



Structured Matrix Estimation and Completion

Olga Klopp, Yu Lu, Alexandre B. Tsybakov, Harrison H. Zhou

► To cite this version:

Olga Klopp, Yu Lu, Alexandre B. Tsybakov, Harrison H. Zhou. Structured Matrix Estimation and Completion. 2017. hal-01557745

HAL Id: hal-01557745

<https://hal.science/hal-01557745>

Preprint submitted on 6 Jul 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Structured Matrix Estimation and Completion

Olga Klopp, *ESSEC and CREST**

Yu Lu, *Yale University*[†]

Alexandre B. Tsybakov, *ENSAE, UMR CNRS 9194*[‡]

Harrison H. Zhou, *Yale University*[§]

July 7, 2017

Abstract

We study the problem of matrix estimation and matrix completion under a general framework. This framework includes several important models as special cases such as the gaussian mixture model, mixed membership model, bi-clustering model and dictionary learning. We consider the optimal convergence rates in a minimax sense for estimation of the signal matrix under the Frobenius norm and under the spectral norm. As a consequence of our general result we obtain minimax optimal rates of convergence for various special models.

Keywords: matrix completion, matrix estimation, minimax optimality

AMS 2000 subject classification: 62J99, 62H12, 60B20, 15A83

1 Introduction

Over the past decade, there have been considerable interest in statistical inference for high-dimensional matrices. A fundamental model in this context is the matrix de-noising model, under which one observes a matrix $\theta^* + W$ where θ^* is an unknown non-random $n \times m$ matrix of interest, and W is a random noise matrix. The aim is to estimate θ^* from such observations. Often in applications a part of elements of M is missing. The problem of reconstructing the signal matrix θ^* given partial observations of its entries is known as matrix completion problem. There has been an important research in the past years devoted to accurate matrix completion methods.

In general, the signal θ^* cannot be recovered consistently from noisy and possibly missing observations. If we only know that θ^* is an arbitrary $n \times m$

*kloppolga@math.cnrs.fr

†yu.lu@yale.edu

‡alexandre.tsybakov@ensae.fr

§huibin.zhou@yale.edu

matrix, the guaranteed error of estimating θ^* from noisy observations can be prohibitively high. However, if θ^* has an additional structure one can expect to estimate it with high accuracy from a moderate number of noisy observations. The algorithmic and analytical tractability of the problem depends on the type of adopted structural model. A popular assumption in the matrix completion literature is that the unknown matrix θ^* is of low rank or can be well approximated by a low rank matrix. Significant progresses have been made on low rank matrix estimation and completion problems, see e.g., [9, 8, 19, 26, 31, 32, 16, 28, 7]. However, in several applications, the signal matrix θ^* can have other than just low rank structure. Some examples are as follows.

- *Biology.* The biological data are sometimes expected to have clustering structures. For example, in the gene microarray data, a large number of gene expression levels are measured under different experimental conditions. It has been observed in the experiments that there is a bi-clustering structure on the genes [13]. This means that, besides being of low rank, the gene microarray data can be rearranged to approximately have a block structure.
- *Computer Vision.* To capture higher-level features in natural images, it is common to represent data as a sparse linear combination of basis elements [33] leading to sparse coding models. Unlike the principle component analysis that looks for low rank decompositions, sparse coding learns useful representations with number of basis vectors, which is often greater than the dimension of the data.
- *Networks.* In network models, such as social networks or citation networks, the links between objects are usually governed by the underlying community structures. To capture such structures, several block models have been recently proposed with the purpose of explaining the network data [22, 2, 25].

While there are some successful algorithmic advancements on adapting new structures in these specific applications, not much is known on the fundamental limits of statistical inference for the corresponding models. A few exceptions are the stochastic block model [18, 29] and the bi-clustering model [17]. However, many other structures of signal matrix are not analyzed.

The aim of this paper is to study a general framework of estimating structured matrices. We consider a unified model that includes gaussian mixture model, mixed membership model [2], bi-clustering model [20], and dictionary learning as special cases. We first study the optimal convergence rates in a minimax sense for estimation of the signal matrix under the Frobenius norm and under the spectral norm from complete observations on the sparsity classes of matrices. Then, we investigate this problem in the partial observations regime (structured matrix completion problem) and study the minimax optimal rates under the same norms. We also establish accurate oracle inequalities for the suggested methods.

2 Notation

This section provides a brief summary of the notation used throughout this paper. Let A, B be matrices in $\mathbb{R}^{n \times m}$.

- For a matrix A , A_{ij} is its (i, j) th entry, $A_{\cdot j}$ is its j th column and $A_{i\cdot}$ is its i th row.
- The scalar product of two matrices A, B of the same dimensions is denoted by $\langle A, B \rangle = \text{tr}(A^T B)$.
- We denote by $\|A\|_2$ the Frobenius norm of A and by $\|A\|_\infty$ the largest absolute value of its entries: $\|A\|_\infty = \max_{i,j} |A_{ij}|$. The spectral norm of A is denoted by $\|A\|$.
- For $x \in \mathbb{R}^k$, we denote by $\|x\|_0$ its l_0 -norm (the number of non-zero components of x), and by $\|x\|_q$ its l_q -norm, $1 \leq q \leq \infty$.
- We denote by $\|A\|_{0,\infty}$ the largest l_0 -norm of the rows of $A \in \mathbb{R}^{n \times k}$:

$$\|A\|_{0,\infty} = \max_{1 \leq i \leq n} \|A_{i\cdot}\|_0.$$

- For any $i \in \mathbb{N}$, we write for brevity $[i] = \{1, \dots, i\}$.
- Given a matrix $A = (A_{ij}) \in \mathbb{R}^{n \times m}$, and a set of indices $I \subset [n] \times [m]$, we define the restriction of A on I as a matrix A_I with elements $(A_I)_{ij} = A_{ij}$ if $(i, j) \in I$ and $(A_I)_{ij} = 0$ otherwise.
- The notation $\mathbf{I}_{k \times k}$ and $\mathbf{0}_{k \times l}$ (abbreviated to \mathbf{I} and $\mathbf{0}$ when there is no ambiguity) stands for the $k \times k$ identity matrix and the $k \times l$ matrix with all entries 0, respectively.
- We denote by $|S|$ the cardinality of a finite set S , by $\lfloor x \rfloor$ the integer part of $x \in \mathbb{R}$, and by $\lceil x \rceil$ the smallest integer greater than $x \in \mathbb{R}$.
- We denote by $\mathcal{N}_\epsilon(\mathcal{A})$ the ϵ -covering number, under the Frobenius norm, of a set \mathcal{A} of matrices.

3 General model and examples

Assume that we observe a matrix $Y = (Y_{ij}) \in \mathbb{R}^{n \times m}$ with entries

$$Y_{ij} = E_{ij} (\theta_{ij}^* + \xi_{ij}), \quad i = 1, \dots, n, \quad j = 1, \dots, m, \quad (1)$$

where θ_{ij}^* are the entries of the unknown matrix of interest $\theta^* = (\theta_{ij}^*) \in \mathbb{R}^{n \times m}$, the values ξ_{ij} are independent random variables representing the noise, and E_{ij} are i.i.d. Bernoulli variables with parameter $p \in (0, 1]$ such that (E_{ij}) is independent of (ξ_{ij}) .

Model (1) is called the matrix completion model. Under this model, an entry of matrix θ^* is observed with noise (independently of the other entries) with probability p , and it is not observed with probability $1 - p$. We can equivalently write (1) in the form

$$Y/p = \theta^* + W, \quad (2)$$

where W is a matrix with entries

$$W_{ij} = \theta_{ij}^*(E_{ij} - p)/p + \xi_{ij}E_{ij}/p.$$

The model with complete noisy observations is a special case of (1) (and equivalently of (2)) corresponding to $p = 1$. In this case, $W_{ij} = \xi_{ij}$.

We denote by \mathbb{P}_{θ^*} the probability distribution of Y satisfying (1) and by \mathbb{E}_{θ^*} the corresponding expectation. When there is no ambiguity, we abbreviate \mathbb{P}_{θ^*} and \mathbb{E}_{θ^*} to \mathbb{P} and \mathbb{E} , respectively.

We assume that ξ_{ij} are independent zero mean sub-Gaussian random variables. The sub-Gaussian property means that the following assumption is satisfied.

Assumption 1. *There exists $\sigma > 0$ such that, for all $(i, j) \in [n] \times [m]$,*

$$\forall \lambda \in \mathbb{R}, \quad \mathbb{E} \exp(\lambda \xi_{ij}) \leq \exp(\lambda^2 \sigma^2 / 2).$$

We assume that the signal matrix θ^* is structured, that is, it can be factorized using sparse factors. Specifically, let s_n, k_n, s_m, k_m be integers such $0 \leq s_n \leq k_n$ and $0 \leq s_m \leq k_m$. We assume that

$$\theta^* \in \Theta(s_n, s_m) \subset \mathbb{R}^{n \times m},$$

where

$$\Theta(s_n, s_m) = \{\theta = X B Z^T : X \in \mathcal{A}_{s_n}, B \in \mathbb{R}^{k_n \times k_m} \text{ and } Z \in \mathcal{A}_{s_m}\}.$$

Here, for $s_n = 0$ we assume that $n = k_n$ and the set \mathcal{A}_{s_n} is a set containing only one element, which is the $n \times n$ identity matrix, and for $1 \leq s_n \leq k_n$,

$$\mathcal{A}_{s_n} = \mathcal{A}_{s_n}(n, k_n) = \{A \in \mathcal{D}_n^{n \times k_n}, \|A_i\|_0 \leq s_n, \text{ for all } i \in [n]\} \quad (3)$$

where the set \mathcal{D}_n is a subset of \mathbb{R} called an alphabet. The set \mathcal{A}_{s_m} is defined analogously by replacing n by m . We will also consider the class $\Theta_*(s_n, s_m)$ defined analogously to $\Theta(s_n, s_m)$, with the only difference that the inequality in (3) is replaced by the equality.

Choosing different values of s_n, k_n, s_m, k_m , and different alphabets we obtain several well-known examples of matrix structures.

- Mixture Model:

$$\begin{aligned} \Theta_{MM} = \{ & \theta \in \mathbb{R}^{n \times m} : \theta = X B \text{ for some } B \in \mathbb{R}^{k \times m} \\ & \text{and } X \in \{0, 1\}^{n \times k} \text{ with } \|X_i\|_0 = 1, \forall i \in [n]\}. \end{aligned}$$

- Sparse Dictionary Learning:

$$\Theta_{SDL} = \{\theta = BZ^T \in \mathbb{R}^{d \times n} : B \in \mathbb{R}^{d \times k}, Z \in \mathbb{R}^{n \times k} \text{ with } \|Z_{i\cdot}\|_0 \leq s, \forall i \in [n]\}.$$

- Stochastic Block Model:

$$\Theta_{SBM} = \{\theta = ZBZ^T \in \mathbb{R}^{n \times n} : B \in [0, 1]^{k \times k}, Z \in \{0, 1\}^{n \times k} \text{ with } \|Z_{i\cdot}\|_0 = 1, \forall i \in [n]\}.$$

- Mixed Membership Model:

$$\begin{aligned} \Theta_{MMM} = \{ \theta = ZBZ^T \in \mathbb{R}^{n \times n} : B \in [0, 1]^{k \times k}, Z \in [0, 1]^{n \times k}, \\ \text{with } \|Z_{i\cdot}\|_1 = 1, \|Z_{i\cdot}\|_0 \leq s, \text{ for all } i \in [n] \}. \end{aligned}$$

- Bi-clustering Model:

$$\begin{aligned} \Theta_{Bi} = \{ \theta = XBZ^T \in \mathbb{R}^{n \times m} : B \in [0, 1]^{k_n \times k_m}, X \in \{0, 1\}^{n \times k_n}, Z \in \{0, 1\}^{m \times k_m} \\ \text{with } \|X_{i\cdot}\|_0 = 1, \forall i \in [n], \|Z_{i\cdot}\|_0 = 1, \forall i \in [m] \}. \end{aligned}$$

Here, the classes Θ_{SBM} and Θ_{MMM} are not exactly equal to but rather subclasses of $\Theta_*(1, 1)$ and $\Theta_*(s, s)$, respectively.

Statistical properties of inference methods under the general model (1) are far from being understood. Some results were obtained in particular settings such as the Mixture Model and Stochastic Block Model.

Gaussian mixture models provide a useful framework for several machine learning problems such as clustering, density estimation and classification. There is a quite long history of research on mixtures of Gaussians. We mention only some of this work including methods for estimating mixtures such as pairwise distances [14, 15], spectral methods [37, 23] or the method of moments [12, 5]. Most of these papers are concerned with construction of computationally efficient methods but do not address the issue of statistical optimality. In [3] authors provide precise information theoretic bounds on the clustering accuracy and sample complexity of learning a mixture of two isotropic Gaussians in high dimensions under small mean separation.

The Stochastic Block Model is a useful benchmark for the task of recovering community structure in graph data. More generally, any sufficiently large graph behaves approximately like a stochastic block model for some k , which can be large. The problem of estimation of the probability matrix θ^* in the stochastic block model under the Frobenius norm was considered by several authors [11, 39, 40, 10, 6] but convergence rates obtained there are suboptimal. More recently, minimax optimal rates of estimation were obtained by Gao et al. [18] in the dense case and by Klopp et al [29] in the sparse case.

Recently, a related problem to ours was studied by Soni et al. [35]. These authors consider the case when the matrix to be estimated is the product of two matrices, one of which, called a sparse factor, has a small number of non-zero entries (in contrast to this, we assume row-sparsity). The estimator studied in [35] is a sieve maximum likelihood estimator penalized by the l_0 -norm of the sparse factor where the sieve is chosen as a specific countable set.

4 Results for the case of finite alphabets

We start by considering the case of finite alphabets \mathcal{D}_n and \mathcal{D}_m and complete observations, that is $p = 1$. In this section, we establish the minimax optimal rates of estimation of θ^* under the Frobenius norm and we show that they are attained by the least squares estimator

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \|Y - \theta\|_2^2 \quad (4)$$

where Θ is a suitable class of structured matrices. We first derive an upper bound on the risk of this estimator uniformly over the classes $\Theta = \Theta(s_n, s_m)$. The following theorem provides an oracle inequality for the Frobenius risk of $\hat{\theta}$. Here and in what follows, we adopt the convention that $0 \log \frac{x}{0} = 0$ for any $x > 0$. We also set for brevity

$$d = n + m, \quad r_n = n \wedge k_n, \quad r_m = m \wedge k_m.$$

Theorem 1. *Let Assumption 1 hold, and let $p = 1$. If the sets \mathcal{D}_n and \mathcal{D}_m are finite, there exists a constant $C > 0$ depending only on the cardinalities of \mathcal{D}_n and \mathcal{D}_m such that, for all $\theta^* \in \mathbb{R}^{n \times m}$ and all $\epsilon > 0$, the risk of the estimator (4) satisfies*

$$\mathbb{E}_{\theta^*} \left\{ \|\hat{\theta} - \theta^*\|_2^2 \right\} \leq (1 + \epsilon) \inf_{\bar{\theta} \in \Theta(s_n, s_m)} \|\bar{\theta} - \theta^*\|_2^2 + \frac{C\sigma^2}{\epsilon} (R_X + R_B + R_Z),$$

where $R_X = nr_m \wedge ns_n \log \frac{ek_n}{s_n}$, $R_B = r_n r_m$, $R_Z = mr_n \wedge ms_m \log \frac{ek_m}{s_m}$.

