\|\widehat{\Theta} - \Theta^*\|_F^2 \leq c'_1 \frac{\nu^2 r (m_1 + m_2)}{n}.$$

The scaling in this error bound is easily interpretable: naturally, the squared error is proportional to the noise variance ν^2 , and the quantity $r(m_1 + m_2)$ counts the number of degrees of freedom of a $m_1 \times m_2$ matrix with rank r . Note that if we did not impose any constraints on Θ^* , then since a $m_1 \times m_2$ matrix has a total of $m_1 m_2$ free parameters, we would expect at best¹ to obtain rates of the order $\|\widehat{\Theta} - \Theta^*\|_F^2 = \Omega\left(\frac{\nu^2 m_1 m_2}{n}\right)$. Note that when Θ^* is low rank—in particular, when $r \ll \min\{m_1, m_2\}$ —then the nuclear norm estimator achieves substantially faster rates.²

It is worth comparing this corollary to a result on multivariate regression due to Rohde and Tsybakov [49]. Their result applies to exactly low-rank matrices (say with rank r), but provides bounds on general Schatten norms (including the Frobenius norm). In this case, it provides a comparable rate when we make the setting $q = 0$ and $R_0 = r$ in the bound (19), namely showing that we require roughly $n \approx r(m_1 + m_2)$ samples, corresponding to the number of degrees of freedom. A significant difference lies in the conditions imposed on the design matrices: whereas their result is derived under RIP conditions on the design matrices, we require only the milder RSC condition. The following example illustrates the distinction for this model.

¹To clarify our use of sample size, we can either view the multivariate regression model as consisting of n samples with a constant SNR, or as N samples with SNR of order $1/m_1$. We adopt the former interpretation here.

²We also note that as stated, the result requires that $(m_1 + m_2)$ tend to infinity in order for the claim to hold with high probability. Although such high-dimensional scaling is the primary focus of this paper, we note that for application to the classical setting of fixed (m_1, m_2) , the same statement (with different constants) holds with $m_1 + m_2$ replaced by $\log n$.

EXAMPLE 4 (Failure of RIP for multivariate regression). Under the random design model for multivariate regression, we have

$$(20) \quad F(\Theta) := \frac{\mathbb{E}[\|\mathfrak{X}(\Theta)\|_2^2]}{n\|\Theta\|_F^2} = \frac{\sum_{j=1}^{m_2} \|\sqrt{\Sigma}\Theta_j\|_2^2}{\|\Theta\|_F^2},$$

where Θ_j is the j^{th} row of Θ . In order for RIP to hold, it is necessary that quantity $F(\Theta)$ is extremely close to 1—certainly less than two—for all low-rank matrices. We now show that this cannot hold unless Σ has a small condition number. Let $v \in \mathbb{R}^{m_2}$ and $v' \in \mathbb{R}^{m_2}$ denote the minimum and maximum eigenvectors of Σ . By setting $\Theta = e_1 v^T$, we obtain a rank one matrix for which $F(\Theta) = \sigma_{\min}(\Sigma)$, and similarly, setting $\Theta' = e_1 (v')^T$ yields another rank one matrix for which $F(\Theta') = \sigma_{\max}(\Sigma)$. The preceding discussion applies to the average $\mathbb{E}[\|\mathfrak{X}(\Theta)\|_2^2]/n$, but since the individual matrices X_i are i.i.d. and Gaussian, we have

$$\frac{\|\mathfrak{X}(\Theta)\|_2^2}{n} = \frac{1}{n} \sum_{i=1}^n \langle X_i, \Theta \rangle^2 \leq 2F(\Theta) = 2\sigma_{\min}(\Sigma)$$

with high probability, using χ^2 -tail bounds. Similarly, $\|\mathfrak{X}(\Theta')\|_2^2/n \geq (1/2)\sigma_{\max}(\Sigma)$ with high probability. Thus, we have exhibited a pair of rank one matrices with $\|\Theta\|_F = \|\Theta'\|_F = 1$ for which

$$\frac{\|\mathfrak{X}(\Theta')\|_2^2}{\|\mathfrak{X}(\Theta)\|_2^2} \geq \frac{1}{4} \frac{\sigma_{\max}(\Sigma)}{\sigma_{\min}(\Sigma)}.$$

Consequently, unless $\sigma_{\max}(\Sigma)/\sigma_{\min}(\Sigma) \leq 64$, it is not possible for RIP to hold with constant $\delta \leq 1/2$. In contrast, as our results show, the RSC will hold w.h.p. whenever $\sigma_{\min}(\Sigma) > 0$, and the error is allowed to scale with the ratio $\sigma_{\max}(\Sigma)/\sigma_{\min}(\Sigma)$.

Next we turn to the case of estimating the system matrix Θ^* of an autoregressive (AR) model, as discussed in Example 2.

COROLLARY 4 (Autoregressive models). *Suppose that we are given n samples $\{Z_t\}_{t=1}^n$ from a m -dimensional autoregressive process (2) that is stationary, based on a system matrix that is stable ($\|\Theta^*\|_{\text{op}} \leq \gamma < 1$), and approximately low-rank ($\Theta^* \in \mathbb{B}_q(R_q) \cap \mathcal{S}$). Then there are universal constants $\{c_i, i = 1, 2, 3\}$ such that if we solve the SDP (9) with regularization parameter $\lambda_N = \frac{2c_0 \|\Sigma\|_{\text{op}}}{m(1-\gamma)} \sqrt{\frac{m}{n}}$, then any solution $\hat{\Theta}$ satisfies*

$$(21) \quad \|\hat{\Theta} - \Theta^*\|_F^2 \leq c_1 \left[\frac{\sigma_{\max}^2(\Sigma)}{\sigma_{\min}^2(\Sigma)(1-\gamma)^2} \right]^{1-q/2} R_q \left(\frac{m}{n} \right)^{1-q/2}$$

with probability greater than $1 - c_2 \exp(-c_3 m)$.

Remarks: Like Corollary 3, the result as stated requires that the matrix dimension m tends to infinity, but the same bounds hold with m replaced by $\log n$, yielding results suitable for classical (fixed dimension) scaling. Second, the factor $(m/n)^{1-q/2}$, like the analogous term³ in Corollary 3, shows that faster rates are obtained if Θ^* can be well-approximated by a low rank matrix, namely for choices of the parameter $q \in [0, 1]$ that are closer to zero. Indeed, in the limit $q = 0$, we again reduce to the case of an exact rank constraint $r = R_0$, and the corresponding squared error scales as rm/n . In contrast to the case of multivariate regression, the error bound (21) also depends on the upper bound $\|\Theta^*\|_{\text{op}} = \gamma < 1$ on the operator norm of the system matrix Θ^* . Such dependence is to be expected since the quantity γ controls the (in)stability and mixing rate of the autoregressive process. As clarified in the proof, the dependence of the sampling in the AR model also presents some technical challenges not present in the setting of multivariate regression.

Finally, we turn to the analysis of the compressed sensing model for matrix recovery, as initially described in Example 3. Although standard compressed sensing is based on observation matrices X_i with i.i.d. elements, here we consider a more general model that allows for dependence between the entries of X_i . First defining the quantity $M = m_1 m_2$, we use $\text{vec}(X_i) \in \mathbb{R}^M$ to denote the vectorized form of the $m_1 \times m_2$ matrix X_i . Given a symmetric positive definite matrix $\Sigma \in \mathbb{R}^{M \times M}$, we say that the observation matrix X_i is sampled from the Σ -ensemble if $\text{vec}(X_i) \sim N(0, \Sigma)$. Finally, we define the quantity

$$(22) \quad \rho^2(\Sigma) := \sup_{\|u\|_2=1, \|v\|_2=1} \text{var}(u^T X v),$$

where the random matrix $X \in \mathbb{R}^{m_1 \times m_2}$ is sampled from the Σ -ensemble. In the special case $\Sigma = I$, corresponding to the usual compressed sensing model, we have $\rho^2(I) = 1$.

The following result applies to any observation model in which the noise vector $\tilde{\varepsilon} \in \mathbb{R}^N$ satisfies the bound $\|\tilde{\varepsilon}\|_2 \leq 2\nu\sqrt{N}$ for some constant ν . This assumption that holds for any bounded noise, and also holds with high probability for any random noise vector with sub-Gaussian entries with parameter ν . (The simplest example is that of Gaussian noise $N(0, \nu^2)$.)

³The term in Corollary 3 has a factor $m_1 + m_2$, since the matrix in that case could be non-square in general.

COROLLARY 5 (Compressed sensing with dependent sampling). *Suppose that the matrices $\{X_i\}_{i=1}^N$ are drawn i.i.d. from the Σ -Gaussian ensemble, and that the unknown matrix $\Theta^* \in \mathbb{B}_q(R_q) \cap \mathcal{S}$ for some $q \in (0, 1]$. Then there are universal constants c_i such that for a sample size $N > c_1 \rho^2(\Sigma) R_q^{1-q/2} (m_1 + m_2)$, any solution $\hat{\Theta}$ to the SDP (9) with regularization parameter $\lambda_N = c_0 \rho(\Sigma) \nu \sqrt{\frac{m_1 + m_2}{N}}$ satisfies the bound*

$$(23) \quad \|\hat{\Theta} - \Theta^*\|_F^2 \leq c_2 R_q \left(\frac{(\nu^2 \vee 1) \frac{\rho^2(\Sigma)}{\sigma_{\min}^2(\Sigma)} (m_1 + m_2)}{N} \right)^{1-\frac{q}{2}}$$

with probability greater than $1 - c_3 \exp(-c_4(m_1 + m_2))$. In the special case $q = 0$ and Θ^* of rank r , we have

$$(24) \quad \|\hat{\Theta} - \Theta^*\|_F^2 \leq c_2 \frac{\rho^2(\Sigma) \nu^2}{\sigma_{\min}^2(\Sigma)} \frac{r (m_1 + m_2)}{N}.$$

The central challenge in proving this result is in proving an appropriate form of the RSC property. The following result on the random operator \mathfrak{X} may be of independent interest here:

PROPOSITION 1. *Consider the random operator $\mathfrak{X} : \mathbb{R}^{m_1 \times m_2} \rightarrow \mathbb{R}^N$ formed by sampling from the Σ -ensemble. Then it satisfies*

$$(25) \quad \frac{\|\mathfrak{X}(\Theta)\|_2}{\sqrt{N}} \geq \frac{1}{4} \|\sqrt{\Sigma} \text{vec}(\Theta)\|_2 - 12\rho(\Sigma) \left(\sqrt{\frac{m_1}{N}} + \sqrt{\frac{m_2}{N}} \right) \|\Theta\|_1 \quad \text{for all } \Theta \in \mathbb{R}^{m_1 \times m_2}$$

with probability at least $1 - 2 \exp(-N/32)$.

The proof of this result, provided in Appendix E, makes use of the Gordon-Slepian inequalities for Gaussian processes, and concentration of measure. As we show in Section 4.5, it implies the form of the RSC property needed to establish Corollary 5.

In concurrent work, Candes and Plan [13] derived a result similar to Corollary 5 for the compressed sensing observation model. Their result applies to matrices with i.i.d. elements with sub-Gaussian tail behavior. While the analysis given here is specific to Gaussian random matrices, it allows for general dependence among the entries. Their result applies only under certain restrictions on the sample size relative to matrix dimension and rank, whereas our result holds more generally without these extra conditions. Moreover, their proof relies on an application of RIP, which is in general

more restrictive than the RSC condition used in our analysis. The following example provides a concrete illustration of a matrix family where the restricted isometry constants are unbounded as the rank r grows, but RSC still holds.

EXAMPLE 5 (RSC holds when RIP violated). Here we consider a family of random operators \mathfrak{X} for which RSC holds with high probability, while RIP fails. Consider generating an i.i.d. collection of design matrices $X_i \in \mathbb{R}^{m \times m}$, each of the form

$$(26) \quad X_i = z_i I_{m \times m} + G_i, \quad \text{for } i = 1, 2, \dots, N,$$

where $z_i \sim N(0, 1)$ and $G_i \in \mathbb{R}^{m \times m}$ is a standard Gaussian random matrix, independent of z_i . Note that we have $\text{vec}(X_i) \sim N(0, \Sigma)$, where the $m^2 \times m^2$ covariance matrix has the form

$$(27) \quad \Sigma = \text{vec}(I_{m \times m}) \text{vec}(I_{m \times m})^T + I_{m^2 \times m^2}.$$

Let us compute the quantity $\rho(\Sigma) = \sup_{\substack{\|u\|_2=1 \\ \|v\|_2=1}} \text{var}(u^T X v)$. By the independence of z and G in the model (26), we have

$$\rho(\Sigma) \leq \text{var}(z) \sup_{u \in S^{m_1-1}, v \in S^{m_2-1}} u^T v + \sup_{u \in S^{m_1-1}, v \in S^{m_2-1}} \text{var}(u^T G v) \leq 2.$$

Letting \mathfrak{X} be the associated random operator, we observe that for any $\Theta \in \mathbb{R}^{m \times m}$, the independence of z_i and G_i implies that

$$\mathbb{E} \left[\frac{\|\mathfrak{X}(\Theta)\|_2^2}{N} \right] = \|\sqrt{\Sigma} \text{vec}(\Theta)\|_2^2 = \text{trace}(\Theta)^2 + \|\Theta\|_F^2 \geq \|\Theta\|_F^2.$$

Consequently, Proposition 1 implies that

$$(28) \quad \frac{\|\mathfrak{X}(\Theta)\|_2}{\sqrt{N}} \geq \frac{1}{4} \|\Theta\|_F - 48 \sqrt{\frac{m}{N}} \|\Theta\|_1 \quad \text{for all } \Theta \in \mathbb{R}^{m \times m},$$

with high probability. As mentioned previously, we show in Section 4.5 how this type of lower bound implies the RSC condition needed for our results.

