

High-dimensional principal component analysis with heterogeneous missingness

Ziwei Zhu*, Tengyao Wang*,[†] and Richard J. Samworth*

*Statistical Laboratory, University of Cambridge

[†]Department of Statistical Science, University College London

June 28, 2019

Abstract

We study the problem of high-dimensional Principal Component Analysis (PCA) with missing observations. In simple, homogeneous missingness settings with a noise level of constant order, we show that an existing inverse-probability weighted (IPW) estimator of the leading principal components can (nearly) attain the minimax optimal rate of convergence. However, deeper investigation reveals both that, particularly in more realistic settings where the missingness mechanism is heterogeneous, the empirical performance of the IPW estimator can be unsatisfactory, and moreover that, in the noiseless case, it fails to provide exact recovery of the principal components. Our main contribution, then, is to introduce a new method for high-dimensional PCA, called **primePCA**, that is designed to cope with situations where observations may be missing in a heterogeneous manner. Starting from the IPW estimator, **primePCA** iteratively projects the observed entries of the data matrix onto the column space of our current estimate to impute the missing entries, and then updates our estimate by computing the leading right singular space of the imputed data matrix. It turns out that the interaction between the heterogeneity of missingness and the low-dimensional structure is crucial in determining the feasibility of the problem. We therefore introduce an incoherence condition on the principal components and prove that in the noiseless case, the error of **primePCA** converges to zero at a geometric rate when the signal strength is not too small. An important feature of our theoretical guarantees is that they depend on average, as opposed to worst-case, properties of the missingness mechanism. Our numerical studies on both simulated and real data reveal that **primePCA** exhibits very encouraging performance across a wide range of scenarios.

1 Introduction

One of the ironies of working with Big Data is that missing data play an ever more significant role, and often present serious difficulties for analysis. For instance, a common approach to handling missing data is to perform a so-called *complete-case analysis* (Little and Rubin, 2014), where we restrict attention to individuals in our study with no missing attributes. When relatively few features are recorded for each individual, one can frequently expect a sufficiently large proportion of complete cases that, under an appropriate missing at random hypothesis, a complete-case analysis may result in only a relatively small loss of efficiency. On the other hand, in high-dimensional regimes where there are many features of interest, there is often such a small proportion of complete cases that this approach becomes infeasible. As a very simple illustration of this phenomenon, imagine an $n \times p$ data matrix in which each entry is missing independently with probability 0.01. When $p = 5$, a complete-case analysis would result in around 95% of the individuals (rows) being retained, but even when we reach $p = 300$, only around 5% of rows will have no missing entries.

The inadequacy of the complete-case approach in many applications has motivated numerous methodological developments in the field of missing data over the past 60 years or so, including imputation (Ford, 1983; Rubin, 2004), factored likelihood (Anderson, 1957) and Expectation-Maximisation approaches (Dempster, Laird and Rubin, 1977); see, e.g., Little and Rubin (2014) for an introduction to the area. Recent years have also witnessed increasing emphasis on understanding the performance of methods for dealing with missing data in a variety of high-dimensional problems, including sparse regression (Loh and Wainwright, 2012; Belloni, Rosenbaum and Tsybakov, 2017), classification (Cai and Zhang, 2018b), sparse principal component analysis (Elsener and van de Geer, 2018) and covariance and precision matrix estimation (Lounici, 2014; Loh and Tan, 2018).

In this paper, we study the effects of missing data in one of the canonical problems of high-dimensional data analysis, namely dimension reduction via Principal Component Analysis (PCA). This is closely related to the topic of *matrix completion*, which has received a great deal of attention in the literature over the last decade or so (e.g. Candès and Recht, 2009; Candès and Plan, 2010; Keshavan, Montanari and Oh, 2010; Mazumder, Hastie and Tibshirani, 2010; Koltchinskii, Lounici and Tsybakov, 2011; Candès et al., 2011; Negahban and Wainwright, 2012). There, the focus is typically on accurate recovery of the missing entries, subject to a low-rank assumption on the signal matrix; by contrast, our focus is on estimation of the principal eigenspaces. Previously proposed methods for low-dimensional PCA with missing data include non-linear iterative partial least squares (Wold and Lytken, 1969), iterative PCA (Kiers, 1997; Josse and Husson, 2012) and its regularised variant (Josse et al., 2009); see Dray and Josse (2015) for a nice survey and comparative study. More broadly, the R-miss-tastic website <https://rmisstastic.netlify.com/> provides a valuable resource on methods for handling missing data.

The importance of the problem of high-dimensional PCA with missing data derives from its manifold applications. For instance, in many commercial settings, one may have a matrix of customers and products, with entries recording the number of purchases. Naturally, there

will typically be a high proportion of missing entries. Nevertheless, PCA can be used to identify items that distinguish the preferences of customers particularly effectively, to make recommendations to users of products they might like and to summarise efficiently customers' preferences. Later, we will illustrate such an application, on the Million Song Dataset, where we are able to identify particular songs that have substantial discriminatory power for users' preferences as well as other interesting characteristics of the user database. Other potential application areas include health data, where one may seek features that best capture the variation in a population, and where the corresponding principal component scores may be used to cluster individuals into subgroups (that may, for instance, receive different treatment regimens).

To formalise the problem we consider, suppose that the (partially observed) matrix $\mathbf{Y} \in \mathbb{R}^{n \times d}$ is of the form

$$\mathbf{Y} = \mathbf{X} + \mathbf{Z}, \tag{1}$$

for independent random matrices \mathbf{X} and \mathbf{Z} , where \mathbf{X} is a low-rank matrix and \mathbf{Z} is a noise matrix with independent and identically distributed subgaussian entries having zero mean and unit variance. The low-rank property of \mathbf{X} is encoded through the assumption that it is generated via

$$\mathbf{X} = \mathbf{U}\mathbf{V}_K^\top, \tag{2}$$

where $\mathbf{V}_K \in \mathbb{R}^{d \times K}$ has orthonormal columns and \mathbf{U} is a random $n \times K$ matrix (with $n > K$) having independent and identically distributed rows with mean zero.

We are interested in estimating the column space of \mathbf{V}_K , denoted by $\text{Col}(\mathbf{V}_K)$, which is also the K -dimensional leading eigenspace of $\Sigma_{\mathbf{y}} := n^{-1}\mathbb{E}\mathbf{Y}^\top\mathbf{Y}$. [Cho, Kim and Rohe \(2017\)](#) considered a different but related model where \mathbf{U} in (2) is deterministic, and is not necessarily centred, so that \mathbf{V}_K is the top K right singular space of $\mathbb{E}(\mathbf{Y})$. (By contrast, in our setting, $\mathbb{E}(\mathbf{Y}) = \mathbf{0}$, so the mean structure is uninformative for recovering \mathbf{V}_K .) In the context of *p-homogeneous missingness*, where each entry of \mathbf{Y} is observed independently with probability $p \in (0, 1)$, [Cho, Kim and Rohe \(2017\)](#) proposed to estimate $\text{Col}(\mathbf{V}_K)$ by $\text{Col}(\widehat{\mathbf{V}}_K)$, where $\widehat{\mathbf{V}}_K$ is a simple estimator formed as the top K eigenvectors of an inverse-probability weighted (IPW) version of the sample covariance matrix (here, the weighting is designed to achieve approximate unbiasedness). Our first contribution, in [Section 2](#), is to provide a detailed, finite-sample analysis of this estimator in the model given by (1) and (2), with a noise level of constant order. The differences between the settings necessitate completely different arguments, and reveal in particular a new phenomenon in the form of a phase transition in the attainable risk bound for the $\sin \Theta$ loss function, i.e. the Frobenius norm of the diagonal matrix of the sines of the principal angles between $\widehat{\mathbf{V}}_K$ and \mathbf{V}_K . Moreover, we also provide a minimax lower bound in the case of estimating a single principal component, which reveals that this estimator achieves the minimax optimal rate up to a poly-logarithmic factor.