This theorem is proved in Section A.

Note that if the set \mathcal{A}_{s_n} and/or \mathcal{A}_{s_m} in the definition of $\Theta(s_n, s_m)$ contains only the identity matrix, the corresponding term R_X and/or R_Z disappears from the upper bound of Theorem 1.

In Theorem 1, the true signal θ^* can be arbitrary. By assuming that $\theta^* \in \Theta(s_n, s_m)$, we immediately deduce from Theorem 1 that the following bound holds.

Corollary 2. *Under the assumptions of Theorem 1,*

$$\sup_{\theta \in \Theta(s_n, s_m)} \mathbb{E}_{\theta} \left\{ \|\hat{\theta} - \theta\|_2^2 \right\} \leq C\sigma^2 (R_X + R_B + R_Z)$$

for a constant $C > 0$ depending only on the cardinalities of \mathcal{D}_n and \mathcal{D}_m .

The next theorem provides a lower bound showing that the convergence rate of Corollary 2 is minimax optimal. This lower bound is valid for the general matrix completion model (1). In what follows, the notation $\inf_{\hat{\vartheta}}$ stands for the infimum over all estimators $\hat{\vartheta}$ taking values in $\mathbb{R}^{n \times m}$.

Theorem 3. Let the entries W_{ij} of matrix W in model (2) be independent random variables with Gaussian distribution $\mathcal{N}(0, \sigma^2)$, and let the alphabets \mathcal{D}_n and \mathcal{D}_m contain the set $\{0, 1\}$. There exists an absolute constant $C > 0$ such that

$$\inf_{\hat{\vartheta}} \sup_{\theta \in \Theta(s_m, s_n)} \mathbb{P}_\theta \left\{ \|\hat{\vartheta} - \theta\|_2^2 \geq \frac{C\sigma^2}{p} (R_X + R_B + R_Z) \right\} \geq 0.1, \quad (5)$$

and

$$\inf_{\hat{\vartheta}} \sup_{\theta \in \Theta(s_m, s_n)} \mathbb{E}_\theta \|\hat{\vartheta} - \theta\|_2^2 \geq \frac{C\sigma^2}{p} (R_X + R_B + R_Z). \quad (6)$$

Furthermore, the same inequalities hold with $\Theta_*(s_m, s_n)$ in place of $\Theta(s_m, s_n)$ if $s_n \in \{0, 1\}$ and $s_m \in \{0, 1\}$.

The proof of Theorem 3 is given in Section B.1.

The three ingredients R_X, R_B , and R_Z of the optimal rate are coming from the ignorance of X, B and Z respectively. The proof is based on constructing subsets of Θ by fixing two of these parameters to get each of the three terms. The choice of B when fixing the pairs (X, B) and (Z, B) is based on a probabilistic method, namely, Lemma 17. Similar techniques have been used in [29] to prove the lower bounds for sparse graphon estimation, and in [18].

Remark 1. Theorem 3 can be extended to more general sub-Gaussian distributions under an additional Kullback-Leibler divergence assumption. Assume that there is a constant c such that the distribution of Y in model (1) satisfies

$$KL(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) \leq \frac{cp}{2\sigma^2} \|\theta - \theta'\|_2^2.$$

Let the alphabets \mathcal{D}_n and \mathcal{D}_m contain the set $\{0, 1\}$. Then there exists an absolute constant $C > 0$ such that

$$\inf_{\hat{\vartheta}} \sup_{\theta \in \Theta(s_m, s_n)} \mathbb{P}_\theta \left\{ \|\hat{\vartheta} - \theta\|_2^2 \geq \frac{C\sigma^2}{p} (R_X + R_B + R_Z) \right\} \geq 0.1,$$

and

$$\inf_{\hat{\vartheta}} \sup_{\theta \in \Theta(s_m, s_n)} \mathbb{E}_\theta \|\hat{\vartheta} - \theta\|_2^2 \geq \frac{C\sigma^2}{p} (R_X + R_B + R_Z).$$

The proof of this result is similar to that of Theorem 3 and only needs to replace the equality in (25) by inequality. In addition, the lower bounds hold with $\Theta_*(s_m, s_n)$ in place of $\Theta(s_m, s_n)$ if $s_n \in \{0, 1\}$ and $s_m \in \{0, 1\}$.

Remark 2. We summarize the minimax rates for some examples introduced in Section 3 in the following table. The case of Θ_{SBM} is due to [18].

| | |
|----------------|---|
| Θ_{MM} | $\min\{n \log(ek) + km, nm\}$ |
| Θ_{SDL} | $\min\{ns \log(ek/s) + kd, nd\}$ |
| Θ_{SBM} | $n \log(ek) + k^2$ |
| Θ_{Bi} | $\min\{n \log(ek_n) + m \log(ek_m) + k_n k_m, nk_m + m \log(ek_m), mk_n + n \log(ek_n), nm\}$ |

5 Optimal rates in the spectral norm

In this section we derive the optimal rates of convergence of estimators of θ^* when the error is measured in the spectral norm. Interestingly, our results imply that these optimal rates coincide with those obtained for estimation of matrices with no structure. That is, the additional structure that we consider in the present paper does not have any impact on the rate of convergence of the minimax risk when the error is measured in the spectral norm.

The lower bound under the spectral norm can be obtained as a corollary of the lower bound under the Frobenius norm given by Theorem 3.

Corollary 4. *Under the assumptions of Theorem 3, there exists a absolute constant $C' > 0$ such that*

$$\inf_{\hat{\vartheta}} \sup_{\theta \in \Theta(s_m, s_n)} \mathbb{P}_{\theta} \left\{ \|\hat{\vartheta} - \theta\|^2 \geq \frac{C' \sigma^2}{p} (n \vee m) \right\} \geq 0.1,$$

and

$$\inf_{\hat{\vartheta}} \sup_{\theta \in \Theta(s_m, s_n)} \mathbb{E}_{\theta} \|\hat{\vartheta} - \theta\|^2 \geq \frac{C' \sigma^2}{p} (n \vee m).$$

The proof of this corollary is given in Section B.2.

To get matching upper bounds we can use the soft thresholding estimator introduced in [30] or the hard thresholding estimator proposed in [27]. These papers deal with the completion problem for low rank matrices in the context of trace regression model, which is a slightly different setting.

Here, we consider the hard thresholding estimator. Set

$$Y' = Y/p.$$

The singular value decomposition of matrix Y' has the form

$$Y' = \sum_{j=1}^{\text{rank}(Y')} \sigma_j(Y') u_j(Y') v_j(Y')^T, \quad (7)$$

where $\text{rank}(Y')$ is the rank of Y' , $\sigma_j(Y')$ are the singular values of Y' indexed in the decreasing order, and $u_j(Y')$ (respectively, $v_j(Y')$) are the left (respectively,

the right) singular vectors of Y' . The hard thresholding estimator is defined by the formula

$$\tilde{\theta} = \sum_{j: \sigma_j(Y') \geq \lambda} \sigma_j(Y') u_j(Y') v_j(Y')^T \quad (8)$$

where $\lambda > 0$ is the regularization parameter. In this section, we assume that the noise variables W_{ij} are bounded as stated in the next assumption.

Assumption 2. *For all i, j we have $\mathbb{E}(W_{ij}) = 0$, $\mathbb{E}(W_{ij}^2) = \sigma^2$ and there exists a positive constant $b > 0$ such that*

$$\max_{i,j} |W_{ij}| \leq b.$$

A more general case of sub-Gaussian noise can be treated as well; in this case, we can work on the event \mathcal{E}_b where $\|W\|_\infty$ is bounded by a suitable constant b and show that the probability of the complement of \mathcal{E}_b is small.

The following theorem gives the upper bound on the estimation error of the hard thresholding estimator (8).

Theorem 5. *Assume that $\|\theta^*\|_\infty \leq \theta_{\max}$ and let Assumption 2 hold. Let $\lambda = c(b + \theta_{\max})\sqrt{\frac{n \vee m}{p}}$ where $c > 0$ is a sufficiently large absolute constant. Assume that $p \geq \log(n + m)/(n \vee m)$. Then, with \mathbb{P}_{θ^*} probability at least $1 - 2/(n + m)$, the hard thresholding estimator $\tilde{\theta}$ satisfies*

$$\|\tilde{\theta} - \theta^*\|^2 \leq C(b + \theta_{\max})^2 \frac{n \vee m}{p}$$

where $C > 0$ is an absolute constant.

The proof of Theorem 5 is close to the argument in [27]. It is given in Section C.

6 A general oracle inequality under incomplete observations

The aim of this section is to present a general theorem about the behavior of least squares estimators in the setting with incomplete observations. This theorem will be applied in the next section to obtain an analog of the upper bound of Theorem 1 for general alphabets. To state the theorem, it does not matter whether we consider a vector or matrix setting. Therefore, in this section, we will deal with the vector model. Assume that we observe a vector $Y = (Y_1, \dots, Y_N)$ with entries

$$Y_i = E_i(\theta_i^* + \xi_i), \quad i = 1, \dots, N, \quad (9)$$

for some unknown $\theta^* = (\theta_1^*, \dots, \theta_N^*)$. Our goal is to estimate θ^* . Here, ξ_i are independent random noise variables, and E_i are i.i.d. Bernoulli variables with parameter $p \in (0, 1]$ such that (E_1, \dots, E_N) is independent of (ξ_1, \dots, ξ_N) .

When p is known we can equivalently write (9) in the form

$$Y' = \theta^* + W, \quad (10)$$

where now W is a vector with entries

$$W_i = \theta_i^*(E_i - p)/p + \xi_i E_i/p,$$

and $Y' = Y/p$. In this section, we denote by \mathbb{P}_{θ^*} the probability distribution of Y' satisfying (10).

Consider the least squares estimator of θ^* :

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \|Y' - \theta\|_2^2, \quad (11)$$

where Θ is a subset of \mathbb{R}^N . For some element θ_0 of $\arg \min_{\theta \in \Theta} \|\theta - \theta^*\|_2^2$ we set $\Theta_1 = \{\theta \in \Theta : \|\theta - \theta_0\|_2 \leq 1\}$.

Set

$$\epsilon_0 = \frac{1}{2} \left(\inf\{\epsilon \in (0, 1] : N\epsilon^2 > \log \mathcal{N}_\epsilon(\Theta_1)\} + \sup\{\epsilon \in (0, 1] : N\epsilon^2 < \log \mathcal{N}_\epsilon(\Theta_1)\} \right).$$

Since $\mathcal{N}_\epsilon(\Theta_1)$ is a decreasing left-continuous function of $\epsilon \in (0, 1]$, we have

$$\frac{1}{2} \log \mathcal{N}_{\epsilon_0}(\Theta_1) \leq N\epsilon_0^2 \leq \log \mathcal{N}_{\epsilon_0}(\Theta_1). \quad (12)$$

Theorem 6. *Let ξ_i be independent random variables satisfying $\mathbb{E}e^{\lambda \xi_i} \leq e^{\lambda^2 \sigma^2/2}$ for some $\sigma > 0$ and all $\lambda \in \mathbb{R}$. Assume that there exists a constant θ_{\max} such that $\|\theta\|_\infty \leq \theta_{\max}$ for all $\theta \in \Theta$. Then, for any $\theta^* \in \mathbb{R}^N$, with \mathbb{P}_{θ^*} -probability at least $1 - 4/\mathcal{N}_{\epsilon_0}(\Theta_1) - \exp(-pN/6)$, the least squares estimator (11) satisfies the oracle inequality*

$$\|\hat{\theta} - \theta^*\|_2^2 \leq 3 \inf_{\theta \in \Theta} \|\theta - \theta^*\|_2^2 + C \frac{\theta_{\max}^2 + \sigma^2}{p} N\epsilon_0^2,$$

where $C > 0$ is an absolute constant.

The proof of this theorem is given in Section D.

Note that Theorem 6 has no assumption on the true signal θ^* . Using Theorem 6 with $\theta^* \in \Theta$ we immediately deduce that

$$\inf_{\theta \in \Theta} \mathbb{P}_\theta \left(\|\hat{\theta} - \theta\|_2^2 \leq C \frac{\theta_{\max}^2 + \sigma^2}{p} N\epsilon_0^2 \right) \geq 1 - 4/\mathcal{N}_{\epsilon_0}(\Theta_1) - \exp(-pN/6),$$

where $C > 0$ is an absolute constant.

Theorem 6 shows that the rate of convergence of the least squares estimator is determined by the value of ϵ_0 satisfying the global entropy condition (12). This quantity is the critical covering radius that appeared in the literature in different contexts, see, e.g., [41]. In particular, this critical radius has been shown to determine the minimax optimal rates in nonparametric estimation problems. However, it may lead to slightly suboptimal rates (with deterioration by a logarithmic factor) for parametric estimation problems.

7 Structured matrix completion with general alphabets

For the structured matrix completion over infinite alphabets we consider the following parameter spaces:

$$\tilde{\Theta}(s_n, s_m) = \{\theta = XBZ^T : X \in \tilde{\mathcal{A}}_n, \|B\|_\infty \leq B_{max}, Z \in \tilde{\mathcal{A}}_m, \|\theta\|_\infty \leq \theta_{mx}\}.$$

Here, B_{max} and θ_{mx} are positive constants, and for $1 \leq s_n \leq k_n$,

$$\tilde{\mathcal{A}}_n = \{A \in \mathcal{D}_n^{n \times k_n} : \|A_{i\cdot}\|_0 \leq s_n, \text{ for all } i \in [n] \text{ and } \|A\|_\infty \leq 1\}.$$

If $s_n = 0$, we assume that $n = k_n$ and we define $\tilde{\mathcal{A}}_{s_n}$ as the set containing only one element, which is the $n \times n$ identity matrix.

The difference from the class $\Theta(s_n, s_m)$ is only in the fact that the elements of matrix $\theta \in \tilde{\Theta}(s_n, s_m)$ and those of the corresponding factor matrices X, B, Z are assumed to be uniformly bounded. This assumption is natural in many situations, for example, in the Stochastic Block Model or in recommendation systems, where the entries of the matrix are ratings. We introduce the bounds of the entries of the factor matrices in order to fix ambiguities associated with the factorization structure.

A key ingredient in applying Theorem 6 to this particular case is to find the covering number $\log \mathcal{N}_\epsilon(\Theta_1)$ when $\Theta = \tilde{\Theta}(s_n, s_m)$. For any $\Theta \subset \mathbb{R}^{n \times m}$, any $\theta_0 \in \Theta$, and any $u > 0$, set

$$\Theta_u = \{\theta \in \Theta : \|\theta - \theta_0\|_2 \leq u\}.$$

The following result is proved in Section E.

Proposition 7. *For any $\theta_0 \in \tilde{\Theta}(s_n, s_m)$, $0 < \epsilon < 1$, and $u \leq 1$ we have*

$$\log \mathcal{N}_\epsilon(\tilde{\Theta}_u(s_n, s_m)) \leq R_1(\epsilon) \wedge R_2(\epsilon) \wedge R_3(\epsilon) \wedge R_4(\epsilon),$$

where

$$\begin{aligned} R_1(\epsilon) &= ns_n \log \frac{ek_n}{s_n} + ms_m \log \frac{ek_m}{s_m} + (ns_n + ms_m) \log \frac{6B_{max}\sqrt{mn}s_ms_n}{\epsilon} + r_nr_m \log \frac{9u}{\epsilon}, \\ R_2(\epsilon) &= nr_m \log \frac{6u}{\epsilon} + ms_m \log \frac{ek_m}{s_m} + ms_m \log \frac{2B_{max}\sqrt{mn}s_ms_n}{\epsilon}, \\ R_3(\epsilon) &= mr_n \log \frac{6u}{\epsilon} + ns_n \log \frac{ek_n}{s_n} + ns_n \log \frac{2B_{max}\sqrt{mn}s_ms_n}{\epsilon}, \\ R_4(\epsilon) &= mn \log \frac{3u}{\epsilon}. \end{aligned}$$

Theorem 6, together with Proposition 7, imply implies the following upper bound on the estimation error in structured matrix completion.