On the other hand, the random operator can never satisfy RIP (with the rank r increasing), as the following calculation shows. In this context, RIP requires that bounds of the form

$$\frac{\|\mathfrak{X}(\Theta)\|_2}{N \|\Theta\|_F^2} \in [1 - \delta, 1 + \delta] \quad \text{for all } \Theta \text{ with rank at most } r,$$

where $\delta \in (0, 1)$ is a constant independent of r . Note that the bound (28) implies that a *lower bound* of this form holds as long as $N = \Omega(rm)$. Moreover, this lower bound cannot be substantially sharpened, since the trace term plays no role for matrices with zero diagonals.

We now show that no such upper bound can ever hold. For a rank $1 \leq r < m$, consider the $m \times m$ matrix of the form

$$\Gamma := \begin{bmatrix} I_{r \times r} / \sqrt{r} & 0_{r \times (m-r)} \\ 0_{(m-r) \times r} & 0_{(m-r) \times (m-r)} \end{bmatrix}.$$

By construction, we have $\|\Gamma\|_F = 1$ and $\text{trace}(\Gamma) = \sqrt{r}$. Consequently, we have

$$\mathbb{E} \left[\frac{\|\mathfrak{X}(\Gamma)\|_2^2}{N} \right] = \text{trace}(\Gamma)^2 + \|\Gamma\|_F^2 = r + 1.$$

The independence of the matrices X_i implies that $\frac{\|\mathfrak{X}(\Gamma)\|_2^2}{N}$ is sharply concentrated around this expected value, so that we conclude that

$$\frac{\|\mathfrak{X}(\Gamma)\|_2^2}{N \|\Gamma\|_F^2} \geq \frac{1}{2} [1 + r],$$

with high probability, showing that RIP cannot hold with upper and lower bounds of the same order.

Finally, we note that Proposition 1 also implies an interesting property of the null space of the operator \mathfrak{X} , one which can be used to establish a corollary about recovery of low-rank matrices when the observations are noiseless. In particular, suppose that we are given the noiseless observations $y_i = \langle X_i, \Theta^* \rangle$ for $i = 1, \dots, N$, and that we try to recover the unknown matrix Θ^* by solving the SDP

$$(29) \quad \min_{\Theta \in \mathbb{R}^{m_1 \times m_2}} \|\Theta\|_1 \quad \text{such that } \langle X_i, \Theta \rangle = y_i \text{ for all } i = 1, \dots, N.$$

This recovery procedure was studied by Recht and colleagues [47, 48] in the special case that X_i is formed of i.i.d. $N(0, 1)$ entries. Proposition 1 allows us to obtain a sharp result on recovery using this method for Gaussian matrices with general dependencies.

COROLLARY 6 (Exact recovery with dependent sampling). *Suppose that the matrices $\{X_i\}_{i=1}^N$ are drawn i.i.d. from the Σ -Gaussian ensemble, and*

that $\Theta^* \in \mathcal{S}$ has rank r . Given $N > c_0 \rho^2(\Sigma) r(m_1 + m_2)$ noiseless samples, then with probability at least $1 - 2 \exp(-N/32)$, the SDP (29) recovers the matrix Θ^* exactly.

This result removes some extra logarithmic factors that were included in initial work [47] and provides the appropriate analog to compressed sensing results for sparse vectors [18, 14]. Note that (like in most of our results) we have made little effort to obtain good constants in this result: the important property is that the sample size N scales linearly in both r and $m_1 + m_2$. We refer the reader to Recht et al. [48], who study the standard Gaussian model under the scaling $r = \Theta(m)$ and obtain sharp results on the constants.

4. Proofs. We now turn to the proofs of Theorem 1, and Corollaries 1 through 6. In each case, we provide the primary steps in the main text, with more technical details stated as lemmas and proved in the Appendix.

4.1. *Proof of Theorem 1.* By the optimality of $\widehat{\Theta}$ and feasibility of Θ^* for the SDP (9), we have

$$\frac{1}{2N} \|\vec{y} - \mathfrak{X}(\widehat{\Theta})\|_2^2 + \lambda_N \|\widehat{\Theta}\|_1 \leq \frac{1}{2N} \|\vec{y} - \mathfrak{X}(\Theta^*)\|_2^2 + \lambda_N \|\Theta^*\|_1.$$

Defining the error matrix $\Delta = \Theta^* - \widehat{\Theta}$ and performing some algebra yields the inequality

$$(30) \quad \frac{1}{2N} \|\mathfrak{X}(\Delta)\|_2^2 \leq \frac{1}{N} \langle \vec{\varepsilon}, \mathfrak{X}(\Delta) \rangle + \lambda_N \{ \|\widehat{\Theta} + \Delta\|_1 - \|\widehat{\Theta}\|_1 \}.$$

By definition of the adjoint and Hölder's inequality, we have

$$(31) \quad \frac{1}{N} |\langle \vec{\varepsilon}, \mathfrak{X}(\Delta) \rangle| = \frac{1}{N} |\langle \mathfrak{X}^*(\vec{\varepsilon}), \Delta \rangle| \leq \frac{1}{N} \|\mathfrak{X}^*(\vec{\varepsilon})\|_{\text{op}} \|\Delta\|_1.$$

By the triangle inequality, we have $\|\widehat{\Theta} + \Delta\|_1 - \|\widehat{\Theta}\|_1 \leq \|\Delta\|_1$. Substituting this inequality and the bound (31) into the inequality (30) yields

$$\frac{1}{2N} \|\mathfrak{X}(\Delta)\|_2^2 \leq \frac{1}{N} \|\mathfrak{X}^*(\vec{\varepsilon})\|_{\text{op}} \|\Delta\|_1 + \lambda_N \|\Delta\|_1 \leq 2\lambda_N \|\Delta\|_1,$$

where the second inequality makes use of our choice $\lambda_N \geq \frac{2}{N} \|\mathfrak{X}^*(\vec{\varepsilon})\|_{\text{op}}$.

It remains to lower bound the term on the left-hand side, while upper bounding the quantity $\|\Delta\|_1$ on the right-hand side. The following technical result allows us to do so. Recall our earlier definition (11) of the sets \mathcal{A} and \mathcal{B} associated with a given subspace pair.

LEMMA 1. *Let $U^r \in \mathbb{R}^{m_1 \times r}$ and $V^r \in \mathbb{R}^{m_2 \times r}$ be matrices consisting of the top r left and right (respectively) singular vectors of Θ^* . Then there exists a matrix decomposition $\Delta = \Delta' + \Delta''$ of the error Δ such that*

- (a) *The matrix Δ' satisfies the constraint $\text{rank}(\Delta') \leq 2r$, and*
- (b) *If $\lambda_N \geq 2\|\mathfrak{X}^*(\tilde{\varepsilon})\|_{\text{op}}/N$, then the nuclear norm of Δ'' is bounded as*

$$(32) \quad \|\Delta''\|_1 \leq 3\|\Delta'\|_1 + 4 \sum_{j=r+1}^m \sigma_j(\Theta^*)$$

See Appendix A for the proof of this claim. Using Lemma 1, we can complete the proof of the theorem. In particular, from the bound (32) and the RSC assumption, we find that for $\|\Delta\|_F \geq \delta$, we have

$$\frac{1}{2N} \|\mathfrak{X}(\Delta)\|_2^2 \geq \kappa(\mathfrak{X}) \|\Delta\|_F^2.$$

Using the triangle inequality together with inequality (32), we obtain

$$\|\Delta\|_1 \leq \|\Delta'\|_1 + \|\Delta''\|_1 \leq 4\|\Delta'\|_1 + 4 \sum_{j=r+1}^m \sigma_j(\Theta^*).$$

From the rank constraint in Lemma 1(a), we have $\|\Delta'\|_1 \leq \sqrt{2r}\|\Delta'\|_F$. Putting together the pieces, we find that either $\|\Delta\|_F \leq \delta$, or

$$\kappa(\mathfrak{X}) \|\Delta\|_F^2 \leq \max \left\{ 32\lambda_N \sqrt{r} \|\Delta\|_F, 16 \lambda_N \sum_{j=r+1}^m \sigma_j(\Theta^*) \right\},$$

which implies that

$$\|\Delta\|_F \leq \max \left\{ \delta, \frac{32\lambda_N \sqrt{r}}{\kappa(\mathfrak{X})}, \left(\frac{16 \lambda_N \sum_{j=r+1}^m \sigma_j(\Theta^*)}{\kappa(\mathfrak{X})} \right)^{1/2} \right\},$$

as claimed.

4.2. *Proof of Corollary 2.* Let $m = \min\{m_1, m_2\}$. In this case, we consider the singular value decomposition $\Theta^* = UDV^T$, where $U \in \mathbb{R}^{m_1 \times m}$ and $V \in \mathbb{R}^{m_2 \times m}$ are orthogonal, and we assume that D is diagonal with the singular values in non-increasing order $\sigma_1(\Theta^*) \geq \sigma_2(\Theta^*) \geq \dots \geq \sigma_m(\Theta^*) \geq 0$. For a parameter $\tau > 0$ to be chosen, we define

$$K := \{i \in \{1, 2, \dots, m\} \mid \sigma_i(\Theta^*) > \tau\},$$

and we let U^K (respectively V^K) denote the $m_1 \times |K|$ (respectively the $m_2 \times |K|$) orthogonal matrix consisting of the first $|K|$ columns of U (respectively V). With this choice, the matrix $\Theta_{K^c}^* := \Pi_{\mathcal{B}^{|K|}}(\Theta^*)$ has rank at most $m - |K|$, with singular values $\{\sigma_i(\Theta^*), i \in K^c\}$. Moreover, since $\sigma_i(\Theta^*) \leq \tau$ for all $i \in K^c$, we have

$$\|\Theta_{K^c}^*\|_1 = \tau \sum_{i=|K|+1}^m \frac{\sigma_i(\Theta^*)}{\tau} \leq \tau \sum_{i=|K|+1}^m \left(\frac{\sigma_i(\Theta^*)}{\tau}\right)^q \leq \tau^{1-q} R_q.$$

On the other hand, we also have $R_q \geq \sum_{i=1}^m |\sigma_i(\Theta^*)|^q \geq |K| \tau^q$, which implies that $|K| \leq \tau^{-q} R_q$. From the general error bound with $r = |K|$, we obtain

$$\|\hat{\Theta} - \Theta^*\|_F \leq \max\left\{ \delta, \frac{32\lambda_N \sqrt{R_q} \tau^{-q/2}}{\kappa(\mathcal{X})}, \left[\frac{16 \lambda_N \tau^{1-q} R_q}{\kappa(\mathcal{X})} \right]^{1/2} \right\},$$

Setting $\tau = \lambda_N / \kappa$ yields that

$$\begin{aligned} \|\hat{\Theta} - \Theta^*\|_F &\leq \max\left\{ \delta, \frac{32\lambda_N^{1-q/2} \sqrt{R_q}}{\kappa^{1-q/2}}, \left[\frac{16 \lambda_N^{2-q} R_q}{\kappa^{2-q}} \right]^{1/2} \right\} \\ &= \max\left\{ \delta, 32 \sqrt{R_q} \left(\frac{\lambda_N}{\kappa(\mathcal{X})} \right)^{1-q/2} \right\}, \end{aligned}$$

as claimed.

4.3. *Proof of Corollary 3.* For the proof of this corollary, we adopt the following notation. We first define the three matrices

$$(33) \quad X = \begin{bmatrix} Z_1^T \\ Z_2^T \\ \dots \\ Z_n^T \end{bmatrix} \in \mathbb{R}^{n \times m_2}, \quad Y = \begin{bmatrix} Y_1^T \\ Y_2^T \\ \dots \\ Y_n^T \end{bmatrix} \in \mathbb{R}^{n \times m_1}, \quad \text{and} \quad W = \begin{bmatrix} W_1^T \\ W_2^T \\ \dots \\ W_n^T \end{bmatrix} \in \mathbb{R}^{n \times m_1}.$$

With this notation and using the relation $N = nm_1$, the SDP objective function (9) can be written as $\frac{1}{m_1} \left\{ \frac{1}{2n} \|Y - X\Theta^T\|_F^2 + \lambda_n \|\Theta\|_1 \right\}$, where we have defined $\lambda_n = \lambda_N m_1$.