While this appears to be a very encouraging story for the IPW estimator, it turns out that it is really only the starting point for a more complete understanding of high-dimensional PCA. For instance, in the noiseless case, the IPW estimator fails to provide exact recovery of the principal components. Moreover, it is the norm rather than the exception in applications

that the missingness mechanism is *heterogeneous*, in the sense that the probability of observing entries of \mathbf{Y} varies (often significantly) across columns. For instance, in recommendation systems, some products will typically be more popular than others, and hence we observe more ratings in those columns. As another example, in meta-analyses of data from several studies, it is frequently the case that some covariates are common across all studies, while others appear only in a reduced proportion of them. In Section 2.2, we present an example to show that PCA algorithms can break down entirely for such heterogeneous missingness mechanisms when individual coordinates in \mathbf{V}_K may be large in absolute value. Intuitively, if we do not observe the interaction between the j th and k th columns of \mathbf{Y} , then we cannot hope to estimate the j th or k th rows of \mathbf{V}_K , and this will cause substantial error if these rows of \mathbf{V}_K contain significant signal. This example illustrates that it is only possible to handle heterogeneous missingness in high-dimensional PCA with additional structure, and indicates that it is natural to assume *incoherence* among the entries of \mathbf{V}_K ; i.e. no single coordinate of \mathbf{V}_K is too large in absolute value.

Our main contribution, then, is to propose a new, iterative algorithm, called **primePCA** (short for projected refinement for imputation of missing entries in Projectal Component Analysis), in Section 3, to estimate \mathbf{V}_K under this incoherence assumption, even with heterogeneous missingness. The initialiser for this algorithm is a modified version of the simple estimator discussed above, where the modification accounts for potential heterogeneity. Each iteration of **primePCA** projects the observed entries of \mathbf{Y} onto the column space of the current estimate of \mathbf{V}_K to impute missing entries, and then updates our estimate of \mathbf{V}_K by computing the leading right singular space of the imputed data matrix. We show that in the noiseless setting, i.e., $\mathbf{Z} = \mathbf{0}$, **primePCA** achieves exact recovery of the principal eigenspaces (with a geometric convergence rate) when the initial estimator is close to the truth and a sufficiently large proportion of the data are observed. Moreover, we also provide a performance guarantee for the initial estimator, showing that it satisfies the desired requirement with high probability, conditional on any observed missingness pattern. Code for our algorithm is available in the R package **primePCA** (Zhu, Wang and Samworth, 2019).

To the best of our knowledge, **primePCA** is the first method for high-dimensional PCA that is designed to cope with settings where the missingness mechanism is heterogeneous. Indeed, the previously mentioned works on high-dimensional PCA and other high-dimensional statistical problems with missing data have either focused on a uniform missingness setting or have imposed a lower bound on entrywise observation probabilities, which reduces to this uniform case. A key contribution of our work is to account explicitly for the effect of a heterogeneous missingness mechanism, where the estimation error depends on average entrywise missingness rather than worst-case missingness; see the discussions after Theorem 4 and Proposition 2 below. In Section 4, the empirical performance of **primePCA** is compared both with that of the initialiser, and with a popular method for matrix completion called **softImpute** (Mazumder, Hastie and Tibshirani, 2010; Hastie et al., 2015), which solves a nuclear-norm regularised optimisation problem, and which can be adapted to provide an estimate of \mathbf{V}_K . It turns out that in many settings, **primePCA** outperforms the **softImpute** algorithm, even when the latter is allowed access to the oracle choice of regularisation pa-

parameter for each dataset. Our analysis of the Million Song Dataset is given in Section 5. In Section 6, we illustrate how some of the ideas in this work may be applied to other high-dimensional statistical problems involving missing data. Proofs of our main results are deferred to Section 7; auxiliary results and their proofs are given in Section 8.

1.1 Notation

For a positive integer T , we write $[T] := \{1, \dots, T\}$. For $\mathbf{v} = (v_1, \dots, v_d)^\top \in \mathbb{R}^d$ and $p \in [1, \infty)$, we define $\|\mathbf{v}\|_p := (\sum_{j=1}^d |v_j|^p)^{1/p}$ and $\|\mathbf{v}\|_\infty := \max_{j \in [d]} |v_j|$. We let $\mathcal{S}^{d-1} := \{\mathbf{u} \in \mathbb{R}^d : \|\mathbf{u}\|_2 = 1\}$ denote the unit Euclidean sphere in \mathbb{R}^d . Given $\mathbf{u} = (u_1, \dots, u_d)^\top, \mathbf{v} = (v_1, \dots, v_d)^\top \in \mathbb{R}^d$, we denote their Hamming distance by $d_H(\mathbf{u}, \mathbf{v}) := \sum_{j=1}^d \mathbf{1}_{\{u_j \neq v_j\}}$. Also, we write $\mathbf{e}_j \in \mathbb{R}^d, j \in [d]$, for the standard basis vector along the j th coordinate axis and $\mathbf{1}_d$ for the all-one vector in \mathbb{R}^d .

Given $\mathbf{u} = (u_1, \dots, u_d)^\top \in \mathbb{R}^d$, we write $\text{diag}(\mathbf{u}) \in \mathbb{R}^{d \times d}$ for the diagonal matrix whose j th diagonal entry is u_j . We let $\mathbb{S}^{d \times d}$ denote the set of symmetric matrices in $\mathbb{R}^{d \times d}$ and let $\mathbb{O}^{d_1 \times d_2}$ denote the set of matrices in $\mathbb{R}^{d_1 \times d_2}$ with orthonormal columns. For a matrix $\mathbf{A} = (A_{ij}) \in \mathbb{R}^{d_1 \times d_2}$, and $p, q \in [1, \infty]$, we write $\|\mathbf{A}\|_p := (\sum_{i,j} |A_{ij}|^p)^{1/p}$ if $1 \leq p < \infty$ and $\|\mathbf{A}\|_\infty := \max_{i,j} |A_{ij}|$ for its entrywise ℓ_p norm, and $\|\mathbf{A}\|_{p \rightarrow q} := \sup_{\|\mathbf{v}\|_p=1} \|\mathbf{A}\mathbf{v}\|_q$ for its p -to- q operator norm, where \mathbf{A} is viewed as a representation of a linear map from $(\mathbb{R}^{d_1}, \|\cdot\|_p)$ to $(\mathbb{R}^{d_2}, \|\cdot\|_q)$. We provide special notation for the (Euclidean) operator norm and the Frobenius norm by writing $\|\mathbf{A}\|_{\text{op}} := \|\mathbf{A}\|_{2 \rightarrow 2}$ and $\|\mathbf{A}\|_F := \|\mathbf{A}\|_2$, respectively, and also write $\|\mathbf{A}\|_*$ for the nuclear norm. If $\mathbf{A} \in \mathbb{S}^{d \times d}$ has the eigendecomposition $\mathbf{A} = \mathbf{Q} \text{diag}(\mu_1, \dots, \mu_d) \mathbf{Q}^\top$ for some $\mathbf{Q} \in \mathbb{O}^{d \times d}$ and $\mu_1 \geq \dots \geq \mu_d$, we write $\lambda_k(\mathbf{A}) := \mu_k$ for its k th largest eigenvalue and abuse terminology slightly to refer to the leftmost k columns of \mathbf{Q} as the top k eigenvectors of \mathbf{A} when $\mu_k > \mu_{k+1}$. Also, we write $|\mathbf{A}| := \mathbf{Q} \text{diag}(|\mu_1|, \dots, |\mu_d|) \mathbf{Q}^\top$. For any $\mathbf{A} \in \mathbb{R}^{n \times d}$ with singular value decomposition $\mathbf{A} = \mathbf{U} \text{diag}(\mu_1, \dots, \mu_r) \mathbf{V}^\top$, where $\mathbf{U} \in \mathbb{O}^{n \times r}$, $\mathbf{V} \in \mathbb{O}^{d \times r}$ and $\mu_1 \geq \dots \geq \mu_r > 0$, we write $\sigma_k(\mathbf{A}) := \mu_k$ for its k th largest singular value and refer to the leftmost k columns of \mathbf{U} (resp. \mathbf{V}) as the top k left (resp. right) singular vectors of \mathbf{A} . The Moore–Penrose pseudoinverse of \mathbf{A} is defined as $\mathbf{A}^\dagger := \mathbf{V} \text{diag}(\mu_1^{-1}, \dots, \mu_r^{-1}) \mathbf{U}^\top$. If $S \subseteq [n]$, we write $\mathbf{A}_S \in \mathbb{R}^{|S| \times d}$ for the matrix obtained by extracting the rows of \mathbf{A} that are in S .