Corollary 8. Consider model (1). Let Assumption 1 hold. Then, for any $\theta^* \in \mathbb{R}^{n \times m}$, the least squares estimator (4) with $\Theta = \tilde{\Theta}(s_n, s_m)$ satisfies the inequality

$$\|\hat{\theta} - \theta^*\|_2^2 \leq 3 \inf_{\theta \in \tilde{\Theta}(s_n, s_m)} \|\theta - \theta^*\|_2^2 + C \frac{\theta_{\max}^2 + \sigma^2}{p} (R_1(\epsilon_0) \wedge R_2(\epsilon_0) \wedge R_3(\epsilon_0) \wedge R_4(\epsilon_0))$$

with \mathbb{P}_{θ^*} -probability at least

$$1 - \exp(-c(R_1(\epsilon_0) \wedge R_2(\epsilon_0) \wedge R_3(\epsilon_0) \wedge R_4(\epsilon_0))) - \exp(-pmn/18),$$

where, $C, c > 0$ are absolute constants.

Note that Proposition 7 and (12) imply that $\epsilon_0 \geq c' \sqrt{(m+n)/mn}$ for some numerical constant c' . Then, we have that for the general scheme of matrix completion and general alphabets the upper bound given by Corollary 8 departs from the lower bound of Theorem 3 by a logarithmic factor:

Corollary 9. Let the assumptions of Corollary 8 be satisfied. Then the least squares estimator (4) with $\Theta = \tilde{\Theta}(s_n, s_m)$ satisfies the inequality

$$\begin{aligned} \inf_{\theta \in \tilde{\Theta}(s_n, s_m)} \mathbb{P}_\theta \left(\|\hat{\theta} - \theta\|_2^2 \leq C \frac{\theta_{\max}^2 + \sigma^2}{p} [\log(n+m) + \log(s_n s_m)] (R_X + R_B + R_Z) \right) \\ \geq 1 - \exp(-c[\log(n+m) + \log(s_n s_m)] (R_X + R_B + R_Z)) - \exp(-pmn/18) \end{aligned}$$

where $C, c > 0$ are absolute constants.

8 Adaptation to unknown sparsity

The estimators considered above require the knowledge of the degrees of sparsity s_n and s_m of θ^* . In this section, we suggest a method that does not require such a knowledge and thus it is adaptive to the unknown degree of sparsity. Our approach will be to estimate θ^* using a sparsity penalized least squares estimator. Let

$$\mathcal{X} = \cup_{s_n=1}^{k_n} \cup_{s_m=1}^{k_m} \tilde{\Theta}(s_n, s_m) \quad (13)$$

and set

$$\begin{aligned} R(s_n, s_m) = [nr_m \log(6\sqrt{n \wedge m})] \wedge [ns_n \log(k_n s_m(n \wedge m))] \\ + [mr_n \log(6\sqrt{n \wedge m})] \wedge [ms_m \log(k_n s_m(n \wedge m))] \\ + r_n r_m \log(9\sqrt{n \wedge m}). \end{aligned} \quad (14)$$

For any $\theta = XBZ^T \in \mathcal{X}$ let

$$R(\theta) = R(\|X\|_{0,\infty}, \|Z\|_{0,\infty}). \quad (15)$$

In the following, Ω denotes the random set of observed indices (i, j) in model (1). In this section we denote by $\hat{\theta}$ the following estimator

$$\hat{\theta} \in \arg \min_{\theta = XBZ^T \in \mathcal{X}} \{ \|Y - \theta_\Omega\|_2^2 + \lambda R(\theta) \} \quad (16)$$

where $\lambda > 0$ is a regularization parameter. Note that this estimator does not require the knowledge of p . The following theorem proved in Appendix F gives an upper bound on the estimation error of $\hat{\theta}$.

Theorem 10. *Assume that $nm \log(3\sqrt{n \wedge m}) \geq 6 \log(k_n k_m)$ and $d \geq 10$. Let $\lambda = 8(\sigma \vee \theta_{\max})^2$. Then, for any $\theta^* \in \mathbb{R}^{n \times m}$, with \mathbb{P}_{θ^*} -probability at least $1 - 5 \exp(-d/10) - 2 \exp(-pnm)$ the estimator (16) satisfies*

$$\|\hat{\theta} - \theta^*\|_2^2 \leq C \inf_{\theta \in \mathcal{X}} \left\{ \|\theta - \theta^*\|_2^2 + \frac{(\sigma \vee \theta_{\max})^2}{p} R(\theta) \right\}$$

where $C > 0$ is an absolute constant.

Theorem 10 implies that for the general scheme of matrix completion and general alphabets we obtain the following upper bound which departs from the lower bound of Theorem 3 by a logarithmic factor:

$$\begin{aligned} \inf_{\theta \in \mathcal{X}} \mathbb{P}_\theta \left(\|\hat{\theta} - \theta\|_2^2 \leq C \frac{(\sigma \vee \theta_{\max})^2}{p} [\log(n \wedge m) + \log(s_n s_m)] (R_X + R_B + R_Z) \right) \\ \geq 1 - 5 \exp(-d/6) - 2 \exp(-pnm), \end{aligned}$$

where $C > 0$ is an absolute constant.

We finish this section by two remarks.

1. *Structured matrix estimation.* In the case of complete observations, that is $p = 1$, the estimator (16) coincides with the following estimator

$$\hat{\theta} \in \arg \min_{\theta = XBZ^T \in \mathcal{X}} \{ \|Y - \theta\|_2^2 + \lambda R(\theta) \}. \quad (17)$$

Then, one can show that, with high probability, the following upper bound on the estimation error holds

$$\|\hat{\theta} - \theta^*\|_2^2 \leq C \sigma^2 [\log(n \wedge m) + \log(s_n s_m)] (R_X + R_B + R_Z).$$

Here we do not need an upper bound on $\|\theta^*\|_\infty$. At the same time, the estimator (17) is adaptive to the sparsity parameter (s_n, s_m) .

2. *Sparse Factor Model.* Sparse Factor Model is studied in [35]. With our notation, it corresponds to a particular case of $n = k_n$ and X being the identity matrix with the difference that we consider row-sparse matrix Z while Z is assumed component-wise sparse in [35]. Convergence rates obtained in [35] are of the order $p^{-1}(nk_m + k_m m)$ (up to a logarithmic factor). This is greater than the upper bound given by Theorem 10 which, in this setting, is of the order $p^{-1}[n(k_m \wedge m) + s_m m]$.

9 Proofs

A Proof of Theorem 1

Set

$$\begin{aligned}\bar{R}_1(n, m) &= ns_n \log \left(\frac{ek_n |\mathcal{D}_n|}{s_n} \right) + r_n m, \\ \bar{R}_2(n, m) &= ns_n \log \left(\frac{ek_n |\mathcal{D}_n|}{s_n} \right) + ms_m \log \left(\frac{ek_m |\mathcal{D}_m|}{s_m} \right) + r_n r_m.\end{aligned}$$

Since $\hat{\theta}$ is the least squares estimator on $\Theta(s_n, s_m)$, and $Y = \theta^* + W$, we have that for any $\bar{\theta} \in \Theta(s_n, s_m)$,

$$\|\hat{\theta} - \theta^*\|_2^2 \leq \|\bar{\theta} - \theta^*\|_2^2 + 2\langle \hat{\theta} - \bar{\theta}, W \rangle. \quad (18)$$

Now we use the following lemma proved in Section G.1.

Lemma 11. *Let $W \in \mathbb{R}^{n \times m}$ be a random matrix with independent σ -sub-Gaussian entries. Introduce the notation*

$$U_{\bar{\theta}}^* = \sup_{\theta \in \Theta(s_n, s_m), \theta \neq \bar{\theta}} \frac{\langle \theta - \bar{\theta}, W \rangle^2}{\|\theta - \bar{\theta}\|_2^2}.$$

For any $t > 0$, the following inequalities hold, where $C > 0$ is an absolute constant:

(i)

$$\sup_{\bar{\theta} \in \Theta(s_n, s_m)} \mathbb{P} \{ U_{\bar{\theta}}^* \geq 3\sigma^2 (\bar{R}_1(n, m) + t) \} \leq e^{-t}, \quad \sup_{\bar{\theta} \in \Theta(s_n, s_m)} \mathbb{E}(U_{\bar{\theta}}^*) \leq C\sigma^2 \bar{R}_1(n, m),$$

(ii)

$$\sup_{\bar{\theta} \in \Theta(s_n, s_m)} \mathbb{P} \{ U_{\bar{\theta}}^* \geq 3\sigma^2 (\bar{R}_2(n, m) + t) \} \leq e^{-t}, \quad \sup_{\bar{\theta} \in \Theta(s_n, s_m)} \mathbb{E}(U_{\bar{\theta}}^*) \leq C\sigma^2 \bar{R}_2(n, m),$$

(iii)

$$\sup_{\bar{\theta} \in \Theta(s_n, s_m)} \mathbb{P} \{ U_{\bar{\theta}}^* \geq 3\sigma^2 (nm + t) \} \leq e^{-t}, \quad \sup_{\bar{\theta} \in \Theta(s_n, s_m)} \mathbb{E}(U_{\bar{\theta}}^*) \leq C\sigma^2 nm.$$

Applying Lemma 11 (i) to (18), for any $\epsilon > 0$, we get that

$$\mathbb{E} \left\{ \|\hat{\theta} - \theta^*\|_2^2 \right\} \leq (1 + \epsilon) \|\theta^* - \bar{\theta}\|_2^2 + \frac{C\sigma^2}{\epsilon} \bar{R}_1(n, m). \quad (19)$$

On the other hand, applying Lemma 11 (i) to $\langle \hat{\theta} - \bar{\theta}, W \rangle = \langle (\hat{\theta} - \bar{\theta})^T, W^T \rangle$ with n replaced by m , for all $\epsilon > 0$, we get

$$\mathbb{E} \left\{ \|\hat{\theta} - \theta^*\|_2^2 \right\} \leq (1 + \epsilon) \|\theta^* - \bar{\theta}\|_2^2 + \frac{C\sigma^2}{\epsilon} \bar{R}_1(m, n). \quad (20)$$

Finally using Lemma 11 (ii) and (iii) we get

$$\mathbb{E} \left\{ \|\hat{\theta} - \theta^*\|_2^2 \right\} \leq (1 + \epsilon) \|\theta^* - \bar{\theta}\|_2^2 + \frac{C\sigma^2}{\epsilon} \bar{R}_2(n, m) \quad (21)$$

and

$$\mathbb{E} \left\{ \|\hat{\theta} - \theta^*\|_2^2 \right\} \leq (1 + \epsilon) \|\theta^* - \bar{\theta}\|_2^2 + \frac{C\sigma^2}{\epsilon} nm. \quad (22)$$

Inequalities (19) - (22) imply that for all $\epsilon > 0$ and all $\bar{\theta} \in \Theta(s_n, s_m)$

$$\mathbb{E} \left\{ \|\hat{\theta} - \theta^*\|_2^2 \right\} \leq (1 + \epsilon) \|\theta^* - \bar{\theta}\|_2^2 + \frac{C\sigma^2}{\epsilon} \bar{R}_3(n, m).$$

where $\bar{R}_3(n, m) = \min\{\bar{R}_1(n, m), \bar{R}_1(m, n), \bar{R}_2(n, m), nm\}$. Taking the supremum over $\bar{\theta} \in \Theta(s_n, s_m)$ and simplifying the expression for $\bar{R}_3(n, m)$ we obtain the result of Theorem 1.

B Proof of Lower Bounds

B.1 Proof of Theorem 3

Lower bound with the terms R_X and R_Z . We only prove the lower bound with the term R_Z by fixing $X = X_0$ and $B = B_0$, where X_0 and B_0 are matrices specified below. The bound with R_X is analogous. Fix

$$X_0 = \begin{cases} [\mathbf{I}_{k_n \times k_n}, \mathbf{0}]^T, & \text{if } n \geq k_n, \\ [\mathbf{I}_{n \times n}, \mathbf{0}], & \text{otherwise.} \end{cases}$$

By Lemma 16, for $k_m \geq 2$, we can find $S_0 \subseteq \{0, 1\}^{k_m}$ with the following properties:

- (i) $\log |S_0| \geq c_1^* s_m \log \frac{ek_m}{s_m}$,
- (ii) $c_2^* s_m \leq \|a\|_0 \leq s_m$ for all $a \in S_0$, and $\|a\|_0 = s_m$ for all $a \in S_0$ if $s_m \leq k_m/2$,
- (iii) $\|a - b\|_2^2 \geq c_3^* s_m$ for all $a, b \in S_0$ such that $a \neq b$,

where $c_j^* > 0$, $j = 1, 2, 3$, are absolute constants.

Assume first that $k_m \geq 2$ and $\min\{\frac{r_n}{97}, c_1^* s_m \log \frac{ek_m}{s_m}\} \geq \log 8$. Then, choose an arbitrary subset $\mathcal{S} \subseteq S_0$ of cardinality $|\mathcal{S}| = \lfloor \exp\left(\min\{\frac{r_n}{97}, c_1^* s_m \log \frac{ek_m}{s_m}\}\right) \rfloor$ where $r_n = n \wedge k_n$ and we denote by $\lfloor x \rfloor$ the integer part of x . Since $\log |\mathcal{S}| \leq r_n/96$, Lemma 17 implies that there exists a matrix $Q \in \{-1, 1\}^{r_n \times k_m}$ such that, for any $a, b \in \mathcal{S}$,

$$\frac{r_n}{2} \|a - b\|_2^2 \leq \|Qa - Qb\|_2^2 \leq \frac{3r_n}{2} \|a - b\|_2^2. \quad (23)$$

For this Q , let

$$B_0 = [\delta Q, \mathbf{0}_{(k_n - r_n) \times k_m}]^T$$

with $\delta > 0$ to be specified below. Define $\mathcal{Z} = \{Z \in \{0, 1\}^{m \times k_m}, Z_{i \cdot} \in \mathcal{S} \text{ for all } i \in [m]\}$ and $T_Z = \{\theta = X_0 B_0 Z^T, Z \in \mathcal{Z}\}$. We have $T_Z \subseteq \Theta(s_n, s_m)$ and $\log |T_Z| = \log |\mathcal{Z}| = m \log |\mathcal{S}|$.