In order to establish the RSC property for this model, some algebra shows that we need to establish a lower bound on the quantity

$$\frac{1}{2n} \|X\Delta\|_F^2 = \frac{1}{2n} \sum_{j=1}^m \|(X\Delta)_j\|_2^2 \geq \frac{\sigma_{\min}(X^T X)}{2n} \|\Delta\|_F^2,$$

where σ_{\min} denotes the minimum eigenvalue. The following lemma follows by adapting known concentration results for random matrices (see the paper [54] for details):

LEMMA 2. *Let $X \in \mathbb{R}^{n \times m}$ be a random matrix with i.i.d. rows sampled from a m -variate $N(0, \Sigma)$ distribution. Then for $n \geq 2m$, we have*

$$\mathbb{P} \left[\sigma_{\min} \left(\frac{1}{n} X^T X \right) \geq \frac{\sigma_{\min}(\Sigma)}{9}, \sigma_{\max} \left(\frac{1}{n} X^T X \right) \leq 9\sigma_{\max}(\Sigma) \right] \geq 1 - 4\exp(-n/2).$$

As a consequence, we have $\frac{\sigma_{\min}(X^T X)}{2n} \geq \frac{\sigma_{\min}(\Sigma)}{18}$ with probability at least $1 - 4\exp(-n)$ for all $n \geq 2m$, which establishes that the RSC property holds with $\kappa(\mathfrak{X}) = \frac{1}{20m_1}\sigma_{\min}(\Sigma)$.

Next we need to upper bound the quantity $\|\mathfrak{X}^*(\vec{\varepsilon})\|_{\text{op}}$ for this model, so as to verify that the stated choice for λ_N is valid. Following some algebra, we find that

$$\frac{1}{n} \|\mathfrak{X}^*(\vec{\varepsilon})\|_{\text{op}} = \frac{1}{n} \|X^T W\|_{\text{op}}.$$

The following lemma is proved in Appendix C:

LEMMA 3. *There are constants $c_i > 0$ such that*

(34)

$$\mathbb{P} \left[\left| \frac{1}{n} \|X^T W\|_{\text{op}} \right| \geq 5\nu \sqrt{\sigma_{\max}(\Sigma)} \sqrt{\frac{m_1 + m_2}{n}} \right] \leq c_1 \exp(-c_2(m_1 + m_2)).$$

Using these two lemmas, we can complete the proof of Corollary 3. First, recalling the scaling $N = m_1 n$, we see that Lemma 3 implies that the choice $\lambda_n = 10\nu \sqrt{\sigma_{\max}(\Sigma)} \sqrt{\frac{m_1 + m_2}{n}}$ satisfies the conditions of Corollary 2 with high probability. Lemma 2 shows that the RSC property holds with $\kappa(\mathfrak{X}) = \sigma_{\min}(\Sigma)/(20m_1)$, again with high probability. Consequently, Corollary 2 implies that

$$\begin{aligned} \|\hat{\Theta} - \Theta^*\|_F^2 &\leq 32^2 R_q \left(10\nu \sqrt{\sigma_{\max}(\Sigma)} \sqrt{\frac{m_1 + m_2}{n}} \frac{20}{\sigma_{\min}(\Sigma)} \right)^{2-q} \\ &= c_1 \left(\frac{\nu^2 \sigma_{\max}(\Sigma)}{\sigma_{\min}^2(\Sigma)} \right)^{1-q/2} R_q \left(\frac{m_1 + m_2}{n} \right)^{1-q/2} \end{aligned}$$

with probability greater than $1 - c_2 \exp(-c_3(m_1 + m_2))$, as claimed.

4.4. *Proof of Corollary 4.* For the proof of this corollary, we adopt the notation

$$X = \begin{bmatrix} Z_1^T \\ Z_2^T \\ \dots \\ Z_n^T \end{bmatrix} \in \mathbb{R}^{n \times m}, \quad \text{and} \quad Y = \begin{bmatrix} Z_2^T \\ Z_2^T \\ \dots \\ Z_{n+1}^T \end{bmatrix} \in \mathbb{R}^{n \times m}.$$

Finally, we let $W \in \mathbb{R}^{n \times m}$ be a matrix where each row is sampled i.i.d. from the $N(0, C)$ distribution corresponding to the innovations noise driving the VAR process. With this notation and using the relation $N = nm$, the SDP objective function (9) can be written as $\frac{1}{m} \{ \frac{1}{2n} \|Y - X\Theta^T\|_F^2 + \lambda_n \|\Theta\|_1 \}$, where we have defined $\lambda_n = \lambda_N m$. At a high level, the proof of this corollary is similar to that of Corollary 3, in that we use random matrix theory to establish the required RSC property, and to justify the choice of λ_n , or equivalently λ_N . However, it is considerably more challenging, due to the dependence in the rows of the random matrices, and the cross-dependence between the two matrices X and W (which were independent in the setting of multivariate regression).

The following lemma provides the lower bound needed to establish RSC for the autoregressive model:

LEMMA 4. *The eigenspectrum of the matrix $X^T X/n$ is well-controlled in terms of the stationary covariance matrix: in particular, as long as $n > c_3 m$, we have*

$$(35) \quad \sigma_{\max} \left(\left(\frac{1}{n} X^T X \right) \right) \stackrel{(a)}{\leq} \frac{24 \sigma_{\max}(\Sigma)}{1 - \gamma}, \quad \text{and} \quad \sigma_{\min} \left(\left(\frac{1}{n} X^T X \right) \right) \stackrel{(b)}{\geq} \frac{\sigma_{\min}(\Sigma)}{4},$$

both with probability greater than $1 - 2c_1 \exp(-c_2 m)$.

Thus, from the bound (35)(b), we see with the high probability, the RSC property holds with $\kappa(\mathfrak{X}) = \sigma_{\min}(\Sigma)/(4m_2)$ as long as $n > c_3 m$.

As before, in order to verify the choice of λ_n , we need to control the quantity $\frac{1}{n} \|X^T W\|_{\text{op}}$. The following inequality, proved in Appendix D.2, yields a suitable upper bound:

LEMMA 5. *There exist constants $c_i > 0$, independent of n, m, Σ etc. such that*

$$(36) \quad \mathbb{P} \left[\frac{1}{n} \|X^T W\|_{\text{op}} \geq \frac{c_0 \|\Sigma\|_{\text{op}}}{1 - \gamma} \sqrt{\frac{m}{n}} \right] \leq c_2 \exp(-c_3 m).$$

From Lemma 5, we see that it suffices to choose $\lambda_n = \frac{2c_0 \|\Sigma\|_{\text{op}}}{1-\gamma} \sqrt{\frac{m}{n}}$. With this choice, Corollary 2 of Theorem 1 yields that

$$\|\Theta - \Theta^*\|_F^2 \leq c_1 R_q \left[\frac{\sigma_{\max}(\Sigma)}{\sigma_{\min}(\Sigma)(1-\gamma)} \right]^{2-q} \left(\frac{m}{n} \right)^{1-q/2}$$

with probability greater than $1 - c_2 \exp(-c_3 m)$, as claimed.

4.5. *Proof of Corollary 5.* Recall that for this model, the observations are of the form $y_i = \langle X_i, \Theta^* \rangle + \varepsilon_i$, where $\Theta^* \in \mathbb{R}^{m_1 \times m_2}$ is the unknown matrix, and $\{\varepsilon_i\}_{i=1}^N$ is an associated noise sequence.

We now show how Proposition 1 implies the RSC property with an appropriate tolerance parameter $\delta > 0$ to be defined. Observe that the bound (25) implies that for any $\Delta \in \mathcal{C}$, we have

$$\begin{aligned} \frac{\|\mathfrak{X}(\Delta)\|_2}{\sqrt{N}} &\geq \frac{\sqrt{\sigma_{\min}(\Sigma)}}{4} \|\Delta\|_F - 12\rho(\Sigma) \left(\sqrt{\frac{m_1}{N}} + \sqrt{\frac{m_2}{N}} \right) \|\Delta\|_1 \\ (37) \quad &= \frac{\sqrt{\sigma_{\min}(\Sigma)}}{4} \left\{ \|\Delta\|_F - \underbrace{\frac{48\rho(\Sigma)}{\sqrt{\sigma_{\min}(\Sigma)}} \left(\sqrt{\frac{m_1}{N}} + \sqrt{\frac{m_2}{N}} \right) \|\Delta\|_1}_{\tau} \right\}, \end{aligned}$$

where we have defined the quantity $\tau > 0$. Following the arguments used in the proofs of Theorem 1 and Corollary 2, we find that

$$(38) \quad \|\Delta\|_1 \leq 4\|\Delta'\|_1 + 4 \sum_{j=r+1}^m \sigma_j(\Theta^*) \leq 4\sqrt{2R_q\tau^{-q}} \|\Delta'\|_F + 4R_q\tau^{1-q}.$$

Note that this corresponds to truncating the matrices at effective rank $r = 2R_q\tau^{-q}$. Combining this bound with the definition of τ , we obtain

$$\tau \|\Delta\|_1 \leq 4\sqrt{2R_q\tau^{1-q/2}} \|\Delta'\|_F + 4R_q\tau^{2-q} \leq 4\sqrt{2R_q\tau^{1-q/2}} \|\Delta\|_F + 4R_q\tau^{2-q}.$$

Substituting this bound into equation (37) yields

$$\frac{\|\mathfrak{X}(\Delta)\|_2}{\sqrt{N}} \geq \frac{\sqrt{\sigma_{\min}(\Sigma)}}{4} \left\{ \|\Delta\|_F - 4\sqrt{2R_q\tau^{1-q/2}} \|\Delta'\|_F - 4R_q\tau^{2-q} \right\}.$$

As long $N > c_0 R_q^{2/(2-q)} \frac{\rho^2(\Sigma)}{\sigma_{\min}(\Sigma)} (m_1 + m_2)$ for a sufficiently large constant c_0 , we can ensure that $4\sqrt{2R_q\tau^{1-q/2}} < 1/2$, and hence that

$$\frac{\|\mathfrak{X}(\Delta)\|_2}{\sqrt{N}} \geq \frac{\sqrt{\sigma_{\min}(\Sigma)}}{4} \left\{ \frac{1}{2} \|\Delta\|_F - 4R_q\tau^{2-q} \right\}.$$

Consequently, if we define $\delta := 16 R_q \tau^{2-q}$, then we are guaranteed that for all $\|\Delta\|_F \geq \delta$, we have $4 R_q \tau^{2-q} \leq \|\Delta\|_F/4$, and hence

$$\frac{\|\mathfrak{X}(\Delta)\|_2}{\sqrt{N}} \geq \frac{\sqrt{\sigma_{\min}(\Sigma)}}{16} \|\Delta\|_F$$

for all $\|\Delta\|_F \geq \delta$. We have thus shown that $\mathcal{C}(2 R_q \tau^{-q}; \delta)$ with parameter $\kappa(\mathfrak{X}) = \frac{\sigma_{\min}(\Sigma)}{256}$.

The next step is to control the quantity $\|\mathfrak{X}^*(\vec{\varepsilon})\|_{\text{op}}/N$, required for specifying a suitable choice of λ_N .

LEMMA 6. *If $\|\vec{\varepsilon}\|_2 \leq 2\nu\sqrt{N}$, then there are universal constants c_i such that*

$$(39) \quad \mathbb{P}\left[\frac{\|\mathfrak{X}^*(\vec{\varepsilon})\|_{\text{op}}}{N} \geq c_0 \rho(\Sigma) \nu \left(\sqrt{\frac{m_1}{N}} + \sqrt{\frac{m_2}{N}}\right)\right] \leq c_1 \exp(-c_2(m_1 + m_2)).$$

PROOF. By the definition of the adjoint operator, we have $Z = \frac{1}{N} \mathfrak{X}^*(\vec{\varepsilon}) = \frac{1}{N} \sum_{i=1}^N \varepsilon_i X_i$. Since the observation matrices $\{X_i\}_{i=1}^N$ are i.i.d. Gaussian, if the sequence $\{\varepsilon_i\}_{i=1}^N$ is viewed as fixed (by conditioning as needed), then the random matrix Z is a sample from the Γ -ensemble with covariance matrix $\Gamma = \frac{\|\vec{\varepsilon}\|_2^2}{N^2} \Sigma \preceq \frac{2\nu^2}{N} \Sigma$. Therefore, letting $\tilde{Z} \in \mathbb{R}^{m_1 \times m_2}$ be a random matrix drawn from the $2\nu^2 \Sigma/N$ -ensemble, we have

$$\mathbb{P}[\|Z\|_{\text{op}} \geq t] \leq \mathbb{P}[\|\tilde{Z}\|_{\text{op}} \geq t].$$

Using Lemma 7 from Appendix E, we have

$$\mathbb{E}[\|\tilde{Z}\|_{\text{op}}] \leq \frac{12\sqrt{2}\nu\rho(\Sigma)}{\sqrt{N}} (\sqrt{m_1} + \sqrt{m_2})$$

and

$$\mathbb{P}[\|\tilde{Z}\|_{\text{op}} \geq \mathbb{E}[\|\tilde{Z}\|_{\text{op}}] + t] \leq \exp\left(-c_1 \frac{Nt^2}{\nu^2 \rho^2(\Sigma)}\right)$$

for a universal constant c_1 . Setting $t^2 = \Omega\left(\frac{\nu^2 \rho^2(\Sigma)(\sqrt{m_1} + \sqrt{m_2})^2}{N}\right)$ yields the claim. □

4.6. *Proof of Corollary 6.* This corollary follows from a combination of Proposition 1 and Lemma 1. Let $\widehat{\Theta}$ be an optimal solution to the SDP (29), and let $\Delta = \widehat{\Theta} - \Theta^*$ be the error. Since $\widehat{\Theta}$ is optimal and Θ^* is feasible for the SDP, we have $\|\widehat{\Theta}\|_1 = \|\Theta^* + \Delta\|_1 \leq \|\Theta^*\|_1$. Using the decomposition $\Delta = \Delta' + \Delta''$ from Lemma 1 and applying triangle inequality, we have

$$\|\Theta^* + \Delta' + \Delta''\|_1 \geq \|\Theta^* + \Delta''\|_1 - \|\Delta'\|_1.$$

From the properties of the decomposition in Lemma 1 (see Appendix A), we find that

$$\|\widehat{\Theta}\|_1 = \|\Theta^* + \Delta' + \Delta''\|_1 \geq \|\Theta^*\|_1 + \|\Delta''\|_1 - \|\Delta'\|_1.$$

Combining the pieces yields that $\|\Delta''\|_1 \leq \|\Delta'\|_1$, and hence $\|\Delta\|_1 \leq 2\|\Delta'\|_1$. By Lemma 1(a), the rank of Δ' is at most $2r$, so that we obtain $\|\Delta\|_1 \leq 2\sqrt{2r}\|\Delta\|_F \leq 4\sqrt{r}\|\Delta\|_F$.