For two matrices $\mathbf{A}, \mathbf{B} \in \mathbb{S}^{d \times d}$, we write $\mathbf{A} \preceq \mathbf{B}$ if $\mathbf{B} - \mathbf{A}$ is positive semidefinite. The Kronecker product of two matrices $\mathbf{A} = (A_{ij}) \in \mathbb{R}^{d_1 \times d_2}$ and $\mathbf{B} = (B_{ij}) \in \mathbb{R}^{d'_1 \times d'_2}$ is defined as the block matrix

$$\mathbf{A} \otimes \mathbf{B} := \begin{pmatrix} A_{11} \mathbf{B} & \cdots & A_{1d_2} \mathbf{B} \\ \vdots & \ddots & \vdots \\ A_{d_1 1} \mathbf{B} & \cdots & A_{d_1 d_2} \mathbf{B} \end{pmatrix} \in \mathbb{R}^{d_1 d'_1 \times d_2 d'_2}.$$

When $d'_1 = d_1$ and $d'_2 = d_2$, the Hadamard product of \mathbf{A} and \mathbf{B} , denoted $\mathbf{A} \circ \mathbf{B}$, is defined such that $(\mathbf{A} \circ \mathbf{B})_{ij} = A_{ij} B_{ij}$ for any $i \in [d_1]$ and $j \in [d_2]$. Moreover, we say that \mathbf{B} is the *Hadamard inverse* of \mathbf{A} if $\mathbf{A} \circ \mathbf{B} = \mathbf{1}_{d_1} \mathbf{1}_{d_2}^\top$.

If $\mathbf{U}, \mathbf{V} \in \mathbb{O}^{d \times K}$, then the matrix of principal angles between $\text{Col}(\mathbf{U})$ and $\text{Col}(\mathbf{V})$ is given

by $\Theta(\mathbf{U}, \mathbf{V}) := \text{diag}(\cos^{-1}(\sigma_1), \dots, \cos^{-1}(\sigma_K))$, where $\sigma_j = \sigma_j(\mathbf{U}^\top \mathbf{V})$; we let $\sin \Theta(\mathbf{U}, \mathbf{V})$ be defined entrywise. We define the loss function

$$L(\mathbf{U}, \mathbf{V}) := \|\sin \Theta(\mathbf{U}, \mathbf{V})\|_{\mathbb{F}}.$$

For a real-valued random variable X and $r \in \mathbb{N}$, we define its (Orlicz) ψ_r -norm by

$$\|X\|_{\psi_r} := \sup_{q \in \mathbb{N}} q^{-1/r} (\mathbb{E}|X|^q)^{1/q}.$$

For a random vector \mathbf{x} taking values in \mathbb{R}^d and $r \geq 1$, we define its (Orlicz) ψ_r -norm by

$$\|\mathbf{x}\|_{\psi_r} := \sup_{\mathbf{u} \in \mathcal{S}^{d-1}} \|\mathbf{u}^\top \mathbf{x}\|_{\psi_r},$$

and define a version that is invariant to invertible affine transformations by

$$\|\mathbf{x}\|_{\psi_r^*} := \sup_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{\|\mathbf{u}^\top (\mathbf{x} - \mathbb{E}\mathbf{x})\|_{\psi_r}}{\text{Var}^{1/2}(\mathbf{u}^\top \mathbf{x})}.$$

We say that a d -dimensional random vector \mathbf{x} is *sub-Gaussian* if $\|\mathbf{x}\|_{\psi_2^*} < \infty$. For two distributions P_1 and P_2 defined on the same measurable space $(\mathcal{X}, \mathcal{A})$ and such that P_1 is absolutely continuous with respect to P_2 , the Kullback–Leibler divergence from P_2 to P_1 is given by

$$\text{KL}(P_1, P_2) := \int_{\mathcal{X}} \log \frac{dP_1}{dP_2} dP_1.$$

Finally, for $a, b \geq 0$, we write $a \lesssim b$ if there exists a universal constant $C > 0$ such that $a \leq Cb$, and, where a and b may depend on an additional variable x , say, we write $a \lesssim_x b$ if there exists $C > 0$, depending only on x , such that $a \leq Cb$.

2 The inverse-probability weighted estimator

For notational simplicity, we will write $\lambda_j := \lambda_j(\boldsymbol{\Sigma}_{\mathbf{u}})$ throughout the paper. Let \mathcal{A}_{ij} denote the event that the (i, j) th entry Y_{ij} of \mathbf{Y} is observed. We define the revelation matrix $\boldsymbol{\Omega} = (\omega_{ij}) \in \mathbb{R}^{n \times d}$ by $\omega_{ij} := \mathbb{1}_{\mathcal{A}_{ij}}$ and the partially observed data matrix

$$\mathbf{Y}_{\boldsymbol{\Omega}} := \mathbf{Y} \circ \boldsymbol{\Omega}. \tag{3}$$

In this section, we consider the simple case, where entries of the data matrix \mathbf{Y} are observed independently and completely at random (i.e., independent of (\mathbf{U}, \mathbf{Z})) with p -homogeneous missingness probability. Thus, $\mathbb{P}(\mathcal{A}_{ij}) = p \in (0, 1)$ for all $i \in [n], j \in [d]$, and \mathcal{A}_{ij} and $\mathcal{A}_{i'j'}$ are independent for $(i, j) \neq (i', j')$.

For $i \in [n]$, let \mathbf{y}_i^\top and $\boldsymbol{\omega}_i^\top$ denote the i th rows of \mathbf{Y} and $\boldsymbol{\Omega}$ respectively, and define $\tilde{\mathbf{y}}_i := \mathbf{y}_i \circ \boldsymbol{\omega}_i$. Writing $\mathbf{P} := \mathbb{E}\boldsymbol{\omega}_1 \boldsymbol{\omega}_1^\top$ and \mathbf{W} for its Hadamard inverse, we have that under the

p -homogeneous missingness mechanism, $\mathbf{P} = p^2\{\mathbf{1}_d\mathbf{1}_d^\top - (1 - p^{-1})\mathbf{I}_d\}$ and $\mathbf{W} = p^{-2}\{\mathbf{1}_d\mathbf{1}_d^\top - (1 - p)\mathbf{I}_d\}$. Following [Cho, Kim and Rohe \(2017\)](#), we consider the following weighted sample covariance matrix:

$$\mathbf{G} := \left(\frac{1}{n}\mathbf{Y}_\Omega^\top\mathbf{Y}_\Omega\right) \circ \mathbf{W} = \left(\frac{1}{n}\sum_{i=1}^n\tilde{\mathbf{y}}_i\tilde{\mathbf{y}}_i^\top\right) \circ \mathbf{W}.$$

The reason for including the weight \mathbf{W} is to ensure that $\mathbb{E}(\mathbf{G}|\mathbf{Y}) = n^{-1}\mathbf{Y}^\top\mathbf{Y}$, so that \mathbf{G} is an unbiased estimator of $\Sigma_{\mathbf{y}}$. Related ideas appear in the work of [Cai and Zhang \(2018a\)](#) on high-dimensional covariance matrix estimation with missing data. There, the authors propose a ‘generalised sample covariance matrix’, where the covariance between any two dimensions j and k is estimated using only the observations for which both dimensions j and k were observed. In practice, p is typically unknown and needs to be estimated. It is thus natural to consider the following plug-in estimator $\widehat{\mathbf{G}}$:

$$\widehat{\mathbf{G}} = \left(\frac{1}{n}\mathbf{Y}_\Omega^\top\mathbf{Y}_\Omega\right) \circ \widehat{\mathbf{W}}, \tag{4}$$

where $\widehat{\mathbf{W}} = \widehat{p}^{-2}\{\mathbf{1}_d\mathbf{1}_d^\top - (1 - \widehat{p})\mathbf{I}_d\}$ and $\widehat{p} := (nd)^{-1}\|\Omega\|_1$ denotes the proportion of observed entries in \mathbf{Y} . We let $\widehat{\mathbf{V}}_K$ denote the top K eigenvectors of $\widehat{\mathbf{G}}$.