For any matrices $\theta = X_0 B_0 Z^T \in T_Z$, and $\bar{\theta} = X_0 B_0 \bar{Z}^T \in T_Z$ we have

$$\|\theta - \bar{\theta}\|_2^2 = \delta^2 \|QZ^T - Q\bar{Z}^T\|_2^2 = \delta^2 \sum_{i=1}^m \|QZ_{i \cdot}^T - Q\bar{Z}_{i \cdot}^T\|_2^2.$$

Using (23) and property (iii) of S_0 , we find

$$\|\theta - \bar{\theta}\|_2^2 \geq \frac{r_n \delta^2}{2} \sum_{i=1}^m \|Z_{i \cdot} - \bar{Z}_{i \cdot}\|_2^2 \geq \frac{c_3^* r_n \delta^2 m s_m}{2}. \quad (24)$$

On the other hand, Lemma 15 together with (23) implies that the Kullback-Leibler divergence between \mathbb{P}_θ and $\mathbb{P}_{\bar{\theta}}$ satisfies

$$KL(\mathbb{P}_\theta, \mathbb{P}_{\bar{\theta}}) = \frac{p}{2\sigma^2} \|\theta - \bar{\theta}\|_2^2 \leq \frac{3pr_n \delta^2}{4\sigma^2} \sum_{i=1}^m \|Z_{i \cdot} - \bar{Z}_{i \cdot}\|_2^2 \leq \frac{3pr_n \delta^2 m s_m}{2\sigma^2}. \quad (25)$$

If we choose now $\delta^2 = \frac{C_0 \sigma^2}{pr_n s_m} \log |\mathcal{S}|$ for some absolute constant $C_0 > 0$ small enough, then (24), (25), Theorem 2.5 in [36] and the fact that $|\mathcal{S}| \geq 8$ imply that

$$\inf_{\hat{\theta}} \sup_{\theta \in T_Z} \mathbb{P}_\theta \left\{ \|\hat{\theta} - \theta\|_2^2 \geq C_1 \frac{\sigma^2}{p} \left(mr_n \wedge ms_m \log \frac{ek_m}{s_m} \right) \right\} \geq 0.7 \quad (26)$$

for some absolute constant $C_1 > 0$. This yields the term of the lower bound containing R_Z in the case when $k_m \geq 2$ and $\min\{\frac{r_n}{97}, c_1^* s_m \log \frac{ek_m}{s_m}\} \geq \log 8$. In the complementary case, when $k_m = 1$ or $\min\{\frac{r_n}{97}, c_1^* s_m \log \frac{ek_m}{s_m}\} < \log 8$, the value R_Z is smaller than $C\sigma^2 m/p$ for an absolute constant $C > 0$. Thus, in this case, it suffices to prove the lower bound of order m/p . To do this, let the matrices X_0 and B_0 be such that their (1,1)th entry is equal to 1 and all other entries are 0, and consider the set of matrices Z such that their first column is a binary vector in $\{0, 1\}^m$ and all other columns are 0. This defines a set of matrices $\theta = X_0 B_0 Z^T$ contained in $\Theta(s_n, s_m)$, which is isometric, under the Frobenius norm, to the set of binary vectors $\{0, 1\}^m$ equipped with the Euclidean norm. Therefore, a lower bound of order m/p follows in a standard way as for vector estimation problem. We omit further details.

Analogously, by permuting n and m , we obtain that

$$\inf_{\hat{\theta}} \sup_{\theta \in \Theta(s_n, s_m)} \mathbb{P}_\theta \left\{ \|\hat{\theta} - \theta\|_2^2 \geq C_1 \frac{\sigma^2}{p} \left(nr_m \wedge ns_n \log \frac{ek_n}{s_n} \right) \right\} \geq 0.7, \quad (27)$$

which yields the term of the lower bound containing R_X .

Lower bound with the term R_B . To obtain the term containing R_B in the lower bound (5), we fix $X = X_0$ and $Z = Z_0$ where

$$X_0 = \begin{cases} [\mathbf{I}_{k_n \times k_n}, \mathbf{0}]^T, & \text{if } n \geq k_n, \\ [\mathbf{I}_{n \times n}, \mathbf{0}], & \text{otherwise,} \end{cases} \quad \text{and} \quad Z_0 = \begin{cases} [\mathbf{I}_{k_m \times k_m}, \mathbf{0}]^T, & \text{if } m \geq k_m, \\ [\mathbf{I}_{m \times m}, \mathbf{0}], & \text{otherwise.} \end{cases}$$

We first note that if $r_n r_m < 16$, the lower bound with term R_B is trivially obtained by distinguishing between two matrices. For $r_n r_m \geq 16$, by vectorizing a $r_n \times r_m$ matrix into a $r_n r_m$ dimensional vector, and applying the Varshamov-Gilbert bound [36, Lemma 2.9] we obtain that there exists a subset $\mathcal{B} \subseteq \{0, 1\}^{r_n \times r_m}$ such that for any $Q, \bar{Q} \in \mathcal{B}$,

$$\|Q - \bar{Q}\|_2^2 = \sum_{i,j} \mathbf{1}\{Q_{ij} \neq \bar{Q}_{ij}\} \geq \frac{r_n r_m}{8}$$

and $\log |\mathcal{B}| \geq \frac{r_n r_m}{8}$. We define

$$T_B = \left\{ \theta = X_0 B Z_0^T, B = \delta \begin{bmatrix} Q & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, Q \in \mathcal{B} \right\}.$$

Clearly, $T_B \subseteq \Theta(s_n, s_m)$. For any $\theta = X_0 B Z_0^T \in T_B, \bar{\theta} = X_0 \bar{B} Z_0^T \in T_B$, we have

$$\|\theta - \bar{\theta}\|_2^2 = \|B - \bar{B}\|_2^2 = \delta^2 \|Q - \bar{Q}\|_2^2 \geq \frac{r_n r_m \delta^2}{8}.$$

Lemma 15 implies that $KL(\mathbb{P}_\theta, \mathbb{P}_{\bar{\theta}}) = \frac{p}{2\sigma^2} \|\theta - \bar{\theta}\|_2^2 \leq \frac{p\delta^2 r_n r_m}{2\sigma^2}$. Choosing $\delta^2 = C'_0 \sigma^2 / p$ for some constant $C'_0 > 0$ small enough and using Theorem 2.5 in [36] we obtain

$$\inf_{\hat{\theta}} \sup_{\theta \in T_B} \mathbb{P} \left\{ \|\hat{\theta} - \theta\|_2^2 \geq C_2 \frac{\sigma^2}{p} r_n r_m \right\} \geq 0.7 \quad (28)$$

for some absolute constant $C_2 > 0$.

Combining (26), (27) and (28) proves the lower bound (5). The bound (6) follows from (5) and Markov's inequality.

Finally, the lower bounds for the classes $\Theta_*(s_n, s_m)$ with $s_n, s_m \in \{0, 1\}$ are proved analogously. It suffices to note that, if $s_m = 0$, there is no matrix Z in the definition of the class and thus there is no term R_Z . If $s_m = 1$ we follow the above argument corresponding to R_Z with the only difference that, by (ii) of Lemma 16, we can grant the exact equality $\|a\|_0 = 1$ for all $a \in S_0$ and $k_m \geq 2$. We omit further details.

B.2 Proof of Corollary 4

Define

$$\Gamma_2 = \{\theta = X B Z^T : X \in \mathcal{A}_{s_n}(n, 2), B \in \mathbb{R}^{2 \times 2} \text{ and } Z \in \mathcal{A}_{s_m}(m, 2)\}.$$

For any $\theta = XBZ^T \in \Gamma_2$, let

$$\tilde{X} = [X, \mathbf{0}_{n \times (k_n - 2)}], \tilde{Z} = [Z, \mathbf{0}_{m \times (k_m - 2)}], \tilde{B} = \begin{bmatrix} B & \mathbf{0}_{2 \times (k_m - 2)} \\ \mathbf{0}_{(k_n - 2) \times 2} & \mathbf{0} \end{bmatrix}.$$

We have $\theta = XBZ^T = \tilde{X}\tilde{B}\tilde{Z}^T \in \Theta(s_n, s_m)$, which implies that $\Gamma_2 \subseteq \Theta(s_n, s_m)$. Thus, for any $t > 0$, and any estimator $\hat{\vartheta} \in \mathbb{R}^{n \times m}$ we have

$$\sup_{\theta \in \Theta(s_n, s_m)} \mathbb{P}_\theta(\|\hat{\vartheta} - \theta\|^2 \geq t) \geq \sup_{\theta \in \Gamma_2} \mathbb{P}_\theta(\|\hat{\vartheta} - \theta\|^2 \geq t). \quad (29)$$

For an estimator $\hat{\vartheta} \in \mathbb{R}^{n \times m}$, let $\hat{\vartheta}_2 \in \mathbb{R}^{n \times m}$ be the closest matrix to $\hat{\vartheta}$ in the Frobenius norm among all matrices of rank at most 2. Since $\theta \in \Gamma_2$ is of rank at most 2, we have $\|\hat{\vartheta} - \hat{\vartheta}_2\| \leq \|\hat{\vartheta} - \theta\|$ and

$$\|\theta - \hat{\vartheta}_2\|_2^2 \leq 4\|\theta - \hat{\vartheta}_2\|^2 \leq 8\left(\|\theta - \hat{\vartheta}\|^2 + \|\hat{\vartheta} - \hat{\vartheta}_2\|^2\right) \leq 16\|\theta - \hat{\vartheta}\|^2.$$

Thus,

$$\mathbb{P}_\theta(\|\hat{\vartheta} - \theta\|^2 \geq t) \geq \mathbb{P}_\theta(\|\hat{\vartheta}_2 - \theta\|_2^2 \geq 16t)$$

for any $\theta \in \Gamma_2$ and any estimator $\hat{\vartheta}$. The last inequality and (29) imply

$$\inf_{\hat{\vartheta}} \sup_{\theta \in \Theta(s_n, s_m)} \mathbb{P}_\theta(\|\hat{\vartheta} - \theta\|^2 \geq t) \geq \inf_{\hat{\vartheta}} \sup_{\theta \in \Gamma_2} \mathbb{P}_\theta(\|\hat{\vartheta} - \theta\|_2^2 \geq 16t).$$

The result of Corollary 4 follows now by choosing $t = C_3 \frac{\sigma^2}{p}(n + m)$ for some constant $C_3 > 0$ and using Theorem 3 with $k_m = k_n = 2$.

C Proof of Theorem 5

Note that $Y' - \theta^* = Y/p - \theta^* = W$. It is straightforward to see that if $\lambda \geq \|W\|$, then

$$\|\tilde{\theta} - \theta^*\| \leq 2\lambda.$$

Thus, to prove the theorem, it suffices to show that for $\lambda = c(b + \theta_{\max})\sqrt{\frac{n \vee m}{p}}$ with $c > 0$ large enough we have $\lambda \geq \|W\|$ with high probability.

Since $W_{ij} = \theta_{ij}(E_{ij} - p)/p + \xi_{ij}E_{ij}/p$, we have

$$\|W\| \leq p^{-1}(\|\Sigma_1\| + \|\Sigma_2\|) \quad (30)$$

where $\Sigma_1 \in \mathbb{R}^{n \times m}$ is a matrix with entries $\xi_{ij}E_{ij}$ and $\Sigma_2 \in \mathbb{R}^{n \times m}$ is a matrix with entries $\theta_{ij}(E_{ij} - p)$. The second term in (30) is controlled using the following bound on the spectral norms of random matrices.

Proposition 12 ([4]). *Let A be an $n \times m$ matrix whose entries A_{ij} are independent centered bounded random variables. Then, for any $0 < \epsilon \leq 1/2$ there exists an absolute constant c_ϵ depending only on ϵ such that, for every $t > 0$,*

$$\mathbb{P}\left\{\|A\| \geq (1 + \epsilon)2\sqrt{2}(\sigma_1 \vee \sigma_2) + t\right\} \leq (n \wedge m) \exp\left(-\frac{t^2}{c_\epsilon \sigma_*^2}\right)$$

where

$$\sigma_1 = \max_i \sqrt{\sum_j \mathbb{E}[A_{ij}^2]}, \quad \sigma_2 = \max_j \sqrt{\sum_i \mathbb{E}[A_{ij}^2]}, \quad \sigma_* = \max_{ij} |A_{ij}|.$$

We now apply Proposition 12 with $A_{ij} = (E_{ij} - p)\theta_{ij}$. Then

$$\sigma_1 \leq \theta_{\max} \sqrt{np}, \quad \sigma_2 \leq \theta_{\max} \sqrt{mp} \quad \text{and} \quad \sigma_* \leq \theta_{\max}.$$

Using these bounds and taking in Proposition 12 the values $\epsilon = 1/2$ and $t = \sqrt{c_{1/2}} \theta_{\max} \log(n+m)$ we obtain that there exists an absolute constant $c^* > 0$ such that

$$\|\Sigma_2\| \leq 3\theta_{\max} \sqrt{2(n \vee m)p} + c^* \theta_{\max} \sqrt{2 \log(n+m)}$$

with probability at least $1 - 1/(n+m)$. Similarly, there exists an absolute constants $c^* > 0$ such that, with probability at least $1 - 1/(n+m)$,

$$\|\Sigma_1\| \leq 3\sigma \sqrt{2(n \vee m)p} + c^* b \sqrt{2 \log(n+m)}.$$

Using these remarks, the assumption $p \geq \log(n+m)/(n \vee m)$, and (30) we obtain that the choice $\lambda = c(b + \theta_{\max}) \sqrt{\frac{n \vee m}{p}}$ with $c > 0$ large enough implies the inequality $\lambda \geq \|W\|$ with probability at least $1 - 2/(n+m)$.

D Proof of Theorem 6

Let $R = \sqrt{\frac{3(\theta_{\max}^2 + \sigma^2)}{p}} N \epsilon_0^2$. We consider two cases separately.

Case 1: $\hat{\theta} \in \mathcal{G} \triangleq \{\theta \in \Theta, \|\theta - \theta_0\|_2 \leq 2^* R\}$. Then the desired result follows from the fact that $\|\hat{\theta} - \theta_0\|_2 \leq 2^* R$ and from the inequality

$$\|\hat{\theta} - \theta^*\|_2^2 \leq 2\|\theta_0 - \theta^*\|_2^2 + 2\|\hat{\theta} - \theta_0\|_2^2.$$

Case 2: $\hat{\theta} \notin \mathcal{G}$. The definition of the least squares estimator (11) implies $\|Y' - \hat{\theta}\|_2^2 \leq \|Y' - \theta_0\|_2^2$. Writing Y' as $\theta^* + W$ and rearranging, we obtain

$$\|\hat{\theta} - \theta^*\|_2^2 \leq \|\theta_0 - \theta^*\|_2^2 + 2\langle \hat{\theta} - \theta_0, W \rangle.$$

Since $\hat{\theta} \in \mathcal{G}^c$, Lemma 18 yields

$$\langle \hat{\theta} - \theta_0, W \rangle \leq \frac{1}{8} \|\theta - \theta_0\|_2^2 + 96R^2 \tag{31}$$

with probability greater than $1 - 4 \exp(-\alpha R^2/2) - 2 \exp(-pN/6)$ where $\alpha = \frac{p}{6(\theta_{\max}^2 + \sigma^2)}$. On the event where (31) holds,

$$\begin{aligned} \|\hat{\theta} - \theta^*\|_2^2 &\leq \|\theta_0 - \theta^*\|_2^2 + \frac{1}{4} \|\theta - \theta_0\|_2^2 + 192R^2 \\ &\leq \frac{3}{2} \|\theta_0 - \theta^*\|_2^2 + \frac{1}{2} \|\hat{\theta} - \theta^*\|_2^2 + 192R^2. \end{aligned}$$

This yields $\|\hat{\theta} - \theta^*\|_2^2 \leq 3\|\theta_0 - \theta^*\|_2^2 + 384R^2$ with probability greater than $1 - 4\exp(-\alpha R^2/2) - \exp(-pN/6)$.