Note that $\mathfrak{X}(\Delta) = 0$, since both $\widehat{\Theta}$ and Θ^* agree with the observations. Consequently, from Proposition 1, we have that

$$\begin{aligned} 0 &= \frac{\|\mathfrak{X}(\Delta)\|_2}{\sqrt{N}} \geq \frac{1}{4}\|\Delta\|_F - 12\rho(\Sigma) \left(\sqrt{\frac{m_1}{N}} + \sqrt{\frac{m_2}{N}} \right) \|\Delta\|_1 \\ &\geq \|\Delta\|_F \left(\frac{1}{4} - 12\rho(\Sigma) \sqrt{\frac{rm_1}{N}} + 12\rho(\Sigma) \sqrt{\frac{rm_2}{N}} \right) \\ &\geq \frac{1}{20}\|\Delta\|_F \end{aligned}$$

where the final inequality as long as $N > c_0\rho^2(\Sigma)r(m_1+m_2)$ for a sufficiently large constant c_0 . We have thus shown that $\Delta = 0$, which implies that $\widehat{\Theta} = \Theta^*$ as claimed.

5. Experimental results. In this section, we report the results of various simulations that demonstrate the close agreement between the scaling predicted by our theory, and the actual behavior of the SDP-based M -estimator (9) in practice. In all cases, we solved the convex program (9) by using our own implementation in MATLAB of an accelerated gradient descent method which adapts a non-smooth convex optimization procedure [42] to the nuclear-norm [29]. We chose the regularization parameter λ_N in the manner suggested by our theoretical results; in doing so, we assumed knowledge of the noise variance ν^2 . In practice, one would have to estimate such quantities from the data using methods such as cross-validation, as has been studied in the context of the Lasso, and we leave this as an interesting direction for future research.

We report simulation results for three of the running examples discussed in this paper: low-rank multivariate regression, estimation in vector autoregressive processes, and matrix recovery from random projections (compressed sensing). In each case, we solved instances of the SDP for a square matrix $\Theta^* \in \mathbb{R}^{m \times m}$, where $m \in \{40, 80, 160\}$ for the first two examples, and $m \in \{20, 40, 80\}$ for the compressed sensing example. In all cases, we considered the case of exact low rank constraints, with $\text{rank}(\Theta^*) = r = 10$, and we generated Θ^* by choosing the subspaces of its left and right singular vectors uniformly at random from the Grassman manifold.⁴ The observation noise had variance $\nu^2 = 1$, and we chose $C = \nu^2 I$ for the VAR process. The VAR process was generated by first solving for the covariance matrix Σ using the MATLAB function `dylap` and then generating a sample path. For each setting of (r, m) , we solved the SDP for a range of sample sizes N .

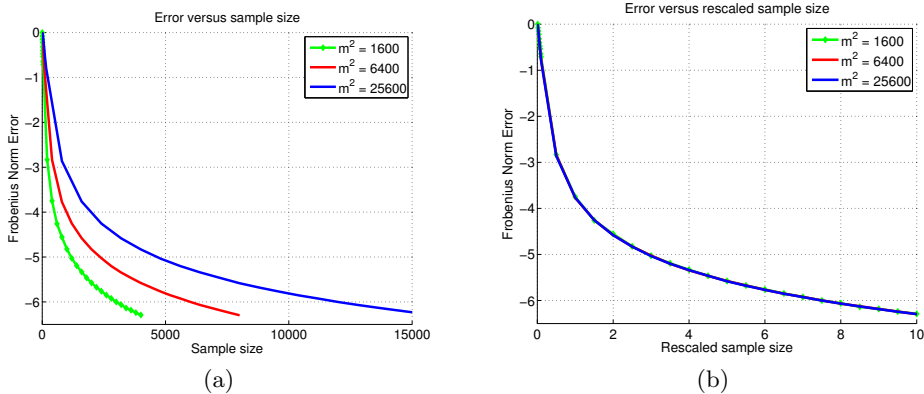


Fig 1. Results of applying the SDP (9) with nuclear norm regularization to the problem of low-rank multivariate regression. (a) Plots of the Frobenius error $\|\hat{\Theta} - \Theta^*\|_F$ on a logarithmic scale versus the sample size N for three different matrix sizes $m^2 \in \{1600, 6400, 25600\}$, all with rank $r = 10$. (b) Plots of the same Frobenius error versus the rescaled sample size $N/(rm)$. Consistent with theory, all three plots are now extremely well-aligned.

Figure 1 shows results for a multivariate regression model with the covariates chosen randomly from a $N(0, I)$ distribution. Panel (a) plots the Frobenius error $\|\hat{\Theta} - \Theta^*\|_F$ on a logarithmic scale versus the sample size N for three different matrix sizes, $m \in \{40, 80, 160\}$. Naturally, in each case, the error decays to zero as N increases, but larger matrices require larger sample sizes, as reflected by the rightward shift of the curves as m is increased. Panel (b) of Figure 1 shows the exact same set of simulation results, but now with

⁴More specifically, we let $\Theta^* = XY^T$, where $X, Y \in \mathbb{R}^{m \times r}$ have i.i.d. $N(0, 1)$ elements.

the Frobenius error plotted versus the rescaled sample size $\tilde{N} := N/(rm)$. As predicted by Corollary 3, the error plots now are all aligned with one another; the degree of alignment in this particular case is so close that the three plots are now indistinguishable. (The blue curve is the only one visible since it was plotted last by our routine.) Consequently, Figure 1 shows that $N/(rm)$ acts as the effective sample size in this high-dimensional setting.

Figure 2 shows similar results for the autoregressive model discussed in Example 2. As shown in panel (a), the Frobenius error again decays as the sample size is increased, although problems involving larger matrices are shifted to the right. Panel (b) shows the same Frobenius error plotted versus the rescaled sample size $N/(rm)$; as predicted by Corollary 4, the errors for different matrix sizes m are again quite well-aligned. In this case, we find (both in our theoretical analysis and experimental results) that the dependence in the autoregressive process slows down the rate at which the concentration occurs, so that the results are not as crisp as the low-rank multivariate setting in Figure 1.

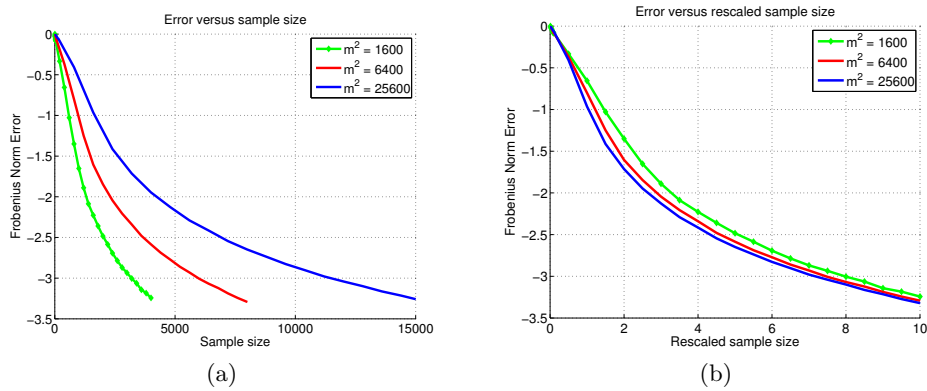


Fig 2. Results of applying the SDP (9) with nuclear norm regularization to estimating the system matrix of a vector autoregressive process. (a) Plots of the Frobenius error $\|\hat{\Theta} - \Theta^*\|_F$ on a logarithmic scale versus the sample size N for three different matrix sizes $m^2 \in \{1600, 6400, 25600\}$, all with rank $r = 10$. (b) Plots of the same Frobenius error versus the rescaled sample size $N/(rm)$. Consistent with theory, all three plots are now reasonably well-aligned.

Finally, Figure 3 presents the same set of results for the compressed sensing observation model discussed in Example 3. Even though the observation matrices X_i here are qualitatively different (in comparison to the multivariate regression and autoregressive examples), we again see the “stacking” phenomenon of the curves when plotted versus the rescaled sample size

N/rm , as predicted by Corollary 5.

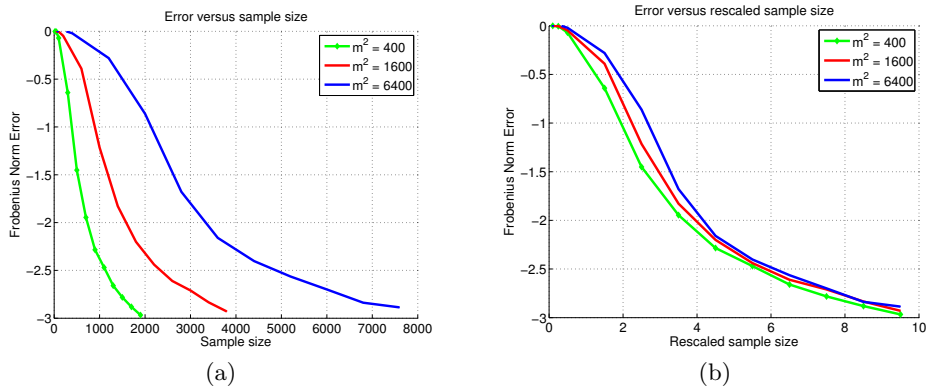


Fig 3. Results of applying the SDP (9) with nuclear norm regularization to recovering a low-rank matrix on the basis of random projections (compressed sensing model) (a) Plots of the Frobenius error $\|\hat{\Theta} - \Theta^*\|_F$ on a logarithmic scale versus the sample size N for three different matrix sizes $m^2 \in \{400, 1600, 6400\}$, all with rank $r = 10$. (b) Plots of the same Frobenius error versus the rescaled sample size $N/(rm)$. Consistent with theory, all three plots are now reasonably well-aligned.

6. Discussion. In this paper, we have analyzed the nuclear norm relaxation for a general class of noisy observation models, and obtained non-asymptotic error bounds on the Frobenius norm that hold under high-dimensional scaling. In contrast to most past work, our results are applicable to both exactly and approximately low-rank matrices. We stated a main theorem that provides high-dimensional rates in a fairly general setting, and then showed how by specializing this result to some specific model classes—namely, low-rank multivariate regression, estimation of autoregressive processes, and matrix recovery from random projections—it yields concrete and readily interpretable rates. Lastly, we provided some simulation results that showed excellent agreement with the predictions from our theory.

This paper has focused on achievable results for low-rank matrix estimation using a particular polynomial-time method. It would be interesting to establish matching lower bounds, showing that the rates obtained by this estimator are minimax-optimal. We suspect that this should be possible, for instance by using the techniques exploited in Raskutti et al. [45] in analyzing minimax rates for regression over ℓ_q -balls.

Acknowledgements. This work was partially supported by a Sloan Foundation Fellowship, an AFOSR-09NL184 grant, and NSF grants CDI-

0941742 and DMS-0907632 to MJW.