2.1 Theory for homogeneous missingness

In order to describe our theoretical performance guarantee for $\widehat{\mathbf{V}}_K$, we first list our conditions on the underlying data generating mechanism. We assume that $(\mathbf{Y}_\Omega, \Omega)$ is generated according to (1), (2) and (3), where:

- (A1) \mathbf{U} , \mathbf{Z} and Ω are independent;
- (A2) \mathbf{U} has independent and identically distributed rows $(\mathbf{u}_i : i \in [n])$ with $\mathbb{E}\mathbf{u}_1 = 0$ and $\|\mathbf{u}_1\|_{\psi_2^*} \leq \tau$;
- (A3) $\mathbf{Z} = (z_{ij})_{i \in [n], j \in [d]}$ has independent and identically distributed entries with $\mathbb{E}z_{11} = 0$, $\text{Var } z_{11} = 1$ and $\|z_{11}\|_{\psi_2^*} \leq \tau$;
- (A4) $\|y_{1j}^2\|_{\psi_1} \leq M$ for all $j \in [d]$;
- (A5) Ω has independent Bern(p) entries.

In many places in this work, we will think intuitively of τ and M as constants. In particular, if \mathbf{U} has multivariate normal rows and \mathbf{Z} has normal entries, then we can simply take $\tau = 1$. For M , under the same normality assumptions, we have $\|y_{1j}^2\|_{\psi_1} = \text{Var}(y_{1j})$, so this intuition amounts to thinking of the variance of each component of our data as being of constant order.

The theorem below gives bounds on the expected proximity between $\text{Col}(\widehat{\mathbf{V}}_K)$ and $\text{Col}(\mathbf{V}_K)$.

Theorem 1. Assume (A1)–(A5) and that $n, d \geq 2$, $dp \geq 1$. Write $R := \lambda_1 + 1$. Then there exists a universal constant $C > 0$ such that

$$\mathbb{E}L(\widehat{\mathbf{V}}_K, \mathbf{V}_K) \leq \frac{CK^{1/2}}{\lambda_{Kp}} \left\{ \left(\frac{Md(R\tau^2p + M \log d) \log^2 d}{n} \right)^{1/2} + \frac{Md \log^2 d \log n}{n} \right\}. \quad (5)$$

In particular, if $n \geq d \log^2 d \log^2 n / (\lambda_1 p + \log d)$, then there exists $C_{M,\tau} > 0$, depending only on M and τ , such that

$$\mathbb{E}L(\widehat{\mathbf{V}}_K, \mathbf{V}_K) \leq \frac{C_{M,\tau}}{\lambda_{Kp}} \left(\frac{Kd(\lambda_1 p + \log d) \log^2 d}{n} \right)^{1/2}. \quad (6)$$

When M, τ are regarded as constants, Theorem 2 below shows that (6) is the minimax rate up to logarithmic factors when $K = 1$. Note that the condition $n \geq d \log^2 d \log^2 n / (\lambda_1 p + \log d)$ is reasonable given the scaling requirement for consistency of the empirical eigenvectors (Shen et al., 2016; Wang and Fan, 2017; Johnstone and Lu, 2009). Indeed, Theorem 5.1 in Shen et al. (2016) shows that when $\lambda_1 \gg 1$, the top eigenvector of the sample covariance matrix estimator is consistent if and only if $d/(n\lambda_1) \rightarrow 0$. If we regard np as the effective sample size in our missing data PCA problem, then it is a sensible analogy to assume $d/(np\lambda_1) \rightarrow 0$ here, which implies that the condition $n \geq d \log^2 d \log^2 n / (\lambda_1 p + \log d)$ holds for large n , up to poly-logarithmic factors.

As mentioned in the introduction, Cho, Kim and Rohe (2017) considered the different but related problem of singular space estimation in a model in which $\mathbf{Y} = \Theta + \mathbf{Z}$, where Θ is a matrix of the form $\mathbf{U}\mathbf{V}_K^\top$ for a *deterministic* matrix \mathbf{U} , whose rows are not necessarily centred. In this setting \mathbf{V}_K is the leading K -dimensional right singular space of Θ , and the same estimator $\widehat{\mathbf{V}}_K$ can be applied. An important distinction is that, when the rows of \mathbf{U} are not centred and the entries of Θ are of comparable magnitude, $\|\Theta\|_F$ is of order \sqrt{nd} , so when K is regarded as a constant, it is natural to think of the singular values of Θ as also being of order \sqrt{nd} . Indeed, this is assumed in Cho, Kim and Rohe (2017). On the other hand, in our model, where the rows of \mathbf{U} have mean zero, assuming that the eigenvalues are of order \sqrt{nd} would amount to an extremely strong requirement, essentially restricting attention to very highly spiked covariance matrices. Removing this condition requires completely different arguments. Moreover, (6) reveals an interesting phase transition phenomenon that has not been observed previously in the literature. Specifically, if the signal strength is large enough that $\lambda_1 \geq p^{-1} \log d$, then we should regard np as the effective sample size, as might intuitively be expected. On the other hand, if $\lambda_1 < p^{-1} \log d$, then the estimation problem is considerably more difficult and the effective sample size is of order np^2 . In fact, by inspecting the proof of Theorem 1, we see that in the high signal case, it is the difficulty of estimating the diagonal entries of $\Sigma_{\mathbf{y}}$ that drives the rate, while when the signal strength is low, the bottleneck is the challenge of estimating the off-diagonal entries. By comparing (6) with the minimax lower bound result in Theorem 2 below, we see that this phase transition phenomenon is an inherent feature of this estimation problem rather than an artefact of the proof techniques we used to derive the upper bound.

In order to state our minimax lower bound, we let $\mathcal{P}_{n,d}(\lambda_1, p)$ denote the class of distributions of pairs $(\mathbf{Y}_\Omega, \Omega)$ satisfying (A1), (A2), (A3) and (A5) with $K = 1$. Since we are now working with vectors instead of matrices, we write \mathbf{v} in place of \mathbf{V}_1 .

Theorem 2. *There exists a universal constant $c > 0$ such that*

$$\inf_{\hat{\mathbf{v}}} \sup_{P \in \mathcal{P}_{n,d}(\lambda_1, p)} \mathbb{E}_P L(\hat{\mathbf{v}}, \mathbf{v}) \geq c \min \left\{ \frac{1}{\lambda_1 p} \left(\frac{d(\lambda_1 p + 1)}{n} \right)^{1/2}, 1 \right\},$$

where the infimum is taken over all estimators $\hat{\mathbf{v}} = \hat{\mathbf{v}}(\mathbf{Y}_\Omega, \Omega)$ of \mathbf{v} .

Theorem 2 reveals that $\hat{\mathbf{V}}_1$ in Theorem 1 achieves the minimax optimal rate of estimation up to a logarithmic factor when M and τ are regarded as constants and $K = 1$.

2.2 General observation mechanism

Although Section 2.1 may appear to indicate that the problem of high-dimensional PCA with missing entries is essentially solved, the aim of this subsection is to show that the situation changes dramatically once the data can be missing heterogeneously.

To this end, consider the following example. Suppose that $\boldsymbol{\omega}$ is equal to $(1, 0, 1, \dots, 1)^\top$ or $(0, 1, 1, \dots, 1)^\top$ with equal probability, so that

$$\mathbf{P} = \mathbb{E} \boldsymbol{\omega} \boldsymbol{\omega}^\top = \begin{pmatrix} 1/2 & 0 & 1/2 & \dots & 1/2 \\ 0 & 1/2 & 1/2 & \dots & 1/2 \\ 1/2 & 1/2 & 1 & \dots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1/2 & 1/2 & 1 & \dots & 1 \end{pmatrix} \in \mathbb{R}^{d \times d}.$$

In other words, for each $i \in [n]$, we observe precisely one of the first two entries of \mathbf{y}_i , together with all of the remaining $(d - 2)$ entries. Let $\boldsymbol{\Sigma} = \mathbf{I}_d + \boldsymbol{\alpha} \boldsymbol{\alpha}^\top$, where $\boldsymbol{\alpha} = (2^{-1/2}, 2^{-1/2}, 0, \dots, 0)^\top \in \mathbb{R}^d$, and $\boldsymbol{\Sigma}' = \mathbf{I}_d + \boldsymbol{\alpha}' (\boldsymbol{\alpha}')^\top$, where $\boldsymbol{\alpha}' = (2^{-1/2}, -2^{-1/2}, 0, \dots, 0)^\top \in \mathbb{R}^d$. Suppose that $\mathbf{y} \sim N_d(\mathbf{0}, \boldsymbol{\Sigma})$ and let $\tilde{\mathbf{y}} := \mathbf{y} \circ \boldsymbol{\omega}$, and similarly assume that $\mathbf{y}' \sim N_d(\mathbf{0}, \boldsymbol{\Sigma}')$ and set $\tilde{\mathbf{y}}' := \mathbf{y}' \circ \boldsymbol{\omega}$. Then $(\tilde{\mathbf{y}}, \boldsymbol{\omega})$ and $(\tilde{\mathbf{y}}', \boldsymbol{\omega})$ are identically distributed. However, the leading eigenvectors of $\boldsymbol{\Sigma}$ and $\boldsymbol{\Sigma}'$ are respectively $\boldsymbol{\alpha}$ and $\boldsymbol{\alpha}'$, which are orthogonal! Thus, it is impossible to simultaneously estimate consistently the leading eigenvector of both $\boldsymbol{\Sigma}$ and $\boldsymbol{\Sigma}'$ from our observations. We note that it is the disproportionate weight of the first two coordinates in the leading eigenvector, combined with the failure to observe simultaneously the first two entries in the data, that makes the estimation problem intractable in this example.