E Proof of Proposition 7

We start by proving the upper bound corresponding to $R_1(\epsilon)$. Let $\tilde{\mathcal{A}}_n(\delta)$ denote the δ -covering set of $\tilde{\mathcal{A}}_n$ under the ℓ_∞ norm. For any given $\theta = XBZ^T \in \tilde{\Theta}_u(s_n, s_m)$, there exist $X_0 \in \tilde{\mathcal{A}}_n(\delta_1)$ and $Z_0 \in \tilde{\mathcal{A}}_m(\delta_2)$ such that $\|X - X_0\|_\infty \leq \delta_1$ and $\|Z - Z_0\|_\infty \leq \delta_2$. For such X_0 and Z_0 , let $\overline{\mathcal{T}}(X_0, Z_0)$ be the $\epsilon/3$ -covering set of $\mathcal{T}_u(X_0, Z_0)$ defined in Lemma 21. For any $B \in \mathcal{T}_u(X_0, Z_0)$, there exists $B_0 \in \overline{\mathcal{T}}(X_0, Z_0)$ such that

$$\begin{aligned} & \|XBZ^T - X_0B_0Z_0^T\|_2 \\ & \leq \|XBZ^T - X_0BZ^T\|_2 + \|X_0BZ^T - X_0BZ_0^T\|_2 + \|X_0BZ_0^T - X_0B_0Z_0^T\|_2 \\ & \leq 2s_n\sqrt{nm}\|X - X_0\|_\infty\|BZ\|_\infty + 2s_m\sqrt{nm}\|X_0B\|_\infty\|Z - Z_0\|_\infty + \frac{\epsilon}{3} \\ & \leq 2B_{max}\sqrt{nm}s_ns_m(\delta_1 + \delta_2) + \frac{\epsilon}{3}, \end{aligned}$$

where in the second inequality we have used Lemma 23 and the last inequality is due to the assumptions that $\|B\|_\infty \leq B_{max}$, $\|X\|_\infty \leq 1$ and $\|Z\|_\infty \leq 1$. Choosing $\delta_1 = \delta_2 = \epsilon/(6B_{max}\sqrt{nm}s_ns_m)$ we get that the set

$$\overline{\Theta}_u := \bigcup_{X_0 \in \tilde{\mathcal{A}}_n(\delta_1), Z_0 \in \tilde{\mathcal{A}}_m(\delta_2)} \overline{\mathcal{T}}(X_0, Z_0)$$

is an ϵ -covering set of $\tilde{\Theta}_u(s_n, s_m)$. Then, Lemmas 21 and 22 imply

$$\begin{aligned} \log \mathcal{N}_\epsilon(\tilde{\Theta}_u(s_n, s_m)) & \leq \log |\tilde{\mathcal{A}}_n(\delta_1)| + \log |\tilde{\mathcal{A}}_m(\delta_2)| + \max_{X_0, Z_0} \log |\overline{\mathcal{T}}(X_0, Z_0)| \\ & \leq ns_n \log \frac{ek_n}{s_n} + ns_n \log \frac{6B_{max}\sqrt{mn}s_ms_n}{\epsilon} + r_nr_m \log \frac{9u}{\epsilon} \\ & \quad + ms_m \log \frac{ek_m}{s_m} + ms_m \log \frac{6B_{max}\sqrt{mn}s_ms_n}{\epsilon}. \end{aligned}$$

To get upper bounds corresponding to $R_2(\epsilon)$, $R_3(\epsilon)$ and $R_4(\epsilon)$ we define

$$\Theta_u^1 = \left\{ \theta = AZ^T, A \in [-s_n B_{max}, s_n B_{max}]^{n \times k_m}, Z \in \tilde{\mathcal{A}}_m, \|\theta - \theta_0\|_2 \leq u \right\},$$

$$\Theta_u^2 = \left\{ \theta = XG^T, X \in \tilde{\mathcal{A}}_n, G \in [-s_m B_{max}, s_m B_{max}]^{k_n \times m}, \|\theta - \theta_0\|_2 \leq u \right\},$$

and

$$\Theta_u^3 = \left\{ \theta \in [-s_ns_mB_{max}, s_ns_mB_{max}]^{n \times m}, \|\theta - \theta_0\|_2 \leq u \right\}.$$

It is easy to verify that $\tilde{\Theta}_u(s_n, s_m) \subseteq \Theta_u^1$, $\tilde{\Theta}_u(s_n, s_m) \subseteq \Theta_u^2$ and $\tilde{\Theta}_u(s_n, s_m) \subseteq \Theta_u^3$. Using the same techniques as above we obtain

$$\log \mathcal{N}_\epsilon(\tilde{\Theta}_u(s_n, s_m)) \leq nr_m \log \frac{6u}{\epsilon} + ms_m \log \frac{ek_m}{s_m} + ms_m \log \frac{2B_{max}\sqrt{mn}s_ms_n}{\epsilon},$$

$$\log \mathcal{N}_\epsilon \left(\tilde{\Theta}_u(s_n, s_m) \right) \leq mr_n \log \frac{6u}{\epsilon} + ns_n \log \frac{ek_n}{s_n} + ns_n \log \frac{2B_{\max} \sqrt{mn} s_m s_n}{\epsilon},$$

and

$$\log \mathcal{N}_\epsilon \left(\tilde{\Theta}_u(s_n, s_m) \right) \leq mn \log \frac{3u}{\epsilon}.$$

Combining these bounds completes the proof of Proposition 7.

F Proof of Theorem 10

Let $\mathcal{I} = \min_{(s_n, s_m)} R(s_n, s_m)$ and $\nu^2 = \frac{(\sigma \vee \theta_{\max})^2}{p} \mathcal{I}$. By the definition of $R(s_n, s_m)$ in (14) we have $\mathcal{I} \geq d$. Note first that if $\|\hat{\theta} - \theta^*\|_2 \leq 2^7 \nu$, then Theorem 10 holds trivially. So, without loss of generality, we can assume that $\hat{\theta} \in \mathcal{X}_\nu \triangleq \{\theta \in \mathcal{X} : \|\theta - \theta^*\|_2 > 2^7 \nu\}$. By the definition (16) of the estimator $\hat{\theta} = \hat{X} \hat{B} \hat{Z}^T$ we have that for any $\theta = XBZ^T \in \mathcal{X}$

$$\|Y - \hat{\theta}_\Omega\|_2^2 + \lambda R(\hat{\theta}) \leq \|Y - \theta_\Omega\|_2^2 + \lambda R(\theta)$$

which implies

$$\|\hat{\theta}_\Omega - \theta_\Omega^*\|_2^2 \leq \|\theta_\Omega - \theta_\Omega^*\|_2^2 - 2\langle \xi_\Omega, \theta - \theta^* \rangle + \lambda R(\theta) + 2\langle \xi_\Omega, \hat{\theta} - \theta^* \rangle - \lambda R(\hat{\theta}) \quad (32)$$

where we set $\xi = (\xi_{ij})$. We will bound each term in (32) separately. Lemma 27 implies that with probability at least $1 - \exp(-pnm) - 2\exp(-d/10)$,

$$\langle \xi_\Omega, \hat{\theta} - \theta^* \rangle \leq 2(\sigma \vee \theta_{\max})^2 R(\hat{\theta}) + \frac{p}{8} \|\hat{\theta} - \theta^*\|_2^2. \quad (33)$$

To control $\langle \xi_\Omega, \theta - \theta^* \rangle$, we use Lemma 28 with $t = p\|\theta - \theta^*\|_2^2 + (\sigma \vee \theta_{\max})^2 R(\theta)$. It follows that, with probability at least $1 - \exp(-d/2)$,

$$\langle \xi_\Omega, \theta - \theta^* \rangle \leq (\sigma \vee \theta_{\max})^2 R(\theta) + p\|\theta - \theta^*\|_2^2 \quad (34)$$

where we have used that $R(\theta) \geq d$. On the other hand, Lemma 24 implies that, with probability at least $1 - \exp(-pnm) - 2\exp(-d/6)$,

$$\|\hat{\theta}_\Omega - \theta_\Omega^*\|_2^2 + 4\theta_{\max}^2 R(\hat{\theta}) \geq \frac{p}{2} \|\hat{\theta} - \theta^*\|_2^2. \quad (35)$$

Finally, using Lemma 26 with $a_{ij} = (\theta - \theta^*)_{ij}^2$ and $t = \frac{p}{2} \|\theta - \theta^*\|_2^2 + 4\theta_{\max}^2 d$ we get that, with probability at least $1 - \exp(-d)$,

$$\|\theta_\Omega - \theta_\Omega^*\|_2^2 \leq 4\theta_{\max}^2 R(\theta) + \frac{3p}{2} \|\theta - \theta^*\|_2^2 \quad (36)$$

where we have used that $R(\theta) \geq d$. Plugging (33) - (36) in (32) we get

$$\begin{aligned} \frac{p}{4} \|\hat{\theta} - \theta^*\|_2^2 &\leq \frac{5p}{2} \|\theta - \theta^*\|_2^2 + 8(\sigma \vee \theta_{\max})^2 R(\hat{\theta}) + 6(\sigma \vee \theta_{\max})^2 R(\theta) \\ &\quad + \lambda R(\theta) - \lambda R(\hat{\theta}) \end{aligned}$$

with probability larger than $1 - 5\exp(-d/10) - 2\exp(-pnm)$ where we have used that $d \geq 10$. Taking here $\lambda = 8(\sigma \vee \theta_{\max})^2$ finishes the proof.

G Proofs of the lemmas

G.1 Lemmas for Theorem 1

Proof of Lemma 11. We start by proving (i). Note that for any fixed $X \in \mathcal{A}_n$, $\theta = XBZ^T$ belongs to a linear space of dimension not greater than $nm \wedge k_n m = r_n m$ as BZ^T belongs to a linear space of dimension not greater than $k_n m$. Thus, $\theta - \theta^*$ belongs to a linear space of dimension not greater than $r_n m + 1$, which we denote by $W_{r_n m}(X)$. We have

$$\sup_{\theta \in \Theta(s_n, s_m), \theta \neq \theta^*} \frac{\langle \theta - \theta^*, W \rangle^2}{\|\theta - \theta^*\|_2^2} \leq \max_{X \in \mathcal{A}_n} U_X$$

where, for a fixed $X \in \mathcal{A}_n$, we define

$$U_X = \sup_{\theta \in \Theta(s_n, s_m), \theta \neq \theta^*, \theta = XBZ^T} \frac{\langle \theta - \theta^*, W \rangle^2}{\|\theta - \theta^*\|_2^2} \leq \sup_{u \in W_{r_n m}(X): \|u\|_2=1} \langle u, W \rangle.$$

It follows from Lemma 13 that

$$\mathbb{P}\{U_X \geq \sigma^2(2(r_n m + 1) + 3v)\} \leq e^{-v}, \quad \forall v > 0.$$

Note that, for any $A \in \mathcal{A}_n$, there are at most s_n non-zero entries in each row of A . This implies that the number of different supports of matrix A is at most $\binom{k_n}{s_n}^n$. For those ns_n non-zero entries, there are at most $|\mathcal{D}_n|^{ns_n}$ choices. Then, we have

$$\log |\mathcal{A}_n| \leq n \log \binom{k_n}{s_n} + ns_n \log |\mathcal{D}_n| \leq ns_n \log \left(\frac{ek_n |\mathcal{D}_n|}{s_n} \right). \quad (37)$$

Applying the union bound and using (37) we get

$$\mathbb{P}\left\{\max_{X \in \mathcal{A}_n} U_X \geq \sigma^2(2(r_n m + 1) + 3v)\right\} \leq \left(\frac{ek_n |\mathcal{D}_n|}{s_n}\right)^{s_n n} e^{-v}, \quad \forall v > 0.$$

which yields the first result of Lemma 11. To get the bound on the expectation, we use the fact that for any non-negative random variable ξ and any $a > 0$

$$\mathbb{P}(\xi \geq a + t) \geq e^{-t}, \quad \forall t > 0$$

implies $\mathbb{E}\xi \leq a + 1$. The proof of (ii) follows the same lines fixing both X and Z . To prove (iii) we use that $\theta = XBZ^T$ belongs to a linear space of dimension not greater than nm . \square

Lemma 13. *Let ξ be a σ -subgaussian random vector in \mathbb{R}^n , and let \mathcal{W} be a linear subspace of \mathbb{R}^n with $\dim(\mathcal{W}) = d$. Consider the Euclidean ball $B(0, 1) = \{u \in \mathcal{W} : \|u\|_2 \leq 1\}$. Then, for any $t > 0$,*

$$\mathbb{P}\left(\max_{u \in B(0, 1)} (u^T \xi)^2 \geq \sigma^2(d + 2\sqrt{dt} + 2t)\right) \leq e^{-t}.$$

Proof. We have

$$\max_{u \in B(0,1)} (u^T \xi)^2 = \max_{u \in B(0,1)} (u^T P_W \xi)^2 = \|P_W \xi\|_2^2$$

where P_W is the orthogonal projector onto \mathcal{W} . Applying the following lemma with $A = P_W$ yields the result.

Lemma 14 (Hsu et al. [24]). *Let ξ be a σ -subgaussian random vector in \mathbb{R}^n , and let $A \in \mathbb{R}^{n \times n}$ be a matrix. Set $\Sigma = A^T A$. Then, for any $t > 0$,*

$$\mathbb{P}(\|A\xi\|_2^2 \geq \sigma^2(\text{Tr}(\Sigma) + 2\sqrt{\text{Tr}(\Sigma^2)t} + 2\lambda_{\max}(\Sigma)t)) \leq e^{-t}$$

where $\text{Tr}(\Sigma)$ and $\lambda_{\max}(\Sigma)$ denote the trace and the maximal eigenvalue of Σ . \square

G.2 Lemmas for Theorem 3

Lemma 15. *Assume that the noise variables W_{ij} in model (1) are i.i.d. Gaussian with distribution $\mathcal{N}(0, \sigma^2)$. Then, the Kullback-Leibler divergence between \mathbb{P}_θ and $\mathbb{P}_{\theta'}$ has the form*

$$KL(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \frac{p}{2\sigma^2} \|\theta - \theta'\|_2^2.$$

Proof of this lemma is straightforward and it is therefore omitted.

Lemma 16. *Let $k \geq 2$ and $s \geq 1$ be integers, $s \leq k$. There exists a subset S_0 of the set of binary sequences $\{0, 1\}^k$ such that*

- (i) $\log |S_0| \geq c_1^* s \log \frac{ek}{s}$,
- (ii) $c_2^* s \leq \|a\|_0 \leq s$ for all $a \in S_0$, and $\|a\|_0 = s$ for all $a \in S_0$ if $s \leq k/2$,
- (iii) $\|a - b\|_2^2 \geq c_3^* s$ for all $a, b \in S_0$ such that $a \neq b$,

where $c_j^* > 0$, $j = 1, 2, 3$, are absolute constants.