APPENDIX A: PROOF OF LEMMA 1

Part (a) of the claim was proved in Recht et al. [47]; we simply provide a proof here for completeness. We write the SVD as $\Theta^* = UDV^T$, where $U \in \mathbb{R}^{m_1 \times m_1}$ and $V \in \mathbb{R}^{m_2 \times m_2}$ are orthogonal matrices, and D is the matrix formed by the singular values of Θ^* . Note that the matrices U^r and V^r are given by the first r columns of U and V respectively. We then define the matrix $\Gamma = U^T \Delta V \in \mathbb{R}^{m_1 \times m_2}$, and write it in block form as

$$\Gamma = \begin{bmatrix} \Gamma_{11} & \Gamma_{12} \\ \Gamma_{21} & \Gamma_{22} \end{bmatrix}, \quad \text{where } \Gamma_{11} \in \mathbb{R}^{r \times r}, \text{ and } \Gamma_{22} \in \mathbb{R}^{(m_1-r) \times (m_2-r)}.$$

We now define the matrices

$$\Delta'' = U \begin{bmatrix} 0 & 0 \\ 0 & \Gamma_{22} \end{bmatrix} V^T, \quad \text{and } \Delta' = \Delta - \Delta''.$$

Note that we have

$$\text{rank}(\Delta') = \text{rank} \begin{bmatrix} \Gamma_{11} & \Gamma_{12} \\ \Gamma_{21} & 0 \end{bmatrix} \leq \text{rank} \begin{bmatrix} \Gamma_{11} & \Gamma_{12} \\ 0 & 0 \end{bmatrix} + \text{rank} \begin{bmatrix} \Gamma_{11} & 0 \\ \Gamma_{21} & 0 \end{bmatrix} \leq 2r,$$

which establishes Lemma 1(a). Moreover, we note for future reference that by construction of Δ'' , the nuclear norm satisfies the decomposition

$$(40) \quad \|\Pi_{\mathcal{A}^r}(\Theta^*) + \Delta''\|_1 = \|\Pi_{\mathcal{A}^r}(\Theta^*)\|_1 + \|\Delta''\|_1.$$

We now turn to the proof of Lemma 1(b). Recall that the error $\Delta = \widehat{\Theta} - \Theta^*$ associated with any optimal solution must satisfy the inequality (30), which implies that

$$(41) \quad 0 \leq \frac{1}{N} \langle \vec{\varepsilon}, \mathfrak{X}(\Delta) \rangle + \lambda_N \{ \|\Theta^*\|_1 - \|\widehat{\Theta}\|_1 \} \leq \frac{1}{N} \mathfrak{X}^*(\vec{\varepsilon})_{\text{op}} \|\Delta\|_1 + \lambda_N \{ \|\Theta^*\|_1 - \|\widehat{\Theta}\|_1 \},$$

where we have used the bound (31).

Note that we have the decomposition $\Theta^* = \Pi_{\mathcal{A}^r}(\Theta^*) + \Pi_{\mathcal{B}^r}(\Theta^*)$. Using this decomposition, the triangle inequality and the relation (40), we have

$$\begin{aligned} \|\widehat{\Theta}\|_1 &= \|(\Pi_{\mathcal{A}^r}(\Theta^*) + \Delta'') + (\Pi_{\mathcal{B}^r}(\Theta^*) + \Delta')\|_1 \\ &\geq \|(\Pi_{\mathcal{A}^r}(\Theta^*) + \Delta'')\|_1 - \|(\Pi_{\mathcal{B}^r}(\Theta^*) + \Delta')\|_1 \\ &\geq \|\Pi_{\mathcal{A}^r}(\Theta^*)\|_1 + \|\Delta''\|_1 - \{ \|(\Pi_{\mathcal{B}^r}(\Theta^*))\|_1 + \|\Delta'\|_1 \}. \end{aligned}$$

Consequently, we have

$$\begin{aligned} \|\Theta^*\|_1 - \|\widehat{\Theta}\|_1 &\leq \|\Theta^*\|_1 - \{\|\Pi_{\mathcal{A}^r}(\Theta^*)\|_1 + \|\Delta''\|_1\} + \{\|(\Pi_{\mathcal{B}^r}(\Theta^*))\|_1 + \|\Delta'\|_1\} \\ &= 2\|\Pi_{\mathcal{B}^r}(\Theta^*)\|_1 + \|\Delta'\|_1 - \|\Delta''\|_1. \end{aligned}$$

Substituting this inequality into the bound (41), we obtain

$$0 \leq \frac{1}{N} \mathfrak{X}^*(\varepsilon) \|\Delta\|_1 + \lambda_N \{2\|\Pi_{\mathcal{B}^r}(\Theta^*)\|_1 + \|\Delta'\|_1 - \|\Delta''\|_1\}.$$

Finally, since $\frac{1}{N} \mathfrak{X}^*(\varepsilon) \|\Delta\|_1 \leq \lambda_N/2$ by assumption, we conclude that

$$0 \leq \lambda_N \{2\|\Pi_{\mathcal{B}^r}(\Theta^*)\|_1 + \frac{3}{2}\|\Delta'\|_1 - \frac{1}{2}\|\Delta''\|_1\}.$$

Since $\|\Pi_{\mathcal{B}^r}(\Theta^*)\|_1 = \sum_{j=r+1}^m \sigma_j(\Theta^*)$, the bound (32) follows.

APPENDIX B: CONSISTENCY IN OPERATOR NORM

In this appendix, we derive a bound on the operator norm error for both the low-rank multivariate regression and auto-regressive model estimation problems. In this statement, it is convenient to specify these models in the form $Y = X\Theta^* + W$, where $Y \in \mathbb{R}^{n \times m_2}$ is a matrix of observations.

PROPOSITION 2 (Operator norm consistency). *Consider the multivariate regression problem and the SDP under the conditions of Corollary 3. Then any solution $\widehat{\Theta}$ to the SDP satisfies the bound*

$$(42) \quad \|\widehat{\Theta} - \Theta^*\|_{\text{op}} \leq c' \frac{\nu \sqrt{\sigma_{\max}(\Sigma)}}{\sigma_{\min}(\Sigma)} \sqrt{\frac{m_1 + m_2}{n}}.$$

We note that a similar bound applies to the auto-regressive model treated in Corollary 4.

PROOF. For any subgradient matrix $Z \in \partial\|\widehat{\Theta}\|_1$, we are guaranteed $\|Z\|_{\text{op}} \leq 1$. Furthermore, by the KKT conditions [11] for the nuclear norm SDP, any solution $\widehat{\Theta}$ must satisfy the condition

$$\frac{1}{n} X^T X \widehat{\Theta} - \frac{X^T Y}{n} + \lambda_n Z = 0.$$

Hence, simple algebra and the triangle inequality yield that

$$\|\widehat{\Theta}\|_{\text{op}} \leq \left\| \left(\frac{1}{n} X^T X \right)^{-1} \right\|_{\text{op}} \left[\|X^T W/n\|_{\text{op}} + \lambda_n \right].$$

Lemma 2 yields that $\|(\frac{1}{n}X^T X)^{-1}\|_{\text{op}} \leq \frac{9}{\sigma_{\min}(\Sigma)}$ with high probability. Combining these inequalities yields

$$\|\hat{\Theta}\|_{\text{op}} \leq c_1 \frac{\lambda_n}{\sigma_{\min}(\Sigma)}.$$

We require that $\lambda_n \geq 2\|X^T W\|_{\text{op}}/n$. From Lemma 3, it suffices to set $\lambda_n \geq c_0 \sqrt{\sigma_{\max}(\Sigma)} \nu \sqrt{\frac{m_1+m_2}{n}}$. Combining the pieces yields the claim. \square

APPENDIX C: PROOF OF LEMMA 3

Let $S^{m-1} = \{u \in \mathbb{R}^m \mid \|u\|_2 = 1\}$ denote the Euclidean sphere in m -dimensions. The operator norm of interest has the variational representation

$$\frac{1}{n}\|X^T W\|_{\text{op}} = \frac{1}{n} \sup_{u \in S^{m_1-1}} \sup_{v \in S^{m_2-1}} v^T X^T W u$$

For positive scalars a and b , define the (random) quantity

$$\Psi(a, b) := \sup_{u \in a S^{m_1-1}} \sup_{v \in b S^{m_2-1}} \langle Xv, Wu \rangle.$$

and note that our goal is to upper bound $\Psi(1, 1)$. Note moreover that $\Psi(a, b) = ab\Psi(1, 1)$, a relation which will be useful in the analysis.

Let $\mathcal{A} = \{u^1, \dots, u^A\}$ and $\mathcal{B} = \{v^1, \dots, v^B\}$ denote $1/4$ coverings of S^{m_1-1} and S^{m_2-1} , respectively. We now claim that we have the upper bound

$$(43) \quad \Psi(1, 1) \leq 4 \max_{u^a \in \mathcal{A}, v^b \in \mathcal{B}} \langle Xv^b, Wu^a \rangle$$

To establish this claim, we note that since the sets \mathcal{A} and \mathcal{B} are $1/4$ -covers, for any pair $(u, v) \in S^{m_1-1} \times S^{m_2-1}$, there exists a pair $(u^a, v^b) \in \mathcal{A} \times \mathcal{B}$ such that $u = u^a + \Delta u$ and $v = v^b + \Delta v$, with $\max\{\|\Delta u\|_2, \|\Delta v\|_2\} \leq 1/4$. Consequently, we can write

$$(44) \quad \langle Xv, Wu \rangle = \langle Xv^b, Wu^a \rangle + \langle Xv^b, W\Delta u \rangle + \langle X\Delta v, Wu^a \rangle + \langle X\Delta v, W\Delta u \rangle.$$

By construction, we have the bound $|\langle Xv^b, W\Delta u \rangle| \leq \Psi(1, 1/4) = \frac{1}{4}\Psi(1, 1)$, and similarly $|\langle X\Delta v, Wu^a \rangle| \leq \frac{1}{4}\Psi(1, 1)$ as well as $|\langle X\Delta v, W\Delta u \rangle| \leq \frac{1}{16}\Psi(1, 1)$. Substituting these bounds into the decomposition (44) and taking suprema over the left and right-hand sides, we conclude that

$$\Psi(1, 1) \leq \max_{u^a \in \mathcal{A}, v^b \in \mathcal{B}} \langle Xv^b, Wu^a \rangle + \frac{9}{16}\Psi(1, 1),$$

from which the bound (43) follows.

We now apply the union bound to control the discrete maximum. It is known (e.g., [33, 39]) that there exists a $1/4$ covering of S^{m_1-1} and S^{m_2-1} with at most $A \leq 8^{m_1}$ and $B \leq 8^{m_2}$ elements respectively. Consequently, we have

$$(45) \quad \mathbb{P}[|\Psi(1, 1)| \geq 4\delta n] \leq 8^{m_1+m_2} \max_{u^a, v^b} \mathbb{P} \left[\frac{|\langle Xv^b, Wu^a \rangle|}{n} \geq \delta \right].$$

It remains to obtain a good bound on the quantity $\frac{1}{n}\langle Xv, Wu \rangle = \frac{1}{n} \sum_{i=1}^n \langle v, X_i \rangle \langle u, W_i \rangle$, where $(u, v) \in S^{m_1-1} \times S^{m_2-1}$ are arbitrary but fixed. Since $W_i \in \mathbb{R}^{m_1}$ has i.i.d. $N(0, \nu^2)$ elements and u is fixed, we have $Z_i := \langle u, W_i \rangle \sim N(0, \nu^2)$ for each $i = 1, \dots, n$. These variables are independent of one another, and of the random matrix X . Therefore, conditioned on X , the sum $Z := \frac{1}{n} \sum_{i=1}^n \langle v, X_i \rangle \langle u, W_i \rangle$ is zero-mean Gaussian with variance

$$\alpha^2 := \frac{\nu^2}{n} \left(\frac{1}{n} \|Xv\|_2^2 \right) \leq \frac{\nu^2}{n} \|X^T X/n\|_{\text{op}}.$$

Define the event $\mathcal{T} = \{\alpha^2 \leq \frac{9\nu^2 \|\Sigma\|_{\text{op}}}{n}\}$. Using Lemma 2, we have $\|X^T X/n\|_{\text{op}} \leq 9\sigma_{\max}(\Sigma)$ with probability at least $1 - 2\exp(-n/2)$, which implies that $\mathbb{P}[\mathcal{T}^c] \leq 2\exp(-n/2)$. Therefore, conditioning on the event \mathcal{T} and its complement \mathcal{T}^c , we obtain

$$\begin{aligned} \mathbb{P}[|Z| \geq t] &\leq \mathbb{P}[|Z| \geq t \mid \mathcal{T}] + \mathbb{P}[\mathcal{T}^c] \\ &\leq \exp\left(-n \frac{t^2}{2\nu^2(4 + \|\Sigma\|_{\text{op}})}\right) + 2\exp(-n/2). \end{aligned}$$

Combining this tail bound with the upper bound (45), we have

$$\mathbb{P}[|\psi(1, 1)| \geq 4\delta n] \leq 8^{m_1+m_2} \left\{ \exp\left(-n \frac{t^2}{18\nu^2 \|\Sigma\|_{\text{op}}}\right) + 2\exp(-n/2) \right\}.$$

Setting $t^2 = 20\nu^2 \|\Sigma\|_{\text{op}} \frac{m_1+m_2}{n}$, this probability vanishes as long as $n > 16(m_1 + m_2)$.

APPENDIX D: TECHNICAL DETAILS FOR COROLLARY 4

In this appendix, we collect the proofs of Lemmas 4 and 5.