The understanding derived from this example motivates us to assume that \mathbf{V}_K satisfies the incoherence condition $\|\mathbf{V}_K\|_\infty \leq \mu/\sqrt{d}$ for some $\mu > 0$. The intuition here is that the maximally incoherent case is where each column of \mathbf{V}_K is a unit vector proportional to a vector whose entries are either 1 or -1 , in which case $\|\mathbf{V}_K\|_\infty = 1/\sqrt{d}$. Our condition asks for the columns of \mathbf{V}_K to have an incoherence of the same order as this maximally

incoherent case. Similar conditions have been invoked in the literature on matrix completion (e.g., Candès and Plan, 2010; Keshavan, Montanari and Oh, 2010), but for a different reason. There, the purpose is to ensure that the true right singular space is not too closely aligned with the standard basis, which allows the missing entries of the matrix to be inferred from relatively few observations. In our case, the incoherence condition ensures that significant estimation error in a few components of the leading eigenvectors does not affect the overall statistical performance too much. Thus, it rules out examples such as the one described above, where heavy corruption in only a few entries spoils any chance of consistent estimation.

3 Our new algorithm, primePCA

We are now in a position to introduce and analyse our iterative algorithm **primePCA** to estimate principal eigenspaces of the covariance matrix $\Sigma_{\mathbf{y}}$. The basic idea is to iteratively refine a current (input) iterate $\widehat{\mathbf{V}}_K^{(\text{in})}$ by first imputing the missing entries of the data matrix \mathbf{Y}_Ω using the current estimate of \mathbf{V}_K , and then applying a singular value decomposition (SVD) to the completed data matrix. More precisely, for $i \in [n]$, we let \mathcal{J}_i denote the indices for which the corresponding entry of \mathbf{y}_i is observed, and regress the observed data $\widetilde{\mathbf{y}}_{i,\mathcal{J}_i} = \mathbf{y}_{i,\mathcal{J}_i}$ on $(\widehat{\mathbf{V}}_K^{(\text{in})})_{\mathcal{J}_i}$ to obtain an estimate $\widehat{\mathbf{u}}_i$ of the i th row of \mathbf{U} . This is natural in view of the data generating mechanism $\mathbf{y}_i = \mathbf{V}_K \mathbf{u}_i + \mathbf{z}_i$. We then use $\widehat{\mathbf{y}}_{i,\mathcal{J}_i^c} := (\widehat{\mathbf{V}}_K^{(\text{in})} \widehat{\mathbf{u}}_i)_{\mathcal{J}_i^c}$ to impute the missing values $\mathbf{y}_{i,\mathcal{J}_i^c}$, retain the original observed entries as $\widehat{\mathbf{y}}_{i,\mathcal{J}_i} := \mathbf{y}_{i,\mathcal{J}_i}$, and set our next (output) iterate $\widehat{\mathbf{V}}_K^{(\text{out})}$ to be the top K right singular vectors of the imputed matrix $\widehat{\mathbf{Y}} := (\widehat{\mathbf{y}}_1, \dots, \widehat{\mathbf{y}}_n)^\top$. To motivate this final choice, observe that when $\mathbf{Z} = \mathbf{0}$, we have $\text{rank}(\mathbf{Y}) = K$; we therefore have the SVD $\mathbf{Y} = \mathbf{L}\mathbf{\Gamma}\mathbf{R}^\top$, where $\mathbf{L} \in \mathbb{O}^{n \times K}$, $\mathbf{R} \in \mathbb{O}^{d \times K}$ and $\mathbf{\Gamma} \in \mathbb{R}^{K \times K}$ is diagonal with positive diagonal entries. This means that $\mathbf{R} = \mathbf{V}_K \mathbf{U}^\top \mathbf{\Gamma}^{-1}$, so the column spaces of \mathbf{R} and \mathbf{V}_K coincide. For convenience, pseudocode of a single iteration of refinement in this algorithm is given in Algorithm 1.

Algorithm 1 `refine`($K, \widehat{\mathbf{V}}_K^{(\text{in})}, \Omega, \mathbf{Y}_\Omega$), a single step of refinement of current iterate $\widehat{\mathbf{V}}_K^{(\text{in})}$

Input: $K \in [d]$, $\widehat{\mathbf{V}}_K^{(\text{in})} \in \mathbb{O}^{d \times K}$, $\Omega \in \{0, 1\}^{n \times d}$ with $\min_i \|\omega_i\|_1 \geq 1$, $\mathbf{Y}_\Omega \in \mathbb{R}^{n \times d}$

Output: $\widehat{\mathbf{V}}_K^{(\text{out})} \in \mathbb{O}^{d \times K}$

- 1: **for** i in $[n]$ **do**
 - 2: $\mathcal{J}_i \leftarrow \{j \in [d] : \omega_{ij} = 1\}$
 - 3: $\widehat{\mathbf{u}}_i \leftarrow (\widehat{\mathbf{V}}_K^{(\text{in})})_{\mathcal{J}_i}^\dagger \widetilde{\mathbf{y}}_{i,\mathcal{J}_i}$
 - 4: $\widehat{\mathbf{y}}_{i,\mathcal{J}_i^c} \leftarrow \widehat{\mathbf{V}}_K^{(\text{in})} \widehat{\mathbf{u}}_{i,\mathcal{J}_i^c}$
 - 5: $\widehat{\mathbf{y}}_{i,\mathcal{J}_i} \leftarrow \mathbf{y}_{i,\mathcal{J}_i}$
 - 6: **end for**
 - 7: $\widehat{\mathbf{Y}} \leftarrow (\widehat{\mathbf{y}}_1, \dots, \widehat{\mathbf{y}}_n)^\top$
 - 8: $\widehat{\mathbf{V}}_K^{(\text{out})} \leftarrow$ top K right singular vectors of $\widehat{\mathbf{Y}}$
-

We now seek to provide formal justification for Algorithm 1. For any $\mathbf{V}^{(1)}, \mathbf{V}^{(2)} \in \mathbb{O}^{d \times K}$, we let $\mathbf{W}_1 \mathbf{D}_{\mathbf{V}^{(1)}, \mathbf{V}^{(2)}} \mathbf{W}_2^\top$ be an SVD of $(\mathbf{V}^{(2)})^\top \mathbf{V}^{(1)}$ and let $\mathbf{W}_{\mathbf{V}^{(1)}, \mathbf{V}^{(2)}} := \mathbf{W}_1 \mathbf{W}_2^\top$. The two-to-infinity distance between $\mathbf{V}^{(1)}$ and $\mathbf{V}^{(2)}$ is then defined to be

$$\mathcal{T}(\mathbf{V}^{(1)}, \mathbf{V}^{(2)}) := \|\mathbf{V}^{(1)} - \mathbf{V}^{(2)} \mathbf{W}_{\mathbf{V}^{(1)}, \mathbf{V}^{(2)}}\|_{2 \rightarrow \infty}.$$

Note that the definition of $\mathcal{T}(\mathbf{V}^{(1)}, \mathbf{V}^{(2)})$ does not depend on our choice of SVD and that $\mathcal{T}(\mathbf{V}^{(1)}, \mathbf{V}^{(2)}) = \mathcal{T}(\mathbf{V}^{(1)} \mathbf{O}_1, \mathbf{V}^{(2)} \mathbf{O}_2)$ for any $\mathbf{O}_1, \mathbf{O}_2 \in \mathbb{O}^{K \times K}$. The following proposition considers the noiseless setting $\mathbf{Z} = \mathbf{0}$, and provides conditions under which, for any estimator $\widehat{\mathbf{V}}_K^{(\text{in})}$ that is close to \mathbf{V}_K , a single iteration of refinement in Algorithm 1 contracts the two-to-infinity distance between their column spaces. In a slight abuse of notation, we write $\Omega^c := \mathbf{1}_d \mathbf{1}_d^\top - \Omega$.