Proof. For $s \leq k/2$ the result follows from Lemma A.3 in [34]. For $k/2 < s \leq k$ and $k \geq 32$, we restrict the consideration only to binary sequences in $\{0, 1\}^k$ such that the first $m = \lceil k/4 \rceil$ elements can be either 0 or 1, the last $s - m$ elements are 1 and the remaining elements are 0. Then, (i) - (iii) follow from the Varshamov-Gilbert bound [36, Lemma 2.9] applied to the set of binary sequences of length m . For $k/2 < s \leq k$ and $k < 32$, the result is obvious. \square

Lemma 17. *Let $\{a_1, a_2, \dots, a_N\} \subseteq \{0, 1\}^k$. Let r be an integer satisfying $r > 96 \log N$. Then, there exists a matrix $Q \in \{-1, 1\}^{r \times k}$ such that for any $u, v \in [N]$,*

$$\frac{r}{2} \|a_u - a_v\|_2^2 \leq \|Qa_u - Qa_v\|_2^2 \leq \frac{3r}{2} \|a_u - a_v\|_2^2.$$

Proof. The result follows immediately from Johnson - Lindenstrauss Lemma as stated in [1, Theorem 2] by taking there $\beta = 1$ and $\epsilon = 1/2$. \square

G.3 Lemmas for Theorem 6

Lemma 18. Let $R = \sqrt{\frac{3(\theta_{\max}^2 + \sigma^2)}{p}} N \epsilon_0^2$ and $\Theta^R = \{\theta \in \Theta, \|\theta - \theta_0\|_2 \geq 2^{s^*} R\}$. Then we have

$$\mathbb{P} \left\{ \sup_{\theta \in \Theta^R} \left(\langle \theta - \theta_0, W \rangle - \frac{1}{8} \|\theta - \theta_0\|_2^2 \right) > 96R^2 \right\} \leq 4 \exp \left(-\frac{\alpha R^2}{2} \right) + \exp(-pN/6)$$

where $\alpha = \frac{p}{6(\theta_{\max}^2 + \sigma^2)}$.

Proof of Lemma 18. Let $\mathcal{E} = \{\|W\|_2 \leq \sqrt{\frac{N}{2\alpha}}\}$ and for $s \geq 2$ let $\Theta_s^R = \{\theta \in \Theta^R, 2^s R \leq \|\theta - \theta_0\|_2 \leq 2^{s+1} R\}$. Then we have that $\Theta^R = \cup_{s=s^*}^{\infty} \Theta_s^R$. The union bound yields

$$\begin{aligned} & \mathbb{P} \left\{ \sup_{\theta \in \Theta^R} \left(\langle \theta - \theta_0, W \rangle - \frac{1}{8} \|\theta - \theta_0\|_2^2 \right) > 96R^2, \mathcal{E} \right\} \\ & \leq \sum_{s=s^*}^{\infty} \mathbb{P} \left\{ \sup_{\theta \in \Theta_s^R} \left(\langle \theta - \theta_0, W \rangle - \frac{1}{8} \|\theta - \theta_0\|_2^2 \right) > 96R^2, \mathcal{E} \right\} \\ & \leq \sum_{s=s^*}^{\infty} \mathbb{P} \left\{ \sup_{\theta \in \Theta_s^R} \langle \theta - \theta_0, W \rangle \geq (2^{2s-3} + 96) R^2, \mathcal{E} \right\} \\ & \leq \sum_{s=s^*}^{\infty} \mathbb{P} \left\{ \sup_{\theta \in \Theta_{2^{s+1}R}} \langle \theta - \theta_0, W \rangle \geq (2^{2s-3} + 96) R^2, \mathcal{E} \right\}. \end{aligned}$$

where the last step is due to the fact that $\Theta_s^R \subseteq \Theta_{2^{s+1}R}$. Now we are going to apply Lemma 19 with $D = 2^{s+1}R$, $\epsilon = D\epsilon_0$ and $t = (2^{2s-4} + 48)R^2$. It is easy to check that $t \in [DR, D^2]$ for all $s \geq s^*$. Then (38) yields

$$\sum_{s=s^*}^{\infty} \mathbb{P} \left\{ \sup_{\theta \in \Theta_{2^{s+1}R}} \langle \theta - \theta_0, W \rangle \geq (2^{2s-3} + 96) R^2, \mathcal{E} \right\} \leq \sum_{s=2}^{\infty} 2e^{-\frac{\alpha s R^2}{4}} \leq 4e^{-\frac{\alpha R^2}{2}}.$$

The last inequality holds when $\alpha R^2 \geq 4$. By Lemma 20, $P(\mathcal{E}^c) \leq \exp(-pN/6)$, and therefore we obtain the desired result. \square

Lemma 19. Suppose ϵ satisfies $\sqrt{N}\epsilon \leq 2D\sqrt{\log \mathcal{N}_\epsilon(\Theta_D)}$. Then, for $D \geq \sqrt{2 \log \mathcal{N}_\epsilon(\Theta_D)/\alpha}$ and for any $t \in [D\sqrt{2 \log \mathcal{N}_\epsilon(\Theta_D)/\alpha}, D^2]$, we have

$$\mathbb{P} \left\{ \sup_{\theta \in \Theta_D} \langle \theta - \theta_0, W \rangle \geq 2t, \mathcal{E} \right\} \leq 2 \exp \left(-\frac{\alpha t^2}{2D^2} \right) \quad (38)$$

where $\mathcal{E} = \{\|W\|_2 \leq \sqrt{\frac{N}{2\alpha}}\}$.

Proof of Lemma 19. Let \mathcal{C}_D be an ϵ -covering set of Θ_D under the Frobenius norm. That is, for any $\theta \in \Theta_D$, there exists $\bar{\theta} \in \mathcal{C}_D$ such that $\|\theta - \bar{\theta}\|_2 \leq \epsilon$.

Denote by $N \triangleq \mathcal{N}_\epsilon(\Theta_D)$ the minimum cardinality of such \mathcal{C}_D . This yields

$$\begin{aligned}\langle \theta - \theta_0, W \rangle &= \langle \bar{\theta} - \theta_0, W \rangle + \langle \theta - \bar{\theta}, W \rangle \\ &\leq \max_{\bar{\theta} \in \mathcal{C}_D} \langle \bar{\theta} - \theta_0, W \rangle + \epsilon \|W\|_2,\end{aligned}$$

where we use Cauchy-Schwarz for the last inequality. On the event \mathcal{E} , we have that $\epsilon \|W\|_2 \leq \epsilon \sqrt{\frac{N}{2\alpha}} \leq D \sqrt{2 \log(\mathcal{N}_\epsilon(\Theta_D)) / \alpha}$. It implies

$$\mathbb{P} \left\{ \sup_{\theta \in \Theta_D} \langle \theta - \theta_0, W \rangle \geq 2t, \mathcal{E} \right\} \leq \mathbb{P} \left\{ \max_{\bar{\theta} \in \mathcal{C}_D} \langle \bar{\theta} - \theta_0, W \rangle \geq t \right\}$$

for $t \geq D \sqrt{2 \log \mathcal{N}_\epsilon(\Theta_D) / \alpha}$. By union bound and Lemma 20, the right hand side of the above inequality can be bounded from above by

$$\sum_{\bar{\theta} \in \mathcal{C}_D} \mathbb{P} \{ \langle \bar{\theta} - \theta_0, W \rangle \geq t \} \leq \exp(-\alpha \min\{t^2/D^2, t\}) + \log \mathcal{N}_\epsilon(\Theta_D).$$

Then the desired result (38) holds when $D \sqrt{2 \log N / \alpha} \leq t \leq D^2$. \square

Lemma 20. Let $\alpha = \frac{p}{6(\theta_{\max}^2 + \sigma^2)}$. Then, for any $t > 0$ and $a \in \mathbb{R}^n$, we have

$$\mathbb{P} \{ \langle a, W \rangle > t \} \leq \exp \left(-\alpha \min \left\{ \frac{t^2}{\|a\|^2}, \frac{\theta_{\max} t}{\|a\|_\infty} \right\} \right), \quad (39)$$

and

$$\mathbb{P} \left\{ \|W\|_2 \geq \sqrt{\frac{N}{2\alpha}} \right\} \leq \exp(-pN/6). \quad (40)$$

Proof of Lemma 20. Since $W_i = \theta_i^* \frac{E_i - p}{p} + \xi_i \frac{E_i}{p}$, for $\lambda(\sigma \vee 2\theta_{\max}) \|a\|_\infty \leq p$, we have

$$\mathbb{E} \left(e^{\lambda a_i W_i} \right) \leq \mathbb{E} \left(e^{\lambda a_i \theta_i^* \frac{E_i - p}{p}} e^{\frac{\lambda^2 a_i^2 \sigma^2 E_i}{2p^2}} \right) \leq e^{\frac{\lambda^2 a_i^2 \sigma^2}{2p}} e^{\left(\frac{\lambda a_i \theta_i^*}{p} + \frac{\lambda^2 a_i^2 \sigma^2}{2p^2} \right)^2 p} \leq e^{\frac{3\lambda^2 a_i^2 (\theta_{\max}^2 + \sigma^2)}{2p}}.$$

Here the second inequality is due to the fact that $\mathbb{E} e^{\lambda(E_i - p)} \leq e^{\lambda p}$ for $|\lambda| \leq 1$. Now, following the Chernoff argument as in the proof of Lemma 28 we get (39).

To prove (40) note first that the variance of W_i satisfies

$$\mathbb{E} W_i^2 = \theta_i^2 \frac{1-p}{p} + \frac{\mathbb{E} \xi_i^2}{p} \leq \frac{\theta_{\max}^2 + \sigma^2}{p}.$$

Then we have

$$W_i^2 - \mathbb{E} W_i^2 = (\theta_i^*)^2 \frac{(E_i - p)(1 - 2p)}{p^2} + \frac{E_i}{p^2} (\xi_i^2 - \mathbb{E} \xi_i^2) + \frac{E_i - p}{p^2} \mathbb{E} \xi_i^2 + \frac{E_i(E_i - p)}{p^2} 2\theta_i^* \xi_i$$

When $\lambda(\theta_{\text{mx}}^2 + \sqrt{2}\sigma^2)/p^2 \leq 1$, we obtain

$$\begin{aligned} \mathbb{E} \left(e^{\lambda(W_i^2 - \mathbb{E}W_i^2)} \right) &\leq \mathbb{E} \left(e^{\lambda(\theta_i^*)^2 \frac{(E_i-p)(1-2p)}{p^2}} e^{\frac{2\lambda^2\sigma^4 E_i}{p^4}} e^{\frac{\lambda\sigma^2(E_i-p)}{p^2}} e^{\frac{2\lambda^2\theta_{\text{mx}}^2\sigma^2}{p^4} E_{ij}(E_{ij}-p)^2} \right) \\ &\leq e^{\frac{6\lambda^2(\theta_{\text{mx}}^4 + \sigma^4)}{p^3}}. \end{aligned}$$

The Chernoff argument yields

$$\mathbb{P} \left\{ \sum_{i=1}^N (W_i^2 - \mathbb{E}W_i^2) \geq t \right\} \leq \exp \left\{ -\lambda t + \frac{6\lambda^2(\theta_{\text{mx}}^4 + \sigma^4)}{p^3} N \right\}.$$

For $t = \frac{2(\theta_{\text{mx}}^2 + \sigma^2)}{p} N$, we choose $\lambda = \frac{p^2}{6(\theta_{\text{mx}}^2 + \sigma^2)}$ to get the desired result. \square

G.4 Lemmas for Proposition 7

Lemma 21. For any fixed $X \in \mathbb{R}^{n \times k_n}$ and $Z \in \mathbb{R}^{m \times k_m}$, let

$$\mathcal{T}_R(X, Z) = \{ \theta = X B Z^T, B \in \mathbb{R}^{k_n \times k_m}, \|\theta - \theta_0\|_2 \leq R \}.$$

Then, for any $0 < \epsilon \leq R$

$$\mathcal{N}_\epsilon(\mathcal{T}_R(X, Z)) \leq \left(\frac{3R}{\epsilon} \right)^{r_n r_m}.$$

Proof. Note that for any fixed $X \in \mathbb{R}^{n \times k_n}$ and $Z \in \mathbb{R}^{m \times k_m}$ the set of matrices $\{ \theta = X B Z^T, B \in \mathbb{R}^{k_n \times k_m} \}$ belongs to a linear subspace of $\mathbb{R}^{n \times m}$ of dimension at most $r_n r_m$. To see it, note that any of such matrices θ can be written as $\theta = \sum_{i=1}^{r_n} \sum_{j=1}^{r_m} a_{ij}(\theta) M_{ij}$ with scalars $a_{ij}(\theta) \in \mathbb{R}$ and matrices $M_{ij} \in \mathbb{R}^{n \times m}$. Now, applying the standard bound on the covering number of the ball in the Euclidean norm (see, e.g. Lemma 5.2 in [38]) we get the result of the lemma. \square

Lemma 22. We have the following upper bound on the ϵ -covering number of \mathcal{A}_n under ℓ_∞ norm:

$$\mathcal{N}_\epsilon(\mathcal{A}_n, \|\cdot\|_\infty) \leq \binom{k_n}{s_n}^n \left(\frac{1}{\epsilon} \right)^{ns_n}.$$

Proof. Note that there are $\binom{k_n}{s_n}^n$ subsets of $\{1, \dots, k_n\}^n$ that satisfy the column sparsity constraint of $A \in \mathcal{A}_n$. For any such subset, the selected ns_n entries lie in the unit Euclidean ball $\mathbb{B}_2(1)$. By the standard volume ratio argument we can find an ϵ -covering set of $\mathbb{B}_2(1)$ with at most $(\frac{1}{\epsilon})^{ns_n}$ elements. Hence, the lemma follows. \square

Lemma 23. Assume that $A \in \mathbb{R}^{n \times k}$, $Z \in \mathbb{R}^{m \times k}$ and that each row of Z is s -sparse. Then,

$$\|AZ^T\|_2 \leq s\sqrt{mn}\|A\|_\infty\|Z\|_\infty.$$

Proof. We have

$$\begin{aligned}\|AZ^T\|_2^2 &= \sum_{i \in [n]} \sum_{j \in [m]} \left(\sum_{l \in [k]} A_{il} Z_{jl} \right)^2 \leq \sum_{i \in [n]} \sum_{j \in [m]} (s \|A\|_\infty \|Z\|_\infty)^2 \\ &\leq s^2 nm \|A\|_\infty^2 \|Z\|_\infty^2.\end{aligned}$$

□

G.5 Lemmas for Theorem 10

Lemma 24. Assume that $nm \log(3\sqrt{nm}) \geq 6 \log(k_n k_m)$. Then, with probability larger than $1 - 2 \exp(-d/6) - \exp(-pnm)$

$$\sup_{\theta \in \mathcal{X}_\nu} \frac{\frac{p}{2} \|\theta - \theta^*\|_2^2 - \|\theta_\Omega - \theta_\Omega^*\|_2^2}{\theta_{\max}^2 R(\theta)} \leq 4$$

where $R(\theta)$ is defined in (15).