D.1. Proof of Lemma 4. Recalling that S^{m-1} denotes the unit-norm Euclidean sphere in m -dimensions, we first observe that $\|X\|_{\text{op}} = \sup_{u \in S^{m-1}} \|Xu\|_2$. Our next step is to reduce the supremum to a maximization over a finite set, using a standard covering argument. Let $\mathcal{A} = \{u^1, \dots, u^A\}$ denote a $1/2$ -cover of it. By definition, for any $u \in S^{m-1}$, there is some $u^a \in \mathcal{A}$ such that $u = u^a + \Delta u$, where $\|\Delta u\|_2 \leq 1/2$. Consequently, for any $u \in S^{m-1}$, the triangle inequality implies that

$$\|Xu\|_2 \leq \|Xu^a\|_2 + \|X\Delta u\|_2,$$

and hence that $\|X\|_{\text{op}} \leq \max_{u^a \in \mathcal{A}} \|Xu^a\|_2 + \frac{1}{2}\|X\|_{\text{op}}$. Re-arranging yields the useful inequality

$$(46) \quad \|X\|_{\text{op}} \leq 2 \max_{u^a \in \mathcal{A}} \|Xu^a\|_2.$$

Using inequality (46), we have

$$(47) \quad \begin{aligned} \mathbb{P}\left[\frac{1}{n}\|X^T X\|_{\text{op}} > t\right] &\leq \mathbb{P}\left[\max_{u^a \in \mathcal{A}} \frac{1}{n} \sum_{i=1}^n (\langle u^a, X_i \rangle)^2 > \frac{t}{2}\right] \\ &\leq 4^m \max_{u^a \in \mathcal{A}} \mathbb{P}\left[\frac{1}{n} \sum_{i=1}^n (\langle u^a, X_i \rangle)^2 > \frac{t}{2}\right]. \end{aligned}$$

where the last inequality follows from the union bound, and the fact [33, 39] that there exists a $1/2$ -covering of S^{m-1} with at most 4^m elements.

In order to complete the proof, we need to obtain a sharp upper bound on the quantity $\mathbb{P}[\frac{1}{n} \sum_{i=1}^n (\langle u, X_i \rangle)^2 > \frac{t}{2}]$, valid for any fixed $u \in S^{m-1}$. Define the random vector $Y \in \mathbb{R}^n$ with elements $Y_i = \langle u, X_i \rangle$. Note that Y is zero mean, and its covariance matrix R has elements $R_{ij} = \mathbb{E}[Y_i Y_j] = u^T \Sigma (\Theta^*)^{|j-i|} u$. In order to bound the spectral norm of R , we note that since it is symmetric, we have $\|R\|_{\text{op}} \leq \max_{i=1, \dots, m} \sum_{j=1}^m |R_{ij}|$, and moreover

$$|R_{ij}| = |u^T \Sigma (\Theta^*)^{|j-i|} u| \leq (\|\Theta^*\|_{\text{op}})^{|j-i|} \Sigma \leq \gamma^{|j-i|} \|\Sigma\|_{\text{op}}.$$

Combining the pieces, we obtain

$$(48) \quad \|R\|_{\text{op}} \leq \max_i \sum_{j=1}^m |\gamma|^{|i-j|} \|\Sigma\|_{\text{op}} \leq 2\|\Sigma\|_{\text{op}} \sum_{j=0}^{\infty} |\gamma|^j \leq \frac{2\|\Sigma\|_{\text{op}}}{1-\gamma}.$$

Moreover, we have $\text{trace}(R)/n = u^T \Sigma u \leq \|\Sigma\|_{\text{op}}$. Applying Lemma 9 with $t = 5\sqrt{\frac{m}{n}}$, we conclude that

$$\mathbb{P}\left[\frac{1}{n} \|Y\|_2^2 > \|\Sigma\|_{\text{op}} + 5\sqrt{\frac{m}{n}} \|R\|_{\text{op}}\right] \leq 2 \exp(-5m) + 2 \exp(-n/2)..$$

Combined with the bound (47), we obtain

$$(49) \quad \left\| \frac{1}{n} X^T X \right\|_{\text{op}} \leq \|\Sigma\|_{\text{op}} \left\{ 2 + \frac{20}{(1-\gamma)} \sqrt{\frac{m}{n}} \right\} \leq \frac{24\|\Sigma\|_{\text{op}}}{(1-\gamma)},$$

with probability at least $1 - c_1 \exp(-c_2 m)$, which establishes the upper bound (35)(a).

Turning to the lower bound (35)(b), we let $\mathcal{B} = \{v^1, \dots, v^B\}$ be an ϵ -cover of S^{m-1} for some $\epsilon \in (0, 1)$ to be chosen. Thus, for any $v \in \mathbb{R}^m$, there exists some v^b such that $v = v^b + \Delta v$, and $\|\Delta v\|_2 \leq \epsilon$. Define the function $\Psi : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ via $\Psi(u, v) = u^T \left(\frac{1}{n} X^T X \right) v$, and note that $\Psi(u, v) = \Psi(v, u)$. With this notation, we have

$$\begin{aligned} v^T \left(\frac{1}{n} X^T X \right) v &= \Psi(v, v) = \Psi(v^b, v^b) + 2\Psi(\Delta v, v) + \Psi(\Delta v, \Delta v) \\ &\geq \Psi(v^b, v^b) + 2\Psi(\Delta v, v), \end{aligned}$$

since $\Psi(\Delta v, \Delta v) \geq 0$. Since $|\Psi(\Delta v, v)| \leq \epsilon \left\| \left(\frac{1}{n} X^T X \right) \right\|_{\text{op}}$, we obtain the lower bound

$$\sigma_{\min} \left(\left(\frac{1}{n} X^T X \right) \right) = \inf_{v \in S^{m-1}} v^T \left(\frac{1}{n} X^T X \right) v \geq \min_{v^b \in \mathcal{B}} \Psi(v^b, v^b) - 2\epsilon \left\| \frac{1}{n} X^T X \right\|_{\text{op}}.$$

By the previously established upper bound(35)(a), have $\left\| \frac{1}{n} X^T X \right\|_{\text{op}} \leq \frac{24\|\Sigma\|_{\text{op}}}{(1-\gamma)}$ with high probability. Hence, choosing $\epsilon = \frac{(1-\gamma)\sigma_{\min}(\Sigma)}{200\|\Sigma\|_{\text{op}}}$ ensures that $2\epsilon \left\| \frac{1}{n} X^T X \right\|_{\text{op}} \leq \sigma_{\min}(\Sigma)/4$.

Consequently, it suffices to lower bound the minimum over the covering set. We first establish a concentration result for the function $\Psi(v, v)$ that holds for any fixed $v \in S^{m-1}$. Note that we can write

$$\Psi(v, v) = \frac{1}{n} \sum_{i=1}^n (\langle v, X_i \rangle)^2,$$

As before, if we define the random vector $Y \in \mathbb{R}^n$ with elements $Y_i = \langle v, X_i \rangle$, then $Y \sim N(0, R)$ with $\|R\|_{\text{op}} \leq \frac{2\|\Sigma\|_{\text{op}}}{1-\gamma}$. Moreover, we have $\text{trace}(R)/n = v^T \Sigma v \geq \sigma_{\min}(\Sigma)$. Consequently, applying Lemma 9 yields

$$\mathbb{P} \left[\frac{1}{n} \|Y\|_2^2 < \sigma_{\min}(\Sigma) - \frac{8t\|\Sigma\|_{\text{op}}}{1-\gamma} \right] \leq 2 \exp(-n(t - 2/\sqrt{n})^2/2) + 2 \exp(-\frac{n}{2}),$$

Note that this bound holds for any fixed $v \in S^{m-1}$. Setting $t^* = \frac{(1-\gamma)\sigma_{\min}(\Sigma)}{16\|\Sigma\|_{\text{op}}}$ and applying the union bound yields that

$$\mathbb{P} \left[\min_{v^b \in \mathcal{B}} \Psi(v^b, v^b) < \sigma_{\min}(\Sigma)/2 \right] \leq \left(\frac{4}{\epsilon} \right)^m \left\{ 2 \exp(-n(t^* - 2/\sqrt{n})^2/2) + 2 \exp(-\frac{n}{2}) \right\},$$

which vanishes as long as $n > \frac{4 \log(4/\epsilon)}{(t^*)^2} m$.

D.2. Proof of Lemma 5. Let $S^{m-1} = \{u \in \mathbb{R}^m \mid \|u\|_2 = 1\}$ denote the Euclidean sphere in m -dimensions, and for positive scalars a and b , define the random variable

$$\Psi(a, b) := \sup_{u \in a S^{m-1}} \sup_{v \in b S^{m-1}} \langle Xv, Wu \rangle.$$

Note that our goal is to upper bound $\Psi(1, 1)$. Let $\mathcal{A} = \{u^1, \dots, u^A\}$ and $\mathcal{B} = \{v^1, \dots, v^B\}$ denote $1/4$ coverings of S^{m-1} and S^{m-1} , respectively. Following the same argument as in the proof of Lemma 3, we obtain the upper bound

$$(50) \quad \Psi(1, 1) \leq 4 \max_{u^a \in \mathcal{A}, v^b \in \mathcal{B}} \langle Xv^b, Wu^a \rangle$$

We now apply the union bound to control the discrete maximum. It is known (e.g., [33, 39]) that there exists a $1/4$ covering of S^{m-1} with at most 8^m elements. Consequently, we have

$$(51) \quad \mathbb{P}[|\psi(1, 1)| \geq 4\delta n] \leq 8^{2m} \max_{u^a, v^b} \mathbb{P}\left[\frac{|\langle Xv^b, Wu^a \rangle|}{n} \geq \delta\right].$$

It remains to obtain a tail bound on the quantity $\mathbb{P}\left[\frac{|\langle Xv, Wu \rangle|}{n} \geq \delta\right]$, for any fixed pair $(u, v) \in \mathcal{A} \times \mathcal{B}$.

For each $i = 1, \dots, n$, let X_i and W_i denote the i^{th} row of X and W . Following some simple algebra, we have the decomposition $\frac{\langle Xv, Wu \rangle}{n} = T_1 - T_2 - T_3$, where

$$\begin{aligned} T_1 &= \frac{1}{2n} \sum_{i=1}^n (\langle u, W_i \rangle + \langle v, X_i \rangle)^2 - \frac{1}{2}(u^T C u + v^T \Sigma v) \\ T_2 &= \frac{1}{2n} \sum_{i=1}^n (\langle u, W_i \rangle)^2 - \frac{1}{2} u^T C u \\ T_3 &= \frac{1}{2n} \sum_{i=1}^n (\langle v, X_i \rangle)^2 - \frac{1}{2} v^T \Sigma v \end{aligned}$$

We may now bound each T_j for $j = 1, 2, 3$ in turn; in doing so, we make repeated use of Lemma 9, which provides concentration bounds for a random variable of the form $\|Y\|_2^2$, where $Y \sim N(0, Q)$ for some matrix $Q \succeq 0$.

Bound on T_3 . We can write the term T_3 as a deviation of $\|Y\|_2^2/n$ from its mean, where in this case the covariance matrix Q is no longer the identity. In concrete terms, let us define a random vector $Y \in \mathbb{R}^n$ with elements $Y_i := \langle v, X_i \rangle$. As seen in the proof of Lemma 4 from Appendix D.1, the vector Y is zero-mean Gaussian with covariance matrix R such that $\|R\|_{\text{op}} \leq \frac{2\|\Sigma\|_{\text{op}}}{1-\gamma}$ (see equation (48)). Since we have $\text{trace}(R)/n = v^T R v$, applying Lemma 9 yields that

$$(52) \quad \mathbb{P}[|T_3| \geq \frac{8\|\Sigma\|_{\text{op}}}{1-\gamma}t] \leq 2 \exp\left(-\frac{n(t - 2/\sqrt{n})^2}{2}\right) + 2 \exp(-n/2).$$

Bound on T_2 . We control the term T_2 in a similar way. Define the random vector $Y' \in \mathbb{R}^n$ with elements $Y'_i := \langle u, W_i \rangle$. Then Y' is a sample from the distribution $N(0, (u^T C u)I_{n \times n})$, so that $\frac{2}{u^T C u}T_2$ is the difference between a rescaled χ^2 variable and its mean. Applying Lemma 9 with $Q = (u^T C u)I$, we obtain

$$(53) \quad \mathbb{P}[|T_2| > 4(u^T C u)t] \leq 2 \exp\left(-\frac{n(t - 2/\sqrt{n})^2}{2}\right) + 2 \exp(-n/2).$$

Bound on T_1 . To control this quantity, let us define a zero-mean Gaussian random vector $Z \in \mathbb{R}^n$ with elements $Z_i = \langle v, X_i \rangle + \langle u, W_i \rangle$. This random vector has covariance matrix S with elements

$$S_{ij} = \mathbb{E}[Z_i Z_j] = (u^T C u)\delta_{ij} + (1 - \delta_{ij})(u^T C u)v^T (\Theta^*)^{|i-j|-1}u + v^T (\Theta^*)^{|i-j|}\Sigma v,$$

where δ_{ij} is the Kronecker delta for the event $\{i = j\}$. As before, by symmetry of S , we have $\|S\|_{\text{op}} \leq \max_{i=1, \dots, n} \sum_{j=1}^n |S_{ij}|$, and hence

$$\begin{aligned} \|S\|_{\text{op}} &\leq (u^T C u) + \|\Sigma\|_{\text{op}} + \sum_{j=1}^{i-1} |(u^T C u) v^T (\Theta^*)^{|i-j|-1}u + v^T (\Theta^*)^{|i-j|}\Sigma v| \\ &\quad + \sum_{j=i+1}^n |(u^T C u) v^T (\Theta^*)^{|i-j|-1}u + v^T (\Theta^*)^{|i-j|}\Sigma v|. \end{aligned}$$

Since $\|\Theta^*\|_{\text{op}} \leq \gamma < 1$, and $(u^T C u) \leq \|C\|_{\text{op}} \leq \|\Sigma\|_{\text{op}}$, we have

$$\begin{aligned} \|S\|_{\text{op}} &\leq \|C\|_{\text{op}} + \|\Sigma\|_{\text{op}} + 2 \sum_{j=1}^{\infty} \|C\|_{\text{op}} \gamma^{j-1} + 2 \sum_{j=1}^{\infty} \|\Sigma\|_{\text{op}} \gamma^j \\ &\leq 4 \|\Sigma\|_{\text{op}} \left(1 + \frac{1}{1-\gamma}\right) \end{aligned}$$

Moreover, we have $\frac{\text{trace}(S)}{n} = (u^T C u) + v^T \Sigma v \leq 2\|\Sigma\|_{\text{op}}$, so that by applying Lemma 9, we conclude that

$$(54) \quad \mathbb{P}\left[|T_1| > \left(\frac{24\|\Sigma\|_{\text{op}}}{1-\gamma}\right)t\right] \leq 2\exp\left(-\frac{n(t - 2/\sqrt{n})^2}{2}\right) + 2\exp(-n/2),$$

which completes the analysis of this term.