Proposition 1. *Let $\widehat{\mathbf{V}}_K^{(\text{out})} := \text{refine}(K, \widehat{\mathbf{V}}_K^{(\text{in})}, \Omega, \mathbf{Y}_\Omega)$ as in Algorithm 1 and let $\Delta := \mathcal{T}(\widehat{\mathbf{V}}_K^{(\text{in})}, \mathbf{V}_K)$. We assume that $\min_{i \in [n]} \|\omega_i\|_1 > K$, that $\min_{i \in [n]} \frac{d^{1/2} \sigma_K((\widehat{\mathbf{V}}_K^{(\text{in})})_{\mathcal{J}_i})}{|\mathcal{J}_i|^{1/2}} \geq 1/\sigma_* > 0$, and write the SVD of \mathbf{Y} as \mathbf{LFR}^\top . Suppose that $\mathbf{Z} = \mathbf{0}$, and that both $\|\mathbf{L}\|_{2 \rightarrow \infty} \leq \mu_1 (K/n)^{1/2}$ and $\|\mathbf{R}\|_{2 \rightarrow \infty} \leq \mu_2 (K/d)^{1/2}$ hold for some $\mu_1, \mu_2 \geq 1$. Then there exist $c_1, C > 0$, depending only on μ_1, μ_2 and σ_* , such that whenever*

$$\begin{aligned} (i) \quad & \Delta \leq \frac{c_1 \sigma_K(\Gamma)}{K^2 \sigma_1(\Gamma) \sqrt{d}}, \\ (ii) \quad & \rho := \frac{CK^2 \sigma_1(\Gamma) \|\Omega^c\|_{1 \rightarrow 1}}{\sigma_K(\Gamma) n} < 1, \end{aligned}$$

we have that

$$\mathcal{T}(\widehat{\mathbf{V}}_K^{(\text{out})}, \mathbf{V}_K) \leq \rho \Delta.$$

In practice, in cases where either of the two conditions on $\min_{i \in [n]} \|\omega_i\|_1$ or σ_* is not satisfied, we first perform a screening step that restricts attention to a set of row indices for which the data contain sufficient information to estimate the K principal components. This screening step is explicitly accounted for in Algorithm 2 below, as well as in the theory that justifies it.

Algorithm 2 provides pseudocode for the iterative primePCA algorithm, given an initial estimator $\widehat{\mathbf{V}}_K^{(0)}$. The iterations continue until either we hit the convergence threshold κ^* or the maximum iteration number n_{iter} . Theorem 3 below guarantees that, in the noiseless setting of Proposition 1, the primePCA estimator converges to \mathbf{V}_K at a geometric rate.

Theorem 3. *For $t \in [n_{\text{iter}}]$, let $\widehat{\mathbf{V}}_K^{(t)}$ be the t^{th} iterate of Algorithm 2 with input $K, \widehat{\mathbf{V}}_K^{(0)}, \Omega \in \{0, 1\}^{n \times d}, \mathbf{Y}_\Omega \in \mathbb{R}^{n \times d}, n_{\text{iter}} \in \mathbb{N}, \sigma_* \in (0, \infty)$ and $\kappa^* = 0$. Write $\Delta := \mathcal{T}(\widehat{\mathbf{V}}_K^{(0)}, \mathbf{V}_K)$ and let*

$$\mathcal{I} := \left\{ i : \|\omega_i\|_1 > K, \sigma_K((\mathbf{V}_K)_{\mathcal{J}_i}) \geq \frac{|\mathcal{J}_i|^{1/2}}{d^{1/2} \sigma_*} \right\},$$

Algorithm 2 primePCA, an iterative algorithm for estimating \mathbf{V}_K given initialiser $\widehat{\mathbf{V}}_K^{(0)}$

Input: $K \in [d]$, $\widehat{\mathbf{V}}_K^{(0)} \in \mathbb{O}^{d \times K}$, $\boldsymbol{\Omega} \in \{0, 1\}^{n \times d}$, $\mathbf{Y}_\Omega \in \mathbb{R}^{n \times d}$, $n_{\text{iter}} \in \mathbb{N}$, $\sigma_* \in (0, \infty)$, $\kappa^* \in [0, \infty)$

Output: $\widehat{\mathbf{V}}_K \in \mathbb{R}^{d \times K}$

```

1: for  $i$  in  $[n]$  do
2:    $\mathcal{J}_i \leftarrow \{j \in [d] : \omega_{ij} = 1\}$ 
3: end for
4: for  $t$  in  $[n_{\text{iter}}]$  do
5:    $\mathcal{I}^{(t-1)} \leftarrow \{i : \|\boldsymbol{\omega}_i\|_1 > K, \sigma_K((\widehat{\mathbf{V}}_K^{(t-1)})_{\mathcal{J}_i}) \geq \frac{|\mathcal{J}_i|^{1/2}}{d^{1/2}\sigma_*}\}$ 
6:    $\widehat{\mathbf{V}}_K^{(t)} \leftarrow \text{refine}(K, \widehat{\mathbf{V}}_K^{(t-1)}, \boldsymbol{\Omega}_{\mathcal{I}^{(t-1)}}, (\mathbf{Y}_\Omega)_{\mathcal{I}^{(t-1)}})$  # refine is defined in Algorithm 1.
7:   if  $L(\widehat{\mathbf{V}}_K^{(t)}, \widehat{\mathbf{V}}_K^{(t-1)}) < \kappa^*$  then break
8:   end if
9: end for
10: return  $\widehat{\mathbf{V}}_K = \widehat{\mathbf{V}}_K^{(t)}$ 

```

where $\mathcal{J}_i := \{j : \omega_{ij} = 1\}$. Let $\mathbf{Y}_\mathcal{I} = \mathbf{L}\mathbf{R}^\top$ be an SVD of $\mathbf{Y}_\mathcal{I}$. Suppose that both $\|\mathbf{L}\|_{2 \rightarrow \infty} \leq \mu_1(K/|\mathcal{I}|)^{1/2}$ and $\|\mathbf{R}\|_{2 \rightarrow \infty} \leq \mu_2(K/d)^{1/2}$ hold, for some $\mu_1, \mu_2 \geq 1$. Let

$$\mathcal{Z} := \left\{ \frac{\sigma_K((\mathbf{V}_K)_{\mathcal{J}_i})d^{1/2}}{|\mathcal{J}_i|^{1/2}} : i \in [n], \|\boldsymbol{\omega}_i\|_1 > K \right\},$$

and assume that $\epsilon := \min_{z \in \mathcal{Z}} |z - \sigma_*^{-1}| > 0$. Then there exist $c_1, C > 0$, depending only on μ_1, μ_2, σ_* and ϵ , such that whenever

$$(i) \quad \Delta \leq \frac{c_1 \sigma_K(\mathbf{Y}_\mathcal{I})}{K^2 \sigma_1(\mathbf{Y}_\mathcal{I}) \sqrt{d}},$$

$$(ii) \quad \rho := \frac{CK^2 \sigma_1(\mathbf{Y}_\mathcal{I}) \|\boldsymbol{\Omega}_\mathcal{I}^\epsilon\|_{1 \rightarrow 1}}{\sigma_K(\mathbf{Y}_\mathcal{I}) |\mathcal{I}|} < 1,$$

we have that for every $t \in [n_{\text{iter}}]$,

$$\mathcal{T}(\widehat{\mathbf{V}}_K^{(t)}, \mathbf{V}_K) \leq \rho^t \Delta.$$

3.1 Initialisation

Theorem 3 provides a general guarantee on the performance of primePCA, but relies on finding an initial estimator $\widehat{\mathbf{V}}_K^{(0)}$ that is sufficiently close to the truth \mathbf{V}_K . The aim of this subsection, then, is to propose a simple initialiser and show that it satisfies the requirement of Theorem 3 with high probability, conditional on the missingness pattern.