Proof. Let $\mathcal{E}_2 = \left\{ \sum_{ij} (E_{ij} - p)^2 \leq 3pnm \right\}$. Lemma 26 implies that $\mathbb{P}(\mathcal{E}_2) \geq 1 - \exp(-pnm)$. Using the definition of \mathcal{X}_ν we have that

$$\mathbb{P} \left\{ \sup_{\theta \in \mathcal{X}_\nu} \frac{\frac{p}{2} \|\theta - \theta^*\|_2^2 - \|\theta_\Omega - \theta_\Omega^*\|_2^2}{\theta_{\max}^2 R(\theta)} \geq 4, \mathcal{E}_2 \right\} \leq \sum_{s_n=1}^{k_n} \sum_{s_m=1}^{k_m} \mathbb{I}_{s_n, s_m} \quad (41)$$

where

$$\mathbb{I}_{s_n, s_m} = \mathbb{P} \left\{ \sup_{\theta \in \tilde{\Theta}^\nu(s_n, s_m)} \frac{p}{2} \|\theta - \theta^*\|_2^2 - \|\theta_\Omega - \theta_\Omega^*\|_2^2 \geq 4 \theta_{\max}^2 R(s_n, s_m), \mathcal{E}_2 \right\}$$

and $\tilde{\Theta}^\nu(s_n, s_m) \triangleq \{\theta \in \tilde{\Theta}(s_n, s_m), \|\theta - \theta^*\|_2 > 2^6 \nu\}$. In order to bound \mathbb{I}_{s_n, s_m} from above, we use a standard peeling argument. Let $\mu = 2$. For $l \in \mathbb{N}$ set

$$S_l = \left\{ \theta \in \tilde{\Theta}(s_n, s_m) : \mu^l \nu \leq \|\theta - \theta^*\|_2 \leq \mu^{l+1} \nu \right\}.$$

Then, $\tilde{\Theta}^\nu(s_n, s_m) = \bigcup_{l=7}^{\infty} S_l$ and the union bound yields

$$\begin{aligned}\mathbb{I}_{s_n, s_m} &\leq \sum_{l=7}^{\infty} \mathbb{P} \left\{ \sup_{\theta \in S_l} \frac{p}{2} \|\theta - \theta^*\|_2^2 - \|\theta_\Omega - \theta_\Omega^*\|_2^2 \geq 4 \theta_{\max}^2 R(s_n, s_m), \mathcal{E}_2 \right\} \\ &\leq \sum_{l=7}^{\infty} \mathbb{P} \left\{ \sup_{\theta \in \tilde{\Theta}_r(s_n, s_m)} p \|\theta - \theta^*\|_2^2 - \|\theta_\Omega - \theta_\Omega^*\|_2^2 \geq 4 \theta_{\max}^2 R(s_n, s_m) + \frac{pr^2}{2\mu^2}, \mathcal{E}_2 \right\} \\ &\quad (42)\end{aligned}$$

where $r = \mu^{l+1}\nu$. We have that $\mathbb{E} \left(E_{ij} (\theta - \theta^*)_{ij}^2 \right) = p (\theta - \theta^*)_{ij}^2$ and

$$p \|\theta - \theta^*\|_2^2 - \|\theta_\Omega - \theta_\Omega^*\|_2^2 = \sum_{(ij)} (p - E_{ij}) (\theta - \theta^*)_{ij}^2.$$

Let $S = \min \left\{ s \geq 2 : 2^{-s} \leq \sqrt{\frac{\mathcal{I}}{nm}} \right\}$ and $\{\tilde{G}_j^S\}_{j=1}^{N_s}$ be a minimal $2^{-S}r$ -covering set of $\Theta_r(s_n, s_m)$ in Frobenius norm given by Proposition 7. Let $G_j^S = \Pi(\tilde{G}_j^S)$ where Π is the projection operator under the Frobenius norm into the set

$$\mathcal{B} = \{\theta \in \mathbb{R}^{n \times m} : \|\theta\|_\infty \leq \theta_{\max}\}.$$

As \mathcal{B} is closed and convex, Π is non-expansive and we have that for any $\theta \in \tilde{\Theta}_r(s_n, s_m)$, $\|\theta - G_j^S\|_2 \leq \|\theta - \tilde{G}_j^S\|_2$. Then, there exists $G_\theta^S \in \{G_j^S\}_{j=1}^{N_s}$ such that on the event \mathcal{E}_2

$$\left| \sum_{(ij)} (p - E_{ij}) \left[(\theta^* - \theta)_{ij}^2 - (\theta^* - G_\theta^S)_{ij}^2 \right] \right| \leq \frac{pr^2}{4\mu^2}$$

where we use $(\theta^* - \theta)_{ij}^2 - (\theta^* - G_\theta^S)_{ij}^2 = (G_\theta^S - \theta)_{ij} (2\theta^* - G_\theta^S - \theta)_{ij}$, Cauchy-Schwarz inequality and $\frac{pr}{16\theta_{\max}\mu^2\sqrt{3pnm}} \geq \sqrt{\frac{\mathcal{I}}{nm}}$. So it suffices to prove the exponential inequality for

$$\mathbb{P} \left\{ \max_{k=1, \dots, N_s} \sum_{(ij)} (p - E_{ij}) (\theta^* - G_k^S)_{ij}^2 \geq 4\theta_{\max}^2 R(s_n, s_m) + \frac{pr^2}{4\mu^2} \right\}.$$

We apply Markov's inequality. Set $t = R(s_n, s_m) + \frac{pr^2}{16\theta_{\max}^2\mu^2}$

$$\begin{aligned} & \mathbb{P} \left\{ \max_{k=1, \dots, N_s} \frac{1}{4\theta_{\max}^2} \sum_{(ij)} (p - E_{ij}) (\theta^* - G_k^S)_{ij}^2 \geq t \right\} \\ & \leq e^{-t} \mathbb{E} \exp \left(\max_{k=1, \dots, N_s} \frac{1}{4\theta_{\max}^2} \sum_{(ij)} (p - E_{ij}) (\theta^* - G_k^S)_{ij}^2 \right) \\ & \leq e^{-t} \sum_{k=1}^{N_s} \mathbb{E} \exp \left(\frac{1}{4\theta_{\max}^2} \sum_{(ij)} (p - E_{ij}) (\theta^* - G_k^S)_{ij}^2 \right) \\ & \leq e^{-t} \sum_{k=1}^{N_s} \prod_{(ij)} \mathbb{E} \exp \left((p - E_{ij}) \frac{(\theta^* - G_k^S)_{ij}^2}{4\theta_{\max}^2} \right). \end{aligned} \quad (43)$$

Now we use the following lemma that follows easily from [21, p.22], see also [17]:

Lemma 25. *Let $E \sim \text{Ber}(p)$, then for any $|\lambda| \leq 1$ we have*

$$\mathbb{E} \exp\{\lambda(p - E)\} \leq e^{\lambda^2 p}.$$

Lemma 25 and (43) imply

$$\begin{aligned}
& \mathbb{P} \left\{ \max_{k=1, \dots, N_S} \frac{1}{4\theta_{\max}^2} \sum_{(ij)} (p - E_{ij})(\theta^* - G_k^S)_{ij}^2 \geq t \right\} \\
& \leq e^{-t} \sum_{k=1}^{N_S} \prod_{(ij)} \exp \left(\frac{p(\theta^* - G_k^S)_{ij}^4}{16\theta_{\max}^4} \right) \\
& \leq \exp \left(\frac{pr^2}{2^8\theta_{\max}^2} + \log(N_S) - t \right) \leq \exp \left(-\frac{3pr^2}{64\mu^2\theta_{\max}^2} \right)
\end{aligned}$$

where we use $S \geq 3$ and Lemma 29 which implies $\log(N_S) \leq R(s_n, s_m)$. Putting this last bound into (42) and (41) we get

$$\begin{aligned}
& \mathbb{P} \left\{ \sup_{\theta \in \mathcal{X}_\nu} \frac{\frac{p}{2} \|\theta - \theta^*\|_2^2 - \|\theta_\Omega - \theta_\Omega^*\|_2^2}{\theta_{\max}^2 R(\theta)} \geq 4, \mathcal{E}_2 \right\} \\
& \leq \sum_{s_n=1}^{k_n} \sum_{s_m=1}^{k_m} \sum_{l=6}^{\infty} \exp \left(-\frac{3l \log(\mu) \mathcal{I}}{32} \right) \\
& \leq 2 \sum_{s_n=1}^{k_n} \sum_{s_m=1}^{k_m} \exp \left(-\frac{\mathcal{I}}{3} \right) \leq 2 \exp \left(-\frac{\mathcal{I}}{6} \right)
\end{aligned}$$

where we used $e^x \geq x$ and $\nu^2 = \frac{(\sigma \vee \theta_{\max})^2}{p} \mathcal{I}$. In the last inequality we use condition $n m \log(3\sqrt{n \wedge m}) \geq 6 \log(k_n k_m)$ which implies $\frac{\mathcal{I}}{6} \geq \log(k_n k_m)$. Finally, $\mathcal{I} = \min_{(s_n, s_m)} R(s_n, s_m) \geq (n + m)$ implies the result of Lemma 24. \square

Lemma 26. Let $a = (a_{ij}) \in \mathbb{R}^{nm}$. Then, for any $t > 0$ we have

$$\mathbb{P} \left\{ \sum_{ij} (E_{ij} - p) a_{ij} > t \right\} \leq \exp \left(-\min \left\{ \frac{t^2}{4p\|a\|_2^2}, \frac{t}{2\|a\|_\infty} \right\} \right), \quad (44)$$

and

$$\mathbb{P} \left\{ \sum_{ij} (E_{ij} - p)^2 \geq 3pnm \right\} \leq \exp(-pnm). \quad (45)$$

Proof. This is a direct consequence of Bernstein's inequality. For (44) we note that

$$\mathbb{E} \sum_{ij} (E_{ij} - p)^2 a_{ij}^2 \leq p(1-p)\|a\|_2^2.$$

On the other hand, for any $(ij) \in [n] \times [m]$ we have $|(E_{ij} - p)a_{ij}| \leq \|a\|_\infty$. Then, Bernstein's inequality implies (44).

For (45) we have

$$\mathbb{E} \sum_{ij} (E_{ij} - p)^4 \leq \mathbb{E} \sum_{ij} (E_{ij} - p)^2 = p(1-p)nm$$

and

$$|E_{ij} - p|^2 \leq p^2 \vee (1-p)^2 \leq 1.$$

Then, using Bernstein's inequality we get

$$\mathbb{P} \left\{ \sum_{ij} (E_{ij} - p)^2 \geq t + pnm(1-p) \right\} \leq \exp \left\{ -\frac{t^2}{2pnm(1-p) + 2t/3} \right\};$$

Choosing $t = 2pnm$ we get the result of Lemma 26. \square

Lemma 27. *Assume that $nm \log(3\sqrt{nm}) \geq 6 \log(k_n k_m)$. Then, with probability greater than $1 - 2 \exp(-d/10) - \exp(-pnm)$,*

$$\sup_{\theta \in \mathcal{X}_\nu} \frac{\langle \xi_\Omega, \theta - \theta^* \rangle - \frac{p}{8} \|\theta - \theta^*\|_2^2}{(\sigma \vee \theta_{\max})^2 R(\theta)} \leq 2.$$

where $R(\theta)$ is defined in (15), and \mathcal{X}_ν is defined in the proof of Theorem 10.

Proof. We proceed as in the proof of Lemma 24. Let $\mathcal{E}_1 = \left\{ \sqrt{\sum_{ij} E_{ij} \xi_{ij}^2} \leq \sigma \sqrt{6pnm} \right\}$. Lemma 28 implies that $\mathbb{P}(\mathcal{E}_1) \geq 1 - \exp(-pnm)$. Using the definition of \mathcal{X}_ν we get

$$\mathbb{P} \left\{ \sup_{\theta \in \mathcal{X}_\nu} \frac{\langle \xi_\Omega, \theta - \theta^* \rangle - \frac{p}{8} \|\theta - \theta^*\|_2^2}{(\sigma \vee \theta_{\max})^2 R(\theta)} \geq 2, \mathcal{E}_1 \right\} \leq \sum_{s_n=1}^{k_n} \sum_{s_m=1}^{k_m} I'_{s_n, s_m} \quad (46)$$

where

$$I'_{s_n, s_m} = \mathbb{P} \left\{ \sup_{\theta \in \tilde{\Theta}^\nu(s_n, s_m)} \langle \xi_\Omega, \theta - \theta^* \rangle - \frac{p}{8} \|\theta - \theta^*\|_2^2 \geq 2 (\sigma \vee \theta_{\max})^2 R(s_n, s_m), \mathcal{E}_1 \right\}.$$

Using a standard peeling argument we find

$$\begin{aligned} I'_{s_n, s_m} &\leq \sum_{l=7}^{\infty} \mathbb{P} \left\{ \sup_{\theta \in S_l} \langle \xi_\Omega, \theta - \theta^* \rangle - \frac{p}{8} \|\theta - \theta^*\|_2^2 \geq 2 (\sigma \vee \theta_{\max})^2 R(s_n, s_m), \mathcal{E}_1 \right\} \\ &\leq \sum_{l=7}^{\infty} \mathbb{P} \left\{ \sup_{\theta \in \tilde{\Theta}_r(s_n, s_m)} \langle \xi_\Omega, \theta - \theta^* \rangle \geq 2 (\sigma \vee \theta_{\max})^2 R(s_n, s_m) + \frac{pr^2}{8\mu^2}, \mathcal{E}_1 \right\} \end{aligned} \quad (47)$$

where $r = \mu^{l+1}\nu$. Let $\{\tilde{G}_j^S\}_{j=1}^{N_s}$ be a minimal $2^{-S}r$ -covering set of $\Theta_r(s_n, s_m)$ in the Frobenius norm and

$$S = \min \left\{ s \geq 1 : 2^{-s} \leq \sqrt{\frac{\mathcal{I}}{nm}} \right\}.$$

As in the proof of Lemma 24 we define $G_j^S = \Pi(\tilde{G}_j^S)$. Then, there exists $G_\theta^S \in \{G_j^S\}_{j=1}^{N_s}$ such that on the event \mathcal{E}_1

$$\left| \sum_{(ij)} E_{ij} \xi_{ij} \left[(\theta - \theta^*)_{ij} - (G_\theta^S - \theta^*)_{ij} \right] \right| \leq \frac{pr^2}{16\mu^2}$$

where we use $\sqrt{\frac{\mathcal{I}}{nm}} \leq \frac{pr}{16\mu^2\sigma\sqrt{6pnm}}$. So, it suffices to prove the exponential inequality for

$$\mathbb{P} \left\{ \max_{k=1, \dots, N_S} \sum_{(ij)} E_{ij} \xi_{ij} (G_k^S - \theta^*)_{ij} \geq 2(\sigma \vee \theta_{\max})^2 R(s_n, s_m) + \frac{pr^2}{16\mu^2} \right\}.$$

We apply Markov's inequality. Set $t = R(s_n, s_m) + \frac{pr^2}{32(\sigma \vee \theta_{\max})^2\mu^2}$, then

$$\begin{aligned} & \mathbb{P} \left\{ \max_{k=1, \dots, N_S} \frac{1}{2(\sigma \vee \theta_{\max})^2} \sum_{(ij)} E_{ij} \xi_{ij} (G_k^S - \theta^*)_{ij} \geq t \right\} \\ & \leq e^{-t} \mathbb{E} \exp \left(\max_{k=1, \dots, N_S} \frac{1}{2(\sigma \vee \theta_{\max})^2} \sum_{(ij)} E_{ij} \xi_{ij} (G_k^S - \theta^*)_{ij} \right) \\ & \leq e^{-t} \sum_{k=1}^{N_S} \mathbb{E} \exp \left(\frac{1}{2(\sigma \vee \theta_{\max})^2} \sum_{(ij)} E_{ij} \xi_{ij} (G_k^S - \theta^*)_{ij} \right) \\ & \leq e^{-t} \sum_{k=1}^{N_S} \prod_{(ij)} \mathbb{E} \exp \left(E_{ij} \xi_{ij} \frac{(G_k^S - \theta^*)_{ij}}{2(\sigma \vee \theta_{\max})^2} \right). \end{aligned} \quad (48)$$