Combining the bounds (52), (53) and (54), we conclude that for all $t > 0$,

$$(55) \quad \mathbb{P}\left[\frac{|\langle Xv, Wu \rangle|}{n} \geq \frac{40(\|\Sigma\|_{\text{op}} t)}{1-\gamma}\right] \leq 6\exp\left(-\frac{n(t - 2/\sqrt{n})^2}{2}\right) + 6\exp(-n/2).$$

Setting $t = 10\sqrt{m/n}$ and combining with the bound (51), we conclude that

$$\mathbb{P}[|\psi(1, 1)| \geq \frac{1600\|\Sigma\|_{\text{op}}}{1-\gamma}\sqrt{\frac{m}{n}}] \leq 8^{2m} \{6\exp(-16m) + 6\exp(-n/2)\} \leq 12\exp(-m)$$

as long as $n > ((4 \log 8) + 1)m$.

APPENDIX E: PROOF OF PROPOSITION 1

We begin by stating and proving a useful lemma. Recall the definition (22) of $\rho(\Sigma)$.

LEMMA 7. *Let $X \in \mathbb{R}^{m_1 \times m_2}$ be a random sample from the Σ -ensemble. Then we have*

$$(56) \quad \mathbb{E}[\|X\|_{\text{op}}] \leq 12\rho(\Sigma) [\sqrt{m_1} + \sqrt{m_2}]$$

and moreover

$$(57) \quad \mathbb{P}[\|X\|_{\text{op}} \geq \mathbb{E}[\|X\|_{\text{op}}] + t] \leq \exp\left(-\frac{t^2}{2\rho^2(\Sigma)}\right).$$

PROOF. We begin by making note of the variational representation

$$\|X\|_{\text{op}} = \sup_{(u,v) \in S^{m_1-1} \times S^{m_2-1}} u^T X v.$$

Since each variable $u^T X v$ is zero-mean Gaussian, we thus recognize $\|X\|_{\text{op}}$ as the supremum of a Gaussian process. The bound (57) thus follows from Theorem 7.1 in Ledoux [32].

We now use a simple covering argument establish the upper bound (56). Let $\{v^1, \dots, v^{M_2}\}$ be a $1/4$ covering of the sphere S^{m_2-1} . For an arbitrary

$v \in S^{m_2-1}$, there exists some v^j in the cover such that $\|v - v^j\|_2 \leq 1/4$, whence

$$\|Xv\|_2 \leq \|Xv^j\|_2 + \|X(v - v^j)\|_2.$$

Taking suprema over both sides, we obtain that $\|X\|_{\text{op}} \leq \max_{j=1, \dots, M_2} \|Xv^j\|_2 + \frac{1}{4}\|X\|_{\text{op}}$. A similar argument using a $1/4$ -covering $\{u^1, \dots, u^{M_1}\}$ of S^{m_1-1} yields that

$$\|Xv^j\|_2 \leq \max_{i=1, \dots, M_1} \langle u^i, Xv^j \rangle + \frac{1}{4}\|X\|_{\text{op}}.$$

Combining the pieces, we conclude that

$$\|X\|_{\text{op}} \leq 2 \max_{\substack{i=1, \dots, M_1 \\ j=1, \dots, M_2}} \langle u^i, Xv^j \rangle.$$

By construction, each variable $\langle u^i, Xv^j \rangle$ is zero-mean Gaussian with variance at most $\rho(\Sigma)$, so that by standard bounds on Gaussian maxima, we obtain

$$\mathbb{E}[\|X\|_{\text{op}}] \leq 4\rho(\Sigma)\sqrt{\log(M_1M_2)} \leq 4\rho(\Sigma)[\sqrt{\log M_1} + \sqrt{\log M_2}].$$

There exist $1/4$ -coverings of S^{m_1-1} and S^{m_2-1} with $\log M_1 \leq m_1 \log 8$ and $\log M_2 \leq m_2 \log 8$, from which the bound (56) follows. \square

We now return to the proof of Proposition 1. To simplify the proof, let us define an operator $T_\Sigma : \mathbb{R}^{m_1 \times m_2} \rightarrow \mathbb{R}^{m_1 \times m_2}$ such that $\text{vec}(T_\Sigma(\Theta)) = \sqrt{\Sigma} \text{vec}(\Theta)$. Let $\mathfrak{X}' : \mathbb{R}^{m_1 \times m_2} \rightarrow \mathbb{R}^N$ be a random Gaussian operator formed with X'_i sampled with i.i.d. $N(0, 1)$ entries. By construction, we then have $\mathfrak{X}(\Theta) = \mathfrak{X}'(T_\Sigma(\Theta))$ for all $\Theta \in \mathbb{R}^{m_1 \times m_2}$. Now by the variational characterization of the ℓ_2 -norm, we have

$$\|\mathfrak{X}'(T_\Sigma(\Theta))\|_2 = \sup_{u \in S^{N-1}} \langle u, \mathfrak{X}'(T_\Sigma(\Theta)) \rangle.$$

Since the original claim (25) is invariant to rescaling, it suffices to prove it for matrices such that $\|T_\Sigma(\Theta)\|_F = 1$. Letting $t \geq 1$ be a given radius, we seek lower bounds on the quantity

$$Z^*(t) := \inf_{\Theta \in \mathcal{R}(t)} \sup_{u \in S^{N-1}} \langle u, \mathfrak{X}'(T_\Sigma(\Theta)) \rangle, \quad \text{where } \mathcal{R}(t) = \{\Theta \in \mathbb{R}^{m_1 \times m_2} \mid \|T_\Sigma(\Theta)\|_F = 1, \|\Theta\|_1 \leq t\}.$$

In particular, our goal is to prove that for any $t \geq 1$, the lower bound

$$(58) \quad \frac{Z^*(t)}{\sqrt{N}} \geq \frac{1}{4} - 12 \rho(\Sigma) \left[\frac{m_1 + m_2}{N} \right]^{1/2} t$$

holds with probability at least $1 - c_1 \exp(-c_2 N)$. By a standard peeling argument (see Raskutti et al. [45] for details), this lower bound implies the claim (25).

We establish the lower bound (58) using Gaussian comparison inequalities [33] and concentration of measure (see Lemma 8). For each pair $(u, \Theta) \in S^{N-1} \times \mathcal{R}(t)$, consider the random variable $Z_{u,\Theta} = \langle u, \mathfrak{X}'(T_\Sigma(\Theta)) \rangle$, and note that it is Gaussian with zero mean. For any two pairs (u, Θ) and (u', Θ') , some calculation yields

$$(59) \quad \mathbb{E}[(Z_{u,\Theta} - Z_{u',\Theta'})^2] = \|u \otimes T_\Sigma(\Theta) - u' \otimes T_\Sigma(\Theta')\|_F^2.$$

We now define a second Gaussian process $\{Y_{u,\Theta} \mid (u, \Theta) \in S^{N-1} \times \mathcal{R}(t)\}$ via

$$Y_{u,\Theta} := \langle g, u \rangle + \langle\langle G, T_\Sigma(\Theta) \rangle\rangle,$$

where $g \in \mathbb{R}^N$ and $G \in \mathbb{R}^{m_1 \times m_2}$ are independent with i.i.d. $N(0, 1)$ entries. By construction, $Y_{u,\Theta}$ is zero-mean, and moreover, for any two pairs (u, Θ) and (u', Θ') , we have

$$(60) \quad \mathbb{E}[(Y_{u,\Theta} - Y_{u',\Theta'})^2] = \|u - u'\|_2^2 + \|T_\Sigma(\Theta) - T_\Sigma(\Theta')\|_F^2.$$

For all pairs $(u, \Theta), (u', \Theta') \in S^{N-1} \times \mathcal{R}(t)$, we have $\|u\|_2 = \|u'\|_2 = 1$, and moreover $\|T_\Sigma(\Theta)\|_F = \|T_\Sigma(\Theta')\|_F = 1$. Using this fact, some algebra yields that

$$(61) \quad \|u \otimes T_\Sigma(\Theta) - u' \otimes T_\Sigma(\Theta')\|_F^2 \leq \|u - u'\|_2^2 + \|T_\Sigma(\Theta) - T_\Sigma(\Theta')\|_F^2.$$

Moreover, equality holds whenever $\Theta = \Theta'$. The conditions of the Gordon-Slepian inequality [33] are satisfied, so that we are guaranteed that

(62)

$$\mathbb{E} \left[\inf_{\Theta \in \mathcal{R}(t)} \|\mathfrak{X}'(T_\Sigma(\Theta))\|_2 \right] = \mathbb{E} \left[\inf_{\Theta \in \mathcal{R}(t)} \sup_{u \in S^{N-1}} Z_{u,\Theta} \right] \geq \mathbb{E} \left[\inf_{\Theta \in \mathcal{R}(t)} \sup_{u \in S^{N-1}} Y_{u,\Theta} \right]$$

We compute

$$\begin{aligned} \mathbb{E} \left[\inf_{\Theta \in \mathcal{R}(t)} \sup_{u \in S^{N-1}} Y_{u,\Theta} \right] &= \mathbb{E} \left[\sup_{u \in S^{N-1}} \langle g, u \rangle \right] + \mathbb{E} \left[\inf_{\Theta \in \mathcal{R}(t)} \langle\langle G, T_\Sigma(\Theta) \rangle\rangle \right] \\ &= \mathbb{E}[\|g\|_2] - \mathbb{E} \left[\sup_{\Theta \in \mathcal{R}(t)} \langle\langle G, T_\Sigma(\Theta) \rangle\rangle \right] \\ &\geq \frac{1}{2} \sqrt{N} - t \mathbb{E}[\|T_\Sigma(G)\|_{\text{op}}], \end{aligned}$$

where we have used the fact that T_Σ is self-adjoint, and Hölder's inequality (involving the operator and nuclear norms). Since $T_\Sigma(G)$ is a random matrix from the Σ -ensemble, Lemma 7 yields the upper bound $\mathbb{E}[\|T_\Sigma(G)\|_{\text{op}}] \leq 12\rho(\Sigma)(\sqrt{m_1} + \sqrt{m_2})$. Putting together the pieces, we conclude that

$$\mathbb{E}\left[\inf_{\Theta \in \mathcal{R}(t)} \frac{\|\mathfrak{X}'(T_\Sigma(\Theta))\|_2}{\sqrt{N}}\right] \geq \frac{1}{2} - 12\rho(\Sigma) \left(\frac{\sqrt{m_1} + \sqrt{m_2}}{\sqrt{N}}\right) t.$$

Finally, we need to establish sharp concentration around the mean. Since $\|T_\Sigma(\Theta)\|_F = 1$ for all $\Theta \in \mathcal{R}(t)$, the function $f(\mathfrak{X}) := \inf_{\Theta \in \mathcal{R}(t)} \|\mathfrak{X}'(T_\Sigma(\Theta))\|_2 / \sqrt{N}$ is Lipschitz with constant $1/\sqrt{N}$, so that Lemma 8 implies that

$$\mathbb{P}\left[\inf_{\Theta \in \mathcal{R}(t)} \frac{\|\mathfrak{X}(\Theta)\|_2}{\sqrt{N}} \leq \frac{1}{2} - 12\rho(\Sigma) \left(\frac{\sqrt{m_1} + \sqrt{m_2}}{\sqrt{N}}\right) t - \delta\right] \leq 2 \exp(-N\delta^2/2) \quad \text{for all } \delta > 0.$$

Setting $\delta = 1/4$ yields the claim.

APPENDIX F: SOME USEFUL CONCENTRATION RESULTS

The following lemma is classical [33, 38], and yields sharp concentration of a Lipschitz function of Gaussian random variables around its mean.