Consider the following modified weighted sample covariance matrix

$$\widetilde{\mathbf{G}} := \frac{1}{n} \sum_{i=1}^n \widetilde{\mathbf{y}}_i \widetilde{\mathbf{y}}_i^\top \circ \widetilde{\mathbf{W}},$$

where for any $j, k \in [d]$,

$$\widetilde{\mathbf{W}}_{jk} := \begin{cases} \frac{n}{\sum_{i=1}^n \omega_{ij}\omega_{ik}} & \text{if } \sum_{i=1}^n \omega_{ij}\omega_{ik} > 0, \\ 0, & \text{otherwise.} \end{cases}$$

Here, the matrix $\widetilde{\mathbf{W}}$ replaces $\widehat{\mathbf{W}}$ in (4) because we no longer wish to assume homogeneous missingness. We take as our initial estimator of \mathbf{V}_K the matrix of top K eigenvectors of $\widetilde{\mathbf{G}}$, denoted $\widetilde{\mathbf{V}}_K$. Theorem 4 below studies the performance of this initialiser, in terms of its two-to-infinity norm error, as required for application in Theorem 3. We write \mathbb{P}^Ω and \mathbb{E}^Ω for probabilities and expectations conditional on Ω .

Theorem 4. *Assume (A1)–(A4) and that $n, d \geq 2$. Suppose further that $\|\mathbf{V}_K\|_\infty \leq \mu/\sqrt{d}$, that $\sum_{i=1}^n \omega_{ij}\omega_{ik} > 0$ for all j, k and let $R := \lambda_1 + 1$. Then there exist $c_{M,\tau,\mu}, C_{M,\tau,\mu} > 0$, depending only on M, τ and μ , such that for every $\xi > 2$, if*

$$\lambda_K > c_{M,\tau,\mu} \left\{ \left(\frac{\max(\|\widetilde{\mathbf{W}}\|_1, R\|\widetilde{\mathbf{W}}\|_{1 \rightarrow 1}) \xi \log d}{n} \right)^{1/2} + \frac{\xi \|\widetilde{\mathbf{W}}\|_{\text{F}} \log^2 d}{n} \right\}, \quad (7)$$

then

$$\begin{aligned} \mathbb{P}^\Omega \left\{ \mathcal{T}(\widetilde{\mathbf{V}}_K, \mathbf{V}_K) \geq \frac{C_{M,\tau,\mu} K^{3/2} R^{1/2}}{\lambda_K d^{1/2}} \left(1 + \frac{d^{1/2}}{K \lambda_K} \right) \left(\frac{\xi^{1/2} \|\widetilde{\mathbf{W}}\|_{\infty \rightarrow \infty}^{1/2} \log^{1/2} d}{n^{1/2}} + \frac{\xi \|\widetilde{\mathbf{W}}\|_{2 \rightarrow \infty} \log d}{n} \right) \right\} \\ \leq 2(e^{K \log 5} + K + 4) d^{-(\xi-1)} + 2d^{-(\xi-2)}. \end{aligned}$$

As a consequence, writing

$$\mathcal{A} := \left\{ \frac{\sigma_K(\mathbf{Y}_{\mathcal{I}})}{\sigma_1(\mathbf{Y}_{\mathcal{I}})} > \frac{C_{M,\tau,\mu} K^{7/2} R^{1/2}}{c_1 \lambda_K} \left(1 + \frac{d^{1/2}}{K \lambda_K} \right) \left(\frac{\xi^{1/2} \|\widetilde{\mathbf{W}}\|_{\infty \rightarrow \infty}^{1/2} \log^{1/2} d}{n^{1/2}} + \frac{\xi \|\widetilde{\mathbf{W}}\|_{2 \rightarrow \infty} \log d}{n} \right) \right\},$$

where c_1 is as in Theorem 3, we have that

$$\mathbb{P}^\Omega \left(\mathcal{T}(\widetilde{\mathbf{V}}_K, \mathbf{V}_K) > \frac{c_1 \sigma_K(\mathbf{Y}_{\mathcal{I}})}{K^2 \sigma_1(\mathbf{Y}_{\mathcal{I}}) d^{1/2}} \right) \leq 2(e^{K \log 5} + K + 4) d^{-(\xi-1)} + 2d^{-(\xi-2)} + \mathbb{P}^\Omega(\mathcal{A}^c).$$

The first part of Theorem 4 provides a general probabilistic bound for $\mathcal{T}(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)$, after conditioning on the missingness pattern. This allows us, in the second part, to provide a guarantee on the probability with which $\widetilde{\mathbf{V}}_K$ is a good enough initialiser for Theorem 3 to apply. For intuition regarding $\mathbb{P}^\Omega(\mathcal{A}^c)$, consider the p -homogenous missingness setting. In that case, by Lemma 6, typical realisations of $\widetilde{\mathbf{W}}$ have $\|\widetilde{\mathbf{W}}\|_{\infty \rightarrow \infty} = O(d/p^2)$ and $\|\widetilde{\mathbf{W}}\|_{2 \rightarrow \infty} = O(d^{1/2}/p^2)$ when $np^2 \gg \log d$, so in the spiked model where λ_1 and λ_K are both of order d , we expect $\mathbb{P}^\Omega(\mathcal{A}^c)$ to be small.

One of the attractions of our analysis is the fact that we are able to provide bounds that only depend on entrywise missingness probabilities in an average sense, as opposed to worst-case missingness probabilities. The refinements conferred by such bounds are particularly

important when the missingness mechanism is heterogeneous, as typically encountered in practice. The averaging of missingness probabilities can be partially seen in Theorem 4, since $\|\widetilde{\mathbf{W}}\|_{\infty \rightarrow \infty}$ and $\|\widetilde{\mathbf{W}}\|_{2 \rightarrow \infty}$ depend only on the ℓ_1 and ℓ_2 norms of each row of $\widetilde{\mathbf{W}}$, but is even more evident in the proposition below, which gives a probabilistic bound on the original $\sin \Theta$ distance between $\widetilde{\mathbf{V}}_K$ and \mathbf{V}_K .

Proposition 2. *Assume the same conditions as in Theorem 4. Then there exists a universal constant $C > 0$ such that for any $\xi > 1$, if*

$$\lambda_K > C \left\{ \left(\frac{M\tau^2 R \|\widetilde{\mathbf{W}}\|_{1 \rightarrow 1} \xi \log d}{n} \right)^{1/2} + \frac{M \|\widetilde{\mathbf{W}}\|_{\text{op}} \xi \log^2 d}{n} \right\}, \quad (8)$$

then

$$\begin{aligned} \mathbb{P}^\Omega \left\{ L(\widetilde{\mathbf{V}}_K, \mathbf{V}_K) \geq \frac{2^{9/2} e \tau \mu}{\lambda_K} \left(\frac{KMR}{d} \right)^{1/2} \left(\frac{\xi^{1/2} \|\widetilde{\mathbf{W}}\|_1^{1/2} \log^{1/2} d}{n^{1/2}} + \frac{\xi \|\widetilde{\mathbf{W}}\|_F \log d}{n} \right) \right\} \\ \leq (2K + 4) d^{-(\xi-1)}. \end{aligned}$$

In this bound, then, we see that $L(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)$ only depends on $\widetilde{\mathbf{W}}$ through the entrywise ℓ_1 and ℓ_2 norms of the whole matrix. Lemma 6 provides probabilistic control of these norms under the p -homogeneous missingness mechanism. In general, if the rows of Ω are independent and identically distributed, but different covariates are missing with different probabilities, then entries of $\widetilde{\mathbf{W}}$ will concentrate around the reciprocals of the simultaneous observation probabilities of pairs of covariates. As such, for a typical realisation of Ω , our bound in Proposition 2 depends only on the harmonic averages of these simultaneous observation probabilities and their squares. Such an averaging effect ensures that our method is effective in a much wider range of heterogeneous settings than previously allowed in the literature.