Lemma 25 implies that $\mathbb{E}(e^{\lambda E_{ij}}) \leq e^{2\lambda p}$ for any $0 \leq \lambda \leq 1$. Then, using Assumption 1 and (48) we get

$$\begin{aligned} & \mathbb{P} \left\{ \max_{k=1, \dots, N_S} \frac{1}{2(\sigma \vee \theta_{\max})^2} \sum_{(ij)} E_{ij} \xi_{ij} (G_k^S - \theta^*)_{ij} \geq t \right\} \\ & \leq e^{-t} \sum_{k=1}^{N_S} \prod_{(ij)} \exp \left(\frac{p(G_k^S - \theta^*)_{ij}^2}{4(\sigma \vee \theta_{\max})^2} \right) \\ & \leq \exp \left(\frac{pr^2}{2^8(\sigma \vee \theta_{\max})^2} + \log(N_S) - t \right) \leq \exp \left(-\frac{3pr^2}{128\mu^2(\sigma \vee \theta_{\max})^2} \right). \end{aligned}$$

The last inequality follows from Lemma 29 which implies $\log(N_S) \leq R(s_n, s_m)$. Combining the last display with (46) and (47) we find

$$\begin{aligned} & \mathbb{P} \left\{ \sup_{\theta \in \mathcal{X}_\nu} \frac{\langle \xi_\Omega, \theta - \theta^* \rangle - \frac{p}{8} \|\theta - \theta^*\|_2^2}{(\sigma \vee \theta_{\max})^2 R(\theta)} \geq 2, \mathcal{E}_1 \right\} \\ & \leq \sum_{s_n=1}^{k_n} \sum_{s_m=1}^{k_m} \sum_{l=7}^{\infty} \exp \left(-\frac{3l \log(2)\mathcal{I}}{64} \right) \\ & \leq 2 \sum_{s_n=1}^{k_n} \sum_{s_m=1}^{k_m} \exp \left(-\frac{\mathcal{I}}{5} \right) \leq 2 \exp \left(-\frac{n+m}{10} \right) \end{aligned}$$

where we have used that $e^x \geq x$, $\nu^2 = \frac{(\sigma \vee \theta_{\max})^2}{p} \mathcal{I}$, $\mathcal{I} = \min_{(s_n, s_m)} R(s_n, s_m) \geq (n+m)$, and $nm \log(3\sqrt{nm}) \geq 6 \log(k_n k_m)$. \square

Lemma 28. *Let $a = (a_{ij}) \in \mathbb{R}^{n \times m}$. Then, for any $t > 0$ we have*

$$\mathbb{P} \left\{ \sum_{ij} a_{ij} E_{ij} \xi_{ij} > t \right\} \leq \exp \left(-\min \left\{ \frac{t^2}{4\sigma^2 p \|a\|_2^2}, \frac{t}{\sqrt{2}\sigma \|a\|_\infty} \right\} \right), \quad (49)$$

and

$$\mathbb{P} \left\{ \sqrt{\sum_{ij} E_{ij} \xi_{ij}^2} \geq \sigma \sqrt{6pnm} \right\} \leq \exp(-pnm). \quad (50)$$

Proof. Lemma 25 implies that $\mathbb{E}(e^{\lambda E_{ij}}) \leq e^{2\lambda p}$ for any $0 \leq \lambda \leq 1$. Then, using Assumption 1 for $0 \leq \lambda \leq \sqrt{2}/(\sigma \|a\|_\infty)$ we obtain

$$\mathbb{E}[\exp(\lambda E_{ij} a_{ij} \xi_{ij})] \leq \mathbb{E}[\exp(\lambda^2 E_{ij} a_{ij}^2 \sigma^2 / 2)] \leq \exp(\lambda^2 a_{ij}^2 \sigma^2 p).$$

The Chernoff argument yields

$$\begin{aligned} \mathbb{P} \left\{ \sum_{ij} a_{ij} E_{ij} \xi_{ij} > t \right\} &= \mathbb{P} \left\{ e^{\sum_{ij} (\lambda a_{ij} E_{ij} \xi_{ij})} > e^{\lambda t} \right\} \\ &\leq e^{-\lambda t} \exp(\lambda^2 \sigma^2 p \|a\|_2^2). \end{aligned}$$

Now, choosing $\lambda = \frac{t}{2\sigma^2 p \|a\|_2^2}$ if $t \leq \frac{2\sqrt{2}\sigma p \|a\|_2^2}{\|a\|_\infty}$ and $\lambda = \frac{\sqrt{2}}{\sigma \|a\|_\infty}$ if $t > \frac{2\sqrt{2}\sigma p \|a\|_2^2}{\|a\|_\infty}$ we get (49).

We prove (50) in a similar way. Using Markov's inequality for $\lambda = (2\sigma^2)^{-1}$

we find

$$\begin{aligned}
\mathbb{P} \left\{ \sum_{ij} E_{ij} \xi_{ij}^2 \geq t \right\} &= \mathbb{P} \left\{ e^{\lambda \sum_{ij} E_{ij} \xi_{ij}^2} > e^{\lambda t} \right\} \\
&\leq e^{-\lambda t} \Pi_{ij} \mathbb{E} \exp(\lambda E_{ij} (\xi_{ij}^2 - \sigma^2) + \lambda E_{ij} \sigma^2) \\
&\leq e^{-\lambda t} \Pi_{ij} \mathbb{E} \exp(2\lambda^2 E_{ij} \sigma^4 + \lambda E_{ij} \sigma^2) \\
&\leq e^{-t/(2\sigma^2) + 2pnm}
\end{aligned}$$

and we take $t = 6\sigma^2 pnm$. □

Lemma 29. *Let $N_S = \mathcal{N}_{r2^{-S}}(\Theta_r(s_n, s_m))$ where*

$$S = \min \left\{ s \geq 3 : 2^{-s} \leq \sqrt{\frac{\mathcal{I}}{nm}} \right\} \quad \text{and} \quad r = 2^{l+1}(\sigma \vee \theta_{\max}) \sqrt{\frac{\mathcal{I}}{p}}$$

for $l \geq 6$. We have that

$$R(s_n, s_m) \geq \log N_S.$$

Proof. We use Proposition 7 and $\sqrt{\frac{\mathcal{I}}{nm}} \geq 2^{-S} \geq \frac{1}{2} \sqrt{\frac{\mathcal{I}}{nm}}$ to get

$$\begin{aligned}
R_1(r2^{-S}) &\leq ns_n \log \frac{ek_n}{s_n} + ms_m \log \frac{ek_m}{s_m} + (ns_n + ms_m) \log \left(\frac{3B_{\max} s_m s_n}{2^6(\sigma \vee \theta_{\max})} \frac{nm\sqrt{p}}{\mathcal{I}} \right) \\
&\quad + r_n r_m \log \left(9 \sqrt{\frac{nm}{\mathcal{I}}} \right) \\
R_2(r2^{-S}) &\leq nr_m \log \left(6 \sqrt{\frac{nm}{\mathcal{I}}} \right) + ms_m \log \frac{ek_m}{s_m} + ms_m \log \left(\frac{B_{\max} s_m s_n}{2^6(\sigma \vee \theta_{\max})} \frac{nm\sqrt{p}}{\mathcal{I}} \right), \\
R_3(r2^{-S}) &\leq mr_n \log \left(6 \sqrt{\frac{nm}{\mathcal{I}}} \right) + ns_n \log \frac{ek_n}{s_n} + ns_n \log \left(\frac{B_{\max} s_m s_n}{2^6(\sigma \vee \theta_{\max})} \frac{nm\sqrt{p}}{\mathcal{I}} \right), \\
R_4(r2^{-S}) &\leq mn \log \left(3 \sqrt{\frac{nm}{\mathcal{I}}} \right)
\end{aligned}$$

and we have that

$$R_1(r2^{-S}) \wedge R_2(r2^{-S}) \wedge R_3(r2^{-S}) \wedge R_4(r2^{-S}) \leq R(s_n, s_m)$$

where we have used that $\mathcal{I} \geq n + m$. □

References

- [1] Dimitris Achlioptas. Database-friendly random projections: Johnson-Lindenstrauss with binary coins. *J. Comput. System Sci.*, 66(4):671–687, 2003. Special issue on PODS 2001 (Santa Barbara, CA).

- [2] Edoardo M Airoldi, David M Blei, Stephen E Fienberg, and Eric P Xing. Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9(Sep):1981–2014, 2008.
- [3] Martin Azizyan, Aarti Singh, and Larry Wasserman. Minimax theory for high-dimensional gaussian mixtures with sparse mean separation. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 2139–2147. Curran Associates, Inc., 2013.
- [4] Afonso S. Bandeira and Ramon van Handel. Sharp nonasymptotic bounds on the norm of random matrices with independent entries. *Ann. Probab.*, 44(4):2479–2506, 07 2016.
- [5] Mikhail Belkin and Kaushik sinha. Polynomial learning of distribution families. In *FOCS 2010: Proceedings of the 51st Annual IEEE Symposium on Foundations of Computer Science*, pages 103–112. IEEE Computer Society, 2010.
- [6] Christian Borgs, Jennifer Chayes, and Adam Smith. Private graphon estimation for sparse graphs. In *Advances in Neural Information Processing Systems*, pages 1369–1377, 2015.
- [7] T. Tony Cai and Wen-Xin Zhou. Matrix completion via max-norm constrained optimization. *Electron. J. Stat.*, 10(1):1493–1525, 2016.
- [8] Emmanuel J Candes and Yaniv Plan. Matrix completion with noise. *Proceedings of the IEEE*, 98(6):925–936, 2010.
- [9] Emmanuel J. Candès and Terence Tao. The power of convex relaxation: near-optimal matrix completion. *IEEE Trans. Inform. Theory*, 56(5):2053–2080, 2010.
- [10] Stanley H. Chan and Edoardo M. Airoldi. A consistent histogram estimator for exchangeable graph models. In *Proceedings of The 31st International Conference on Machine Learning*, pages 208–216, 2014.
- [11] Sourav Chatterjee. Matrix estimation by universal singular value thresholding. *The Annals of Statistics*, 43(1):177–214, 2014.
- [12] Kamalika Chaudhuri, Sanjoy Dasgupta, and Andrea Vattani. Learning mixtures of gaussians using the k-means algorithm. *CoRR*, abs/0912.0086, 2009.
- [13] Yizong Cheng and George M Church. Biclustering of expression data. In *Ismb*, volume 8, pages 93–103, 2000.
- [14] S. Dasgupta. Learning mixtures of gaussians. In *Foundations of Computer Science, 1999. 40th Annual Symposium on*, pages 634–644, 1999.

- [15] Sanjoy Dasgupta and Leonard J. Schulman. A two-round variant of em for gaussian mixtures. In *Proceedings of the Sixteenth Conference on Uncertainty in Artificial Intelligence*, UAI'00, pages 152–159, San Francisco, CA, USA, 2000. Morgan Kaufmann Publishers Inc.
- [16] R. Foygel and N. Srebro. Concentration-based guarantees for low-rank matrix reconstruction. *Journal 24nd Annual Conference on Learning Theory (COLT)*, 2011.
- [17] Chao Gao, Yu Lu, Zongming Ma, and Harrison H. Zhou. Optimal estimation and completion of matrices with biclustering structures. *J. Mach. Learn. Res.*, 17(1):5602–5630, January 2016.
- [18] Chao Gao, Yu Lu, and Harrison H Zhou. Rate-optimal graphon estimation. *The Annals of Statistics*, 43(6):2624–2652, 2015.
- [19] David Gross. Recovering low-rank matrices from few coefficients in any basis. *IEEE Trans. Inform. Theory*, 57(3):1548–1566, 2011.
- [20] John A Hartigan. Direct clustering of a data matrix. *Journal of the american statistical association*, 67(337):123–129, 1972.
- [21] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *J. Amer. Statist. Assoc.*, 58:13–30, 1963.
- [22] Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. *Social networks*, 5(2):109–137, 1983.
- [23] Daniel Hsu and Sham M. Kakade. Learning mixtures of spherical gaussians: Moment methods and spectral decompositions. In *Proceedings of the 4th Conference on Innovations in Theoretical Computer Science*, ITCS '13, pages 11–20, New York, NY, USA, 2013. ACM.
- [24] Daniel Hsu, Sham M. Kakade, and Tong Zhang. A tail inequality for quadratic forms of subgaussian random vectors. *Electron. Commun. Probab.*, 17:no. 52, 6, 2012.
- [25] Brian Karrer and Mark EJ Newman. Stochastic blockmodels and community structure in networks. *Physical Review E*, 83(1):016107, 2011.
- [26] Raghunandan H. Keshavan, Andrea Montanari, and Sewoong Oh. Matrix completion from noisy entries. *J. Mach. Learn. Res.*, 11:2057–2078, 2010.
- [27] O. Klopp. Rank penalized estimators for high-dimensional matrices. *Electron. J. Statist.*, 5:1161–1183, 2011.
- [28] O. Klopp. Noisy low-rank matrix completion with general sampling distribution. *Bernoulli*, 20(1):282–303, 2014.

- [29] Olga Klopp, Alexandre B. Tsybakov, and Nicolas Verzelen. Oracle inequalities for network models and sparse graphon estimation. *Ann. Statist.*, 45(1):316–354, 2017.
- [30] V. Koltchinskii, K. Lounici, and A. B. Tsybakov. Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *Ann. Statist.*, 39(5):2302–2329, 2011.
- [31] Vladimir Koltchinskii, Karim Lounici, and Alexandre B. Tsybakov. Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *Ann. Statist.*, 39(5):2302–2329, 2011.
- [32] S. Negahban and M. J. Wainwright. Restricted strong convexity and weighted matrix completion: Optimal bounds with noise. *Journal of Machine Learning Research*, 13:1665–1697, 2012.
- [33] Bruno A Olshausen and David J Field. Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision research*, 37(23):3311–3325, 1997.
- [34] P. Rigollet and A. Tsybakov. Exponential screening and optimal rates of sparse estimation. *Annals of Statistics*, 39:731–771, 2011.
- [35] A. Soni, S. Jain, J. Haupt, and S. Gonella. Noisy matrix completion under sparse factor models. *IEEE Transactions on Information Theory*, 62(6):3636–3661, June 2016.
- [36] Alexandre B. Tsybakov. *Introduction to nonparametric estimation*. Springer Series in Statistics. Springer, New York, 2009.
- [37] Santosh Vempala and Grant Wang. A spectral algorithm for learning mixture models. *J. Comput. Syst. Sci.*, 68(4):841–860, June 2004.
- [38] R. Vershynin. *Introduction to the non-asymptotic analysis of random matrices*. Cambridge University Press. Chapter 5 of the book Compressed Sensing, Theory and Applications, ed. Y. Eldar and G. Kutyniok.
- [39] Patrick J. Wolfe and Sofia C. Olhede. Nonparametric graphon estimation. *arXiv preprint arXiv:1309.5936*, 2013.
- [40] Jiaming Xu, Laurent Massoulié, and Marc Lelarge. Edge Label Inference in Generalized Stochastic Block Models: from Spectral Theory to Impossibility Results. In *Conference on Learning Theory*, pages 903–920, Barcelona, Spain, June 2014.
- [41] Yuhong Yang and Andrew Barron. Information-theoretic determination of minimax rates of convergence. *Annals of Statistics*, pages 1564–1599, 1999.