LEMMA 8. *Let $X \in \mathbb{R}^n$ have i.i.d. $N(0, 1)$ entries, and let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be Lipschitz with constant L (i.e., $|f(x) - f(y)| \leq L\|x - y\|_2 \forall x, y \in \mathbb{R}^n$). Then for all $t > 0$, we have*

$$\mathbb{P}[|f(X) - Ef(X)| > t] \leq 2 \exp\left(-\frac{t^2}{2L^2}\right).$$

By exploiting this lemma, we can prove the following result, which yields concentration of the squared ℓ_2 -norm of an arbitrary Gaussian vector:

LEMMA 9. *Given a Gaussian random vector $Y \sim N(0, Q)$, for all $t > 2/\sqrt{n}$, we have*

(63)

$$\mathbb{P}\left[\frac{1}{n} \|\|Y\|_2^2 - \text{trace } Q\| > 4t \|Q\|_{\text{op}}\right] \leq 2 \exp\left(-\frac{n(t - \frac{2}{\sqrt{n}})^2}{2}\right) + 2 \exp(-n/2).$$

PROOF. Let \sqrt{Q} be the symmetric matrix square root, and consider the function $f(x) = \|\sqrt{Q}x\|_2 / \sqrt{n}$. Since it is Lipschitz with constant $\|\sqrt{Q}\|_{\text{op}} / \sqrt{n}$, Lemma 8 implies that

(64)

$$\mathbb{P}[|\|\sqrt{Q}X\|_2 - E\|\sqrt{Q}X\|_2| > \sqrt{n}\delta] \leq 2 \exp\left(-\frac{n\delta^2}{2\|Q\|_{\text{op}}}\right) \quad \text{for all } \delta > 0.$$

By integrating this tail bound, we find that the variable $Z = \|\sqrt{Q}X\|_2/\sqrt{n}$ satisfies the bound $\text{var}(Z) \leq 4\|Q\|_{\text{op}}/n$, and hence conclude that

(65)

$$|\sqrt{\mathbb{E}[Z^2]} - \mathbb{E}[Z]| = |\sqrt{\text{trace}(Q)/n} - \mathbb{E}[\|\sqrt{Q}X\|_2/\sqrt{n}]| \leq \frac{2\sqrt{\|Q\|_{\text{op}}}}{\sqrt{n}}.$$

Combining this bound with the tail bound (64), we conclude that

(66)

$$\mathbb{P}\left[\frac{1}{\sqrt{n}}\left|\|\sqrt{Q}X\|_2 - \sqrt{\text{trace}(Q)}\right| > \delta + 2\sqrt{\frac{\|Q\|_{\text{op}}}{n}}\right] \leq 2 \exp\left(-\frac{n\delta^2}{2\|Q\|_{\text{op}}}\right) \quad \text{for all } \delta > 0.$$

Setting $\delta = (t - 2/\sqrt{n})\sqrt{\|Q\|_{\text{op}}}$ in the bound (66) yields that

(67)

$$\mathbb{P}\left[\frac{1}{\sqrt{n}}\left|\|\sqrt{Q}X\|_2 - \sqrt{\text{trace}(Q)}\right| > t\sqrt{\|Q\|_{\text{op}}}\right] \leq 2 \exp\left(-\frac{n(t - 2/\sqrt{n})^2}{2}\right).$$

Similarly, setting $\delta = \sqrt{\|Q\|_{\text{op}}}$ in the tail bound (66) yields that with probability greater than $1 - 2 \exp(-n/2)$, we have

$$(68) \quad \left|\frac{\|Y\|_2}{\sqrt{n}} + \sqrt{\frac{\text{trace}(Q)}{n}}\right| \leq \sqrt{\frac{\text{trace}(Q)}{n}} + 3\sqrt{\|Q\|_{\text{op}}} \leq 4\sqrt{\|Q\|_{\text{op}}}.$$

Using these two bounds, we obtain

$$\left|\frac{\|Y\|_2^2}{n} - \frac{\text{trace}(Q)}{n}\right| = \left|\frac{\|Y\|_2}{\sqrt{n}} - \sqrt{\frac{\text{trace}(Q)}{n}}\right| \left|\frac{\|Y\|_2}{\sqrt{n}} + \sqrt{\frac{\text{trace}(Q)}{n}}\right| \leq 4t\|Q\|_{\text{op}}$$

with the claimed probability. \square

REFERENCES

- [1] J. Abernethy, F. Bach, T. Evgeniou, and J. Stein. Low-rank matrix factorization with attributes. Technical Report Technical Report N-24/06/MM, Ecole des mines de Paris, France, September 2006.
- [2] A. A. Amini and M. J. Wainwright. High-dimensional analysis of semdefinite relaxations for sparse principal component analysis. *Annals of Statistics*, 5B:2877–2921, 2009.
- [3] C. W. Anderson, E. A. Stolz, and S. Shamsunder. Multivariate autoregressive models for classification of spontaneous electroencephalogram during mental tasks. *IEEE Trans. on bio-medical engineering*, 45(3):277, 1998.

- [4] T. W. Anderson. *The statistical analysis of time series*. Wiley Classics Library. John Wiley and Sons, New York, 1971.
- [5] A. Argyriou, T. Evgeniou, and M. Pontil. Multi-task feature learning. In *Neural Information Processing Systems (NIPS)*, Vancouver, Canada, December 2006.
- [6] F. Bach. Consistency of trace norm minimization. *Journal of Machine Learning Research*, 9:1019–1048, June 2008.
- [7] P. Bickel and E. Levina. Covariance estimation by thresholding. *Annals of Statistics*, 36(6):2577–2604, 2008.
- [8] P. Bickel and E. Levina. Regularized estimation of large covariance matrices. *Annals of Statistics*, 36(1):199–227, 2008.
- [9] P. Bickel and B. Li. Regularization in statistics. *TEST*, 15(2):271–344, 2006.
- [10] P. Bickel, Y. Ritov, and A. Tsybakov. Simultaneous analysis of lasso and dantzig selector. *Annals of Statistics*, 37(4):1705–1732, 2009.
- [11] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, Cambridge, UK, 2004.
- [12] E. N. Brown, R. E. Kass, and P. P. Mitra. Multiple neural spike train data analysis: state-of-the-art and future challenges. *Nature Neuroscience*, 7(5), May 2004.
- [13] E. Candès and Y. Plan. Tight oracle bounds for low-rank matrix recovery from a minimal number of random measurements. Technical Report arXiv:1001.0339v1, January 2010.
- [14] E. Candès and T. Tao. Decoding by linear programming. *IEEE Trans. Info Theory*, 51(12):4203–4215, December 2005.
- [15] E. J. Candès and B. Recht. Exact matrix completion via convex optimization. *Found. Comput. Math.*, 9(6):717–772, 2009.
- [16] S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM J. Sci. Computing*, 20(1):33–61, 1998.
- [17] A. Cohen, W. Dahmen, and R. DeVore. Compressed sensing and best k-term approximation. *J. of American Mathematical Society*, 22(1):211–231, July 2008.
- [18] D. Donoho. Compressed sensing. *IEEE Trans. Info. Theory*, 52(4):1289–1306, April 2006.
- [19] N. El-Karoui. Operator norm consistent estimation of large dimensional sparse covariance matrices. *Annals of Statistics*, 36(6):2717–2756, 2008.
- [20] J. Fan and R. Li. Variable selection via non-concave penalized likelihood and its oracle properties. *Jour. Amer. Stat. Ass.*, 96(456):1348–1360, December 2001.
- [21] J. Fan and J. Lv. A selective overview of variable selection in high dimensional feature space. *Statistica Sinica*, 20:101–148, 2010. Invited Review Article.
- [22] M. Fazel. *Matrix Rank Minimization with Applications*. PhD thesis, Stanford, 2002. Available online: <http://faculty.washington.edu/mfazel/thesis-final.pdf>.
- [23] J. Fisher and M. J. Black. Motor cortical decoding using an autoregressive moving average model. *IEEE Engineering in Medicine and Biology Society*, pages 1469–1472, September 2005.
- [24] J. Friedman, T. Hastie, and R. Tibshirani. Sparse inverse covariance estimation with the graphical Lasso. *Biostatistics*, 2007.
- [25] L. Harrison, W. D. Penny, and K. Friston. Multivariate autoregressive modeling of fmri time series. *NeuroImage*, 19:1477–1491, 2003.
- [26] R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, Cambridge, 1985.

- [27] R. A. Horn and C. R. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, Cambridge, 1991.
- [28] J. Huang and T. Zhang. The benefit of group sparsity. Technical Report arXiv:0901.2962, Rutgers University, January 2009.
- [29] S. Ji and J. Ye. An accelerated gradient method for trace norm minimization. In *International Conference on Machine Learning (ICML)*, New York, NY, USA, 2009. ACM.
- [30] I. M. Johnstone. On the distribution of the largest eigenvalue in principal components analysis. *Annals of Statistics*, 29(2):295–327, April 2001.
- [31] R. H. Keshavan, A. Montanari, and S. Oh. Matrix completion from noisy entries. Preprint available at <http://arxiv.org/abs/0906.2027v1>, 2009.
- [32] M. Ledoux. *The Concentration of Measure Phenomenon*. Mathematical Surveys and Monographs. American Mathematical Society, Providence, RI, 2001.
- [33] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: Isoperimetry and Processes*. Springer-Verlag, New York, NY, 1991.
- [34] K. Lee and Y. Bresler. Guaranteed minimum rank approximation from linear observations by nuclear norm minimization with an ellipsoidal constraint. Technical report, UIUC, 2009. Available at arXiv:0903.4742.
- [35] Z. Liu and L. Vandenberghe. Interior-point method for nuclear norm optimization with application to system identification. *SIAM Journal on Matrix Analysis and Applications*, 31(3):1235–1256, 2009.
- [36] K. Lounici, M. Pontil, A. B. Tsybakov, and S. van de Geer. Taking advantage of sparsity in multi-task learning. Technical Report arXiv:0903.1468, ETH Zurich, March 2009.
- [37] H. Lütkepohl. *New introduction to multiple time series analysis*. Springer, New York, 2006.
- [38] P. Massart. *Concentration Inequalities and Model Selection*. Ecole d’Eté de Probabilités, Saint-Flour. Springer, New York, 2003.
- [39] J. Matousek. *Lectures on discrete geometry*. Springer-Verlag, New York, 2002.
- [40] N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the Lasso. *Annals of Statistics*, 34:1436–1462, 2006.
- [41] S. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu. A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers. In *Proceedings of the NIPS Conference*, Vancouver, Canada, December 2009.
- [42] Y. Nesterov. Gradient methods for minimizing composite objective function. Technical Report 2007/76, CORE, Universit’e catholique de Louvain, 2007.
- [43] G. Obozinski, M. J. Wainwright, and M. I. Jordan. Union support recovery in high-dimensional multivariate regression. *Annals of Statistics*, To appear. Presented in part at NIPS 2008 conference.
- [44] D. Paul and I. Johnstone. Augmented sparse principal component analysis for high-dimensional data. Technical report, UC Davis, January 2008.
- [45] G. Raskutti, M. J. Wainwright, and B. Yu. Minimax rates of estimation for high-dimensional linear regression over ℓ_q -balls. Technical Report arXiv:0910.2042, UC Berkeley, Department of Statistics, 2009. Presented in part at Allerton Conference, Sep. 2009.
- [46] P. Ravikumar, M. J. Wainwright, G. Raskutti, and B. Yu. High-dimensional covariance estimation: Convergence rates of ℓ_1 -regularized log-determinant divergence.

- Technical report, Department of Statistics, UC Berkeley, September 2008. Presented in part at NIPS 2008.
- [47] B. Recht, M. Fazel, and P. A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 2007. to appear.
 - [48] B. Recht, W. Xu, and B. Hassibi. Null space conditions and thresholds for rank minimization. Technical report, U. Madison, 2009. Available at <http://pages.cs.wisc.edu/brecht/papers/10.RecXuHas.Thresholds.pdf>.
 - [49] A. Rohde and A. Tsybakov. Estimation of high-dimensional low-rank matrices. Technical Report arXiv:0912.5338v2, January 2010.
 - [50] A. J. Rothman, P. J. Bickel, E. Levina, and J. Zhu. Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics*, 2:494–515, 2009.
 - [51] N. Srebro, J. Rennie, and T. Jaakkola. Maximum-margin matrix factorization. In *Proceedings of the NIPS Conference*, Vancouver, Canada, 2005.
 - [52] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58(1):267–288, 1996.
 - [53] L. Vandenberghe and S. Boyd. Semidefinite programming. *SIAM Review*, 38:49–95, 1996.
 - [54] M. J. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (Lasso). *IEEE Trans. Information Theory*, 55:2183–2202, May 2009.
 - [55] M. Yuan, A. Ekici, Z. Lu, and R. Monteiro. Dimension reduction and coefficient estimation in multivariate linear regression. *Journal Of The Royal Statistical Society Series B*, 69(3):329–346, 2007.
 - [56] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society B*, 1(68):49, 2006.
 - [57] M. Yuan and Y. Lin. Model selection and estimation in the Gaussian graphical model. *Biometrika*, 94(1):19–35, 2007.

DEPARTMENT OF ELECTRICAL ENGINEERING
UNIVERSITY OF CALIFORNIA, BERKELEY
BERKELEY, CA 94720
USA
E-MAIL: sahand_n@eecs.berkeley.edu

DEPARTMENT OF STATISTICS
UNIVERSITY OF CALIFORNIA, BERKELEY
BERKELEY, CA 94720
USA
E-MAIL: wainwrig@stat.berkeley.edu