4 Simulation studies

In this section, we assess the empirical performance of `primePCA` as proposed in Algorithm 2, with initialiser $\widetilde{\mathbf{V}}_K$ from Section 3.1, and denote the output of this algorithm by $\widehat{\mathbf{V}}_K^{\text{prime}}$. We generate observations according to the model described in (1), (2) and (3) where the rows of the matrix \mathbf{U} are independent $N_d(\mathbf{0}, \Sigma_{\mathbf{u}})$ random vectors, for some $\Sigma_{\mathbf{u}} \succeq \mathbf{0}$. We further generate the observation indicator matrix Ω , independently of \mathbf{U} and \mathbf{Z} , and investigate the following four missingness mechanisms that represent different levels of heterogeneity:

- (H1) Homogeneous: $\mathbb{P}(\omega_{ij} = 1) = 0.05$ for all $i \in [n], j \in [d]$;
- (H2) Mildly heterogeneous: $\mathbb{P}(\omega_{ij} = 1) = P_i Q_j$ for $i \in [n], j \in [d]$, where $P_1, \dots, P_n \stackrel{\text{iid}}{\sim} U[0, 0.2]$ and $Q_1, \dots, Q_d \stackrel{\text{iid}}{\sim} U[0.05, 0.95]$ independently;
- (H3) Highly heterogeneous columns: $\mathbb{P}(\omega_{ij} = 1) = 0.19$ for $i \in [n]$ and all odd $j \in [d]$ and $\mathbb{P}(\omega_{ij} = 1) = 0.01$ for $i \in [n]$ and all even $j \in [d]$.

(H4) Highly heterogeneous rows: $\mathbb{P}(\omega_{ij} = 1) = 0.18$ for $j \in [d]$ and all odd $i \in [n]$ and $\mathbb{P}(\omega_{ij} = 1) = 0.02$ for $j \in [d]$ and all even $i \in [n]$.

In Sections 4.1, 4.2 and 4.3 below, we investigate `primePCA` in noiseless, noisy and misspecified settings respectively. In all cases, the average statistical error was estimated from 100 Monte Carlo repetitions of the experiment. For comparison, we also studied the `softImpute` algorithm (Mazumder, Hastie and Tibshirani, 2010; Hastie et al., 2015), which is considered to be state-of-the-art for matrix completion (Chi, Lu and Chen, 2018). This algorithm imputes the missing entries of \mathbf{Y} by solving the following nuclear-norm-regularised optimisation problem:

$$\hat{\mathbf{Y}}^{\text{soft}} := \operatorname{argmin}_{\mathbf{X} \in \mathbb{R}^{n \times d}} \left\{ \frac{1}{2} \|\mathbf{Y}_\Omega - \mathbf{X}_\Omega\|_F^2 + \lambda \|\mathbf{X}\|_* \right\},$$

where $\lambda > 0$ is to be chosen by the practitioner. The `softImpute` estimator of \mathbf{V}_K is then given by the top K right singular space $\hat{\mathbf{V}}_K^{\text{soft}}$ of $\hat{\mathbf{Y}}^{\text{soft}}$.

Figure 4 presents Monte Carlo estimates of $\mathbb{E}L(\hat{\mathbf{V}}_K^{\text{prime}}, \mathbf{V}_K)$ for different choices of σ_* in two different settings. The first uses the noiseless set-up of Section 4.1, together with missingness mechanism (H1); the second uses the noisy setting of Section 4.2 with parameter $\nu = 20$ and missingness mechanism (H2). We see that the error barely changes when σ_* varies within $[2, 10]$; very similar plots were obtained for different data generation and missingness mechanisms, though we omit these for brevity. For definiteness, we therefore fixed $\sigma_* = 3$ throughout our simulation study.

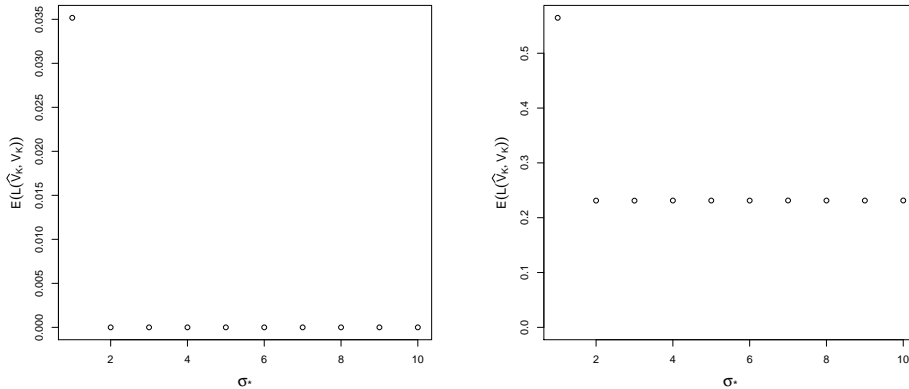


Figure 1: Estimates of $\mathbb{E}L(\hat{\mathbf{V}}_K^{\text{prime}}, \mathbf{V}_K)$ for various choices of σ_* under (H1) in the noiseless setting of Section 4.1 (left) and (H2) in the noisy setting of Section 4.2 with $\nu = 20$ (right).

4.1 Noiseless case

In the noiseless setting, we set $\mathbf{Z} = \mathbf{0}$, and also fix $n = 2000$, $d = 500$, $K = 2$ and $\Sigma_{\mathbf{u}} = 100\mathbf{I}_2$. We set

$$\mathbf{V}_K = \sqrt{\frac{1}{500}} \begin{pmatrix} \mathbf{1}_{250} & \mathbf{1}_{250} \\ \mathbf{1}_{250} & -\mathbf{1}_{250} \end{pmatrix} \in \mathbb{R}^{500 \times 2}.$$

In Figure 2, we present the logarithm of the estimated average loss of `primePCA` and `softImpute` under (H1), (H2), (H3) and (H4). We set the range of y -axis to be the same for each method to facilitate straightforward comparison. We see that the statistical error of `primePCA` decreases geometrically as the number of iterations increases, which confirms the conclusion of Theorem 3 in this noiseless setting. Moreover, after a moderate number of iterations, its performance is a substantial improvement on that of the `softImpute` algorithm, even if this latter algorithm is given access to an oracle choice of the regularisation parameter λ . The high statistical error of `softImpute` in these settings can be partly explained by the default value of the tuning parameter `thresh` in the `softImpute` package in R, namely 10^{-5} , which corresponds to the red curve in the right-hand panels of Figure 2. By reducing the values of `thresh` to 10^{-7} and 10^{-9} , corresponding to the green and blue curves in Figure 2 respectively, we were able to improve the performance of `softImpute` to some extent, though the statistical error is sensitive to the choice of the regularisation parameter λ . Moreover, even with the optimal choice of λ , it is not competitive with `primePCA` (which is also considerably faster to compute, even with 2000 iterations).

4.2 Noisy case

Here, we generate the rows of \mathbf{Z} as independent $N_d(\mathbf{0}, \mathbf{I}_d)$ random vectors, independent of all other data. We maintain the same choices of n , d , K and \mathbf{V}_K as in Section 4.1, set $\Sigma_{\mathbf{u}} = \nu^2 \mathbf{I}_2$ and vary $\nu > 0$ to achieve different signal-to-noise ratios. In particular, defining $\text{SNR} := \text{tr Cov}(\mathbf{x}_1) / \text{tr Cov}(\mathbf{z}_1)$, the choices $\nu = 20, 40, 60$ correspond to the low, medium and high signal-to-noise ratios $\text{SNR} = 1.6, 6.4, 14.4$, respectively. For an additional comparison, we also considered a variant of the `softImpute` algorithm called `hardImpute` (Mazumder, Hastie and Tibshirani, 2010), which retains only a fixed number of top singular values in each iteration of matrix imputation; this can be achieved by setting the argument λ in the `softImpute` function to be 0.

We remark that in general, the choice of λ for `softImpute` is more challenging than in many regularised M -estimation contexts, because in our setting we have no response variable, so cross-validation techniques are less readily available. For our comparisons, therefore, we gave the `softImpute` algorithm a particularly strong form of oracle choice of λ , namely where λ was chosen for each individual repetition of the experiment, so as to minimise the loss function. Naturally, such a choice is not available to the practitioner. Moreover, in order to ensure the range of λ was wide enough to include the best `softImpute` solution, we set the argument `rank.max` in that algorithm to be 20.

In Table 1, we report the statistical error of `primePCA` after 2000 iterations of refinement, together with the corresponding statistical errors of our initial estimator `primePCA_init` and those of `softImpute(oracle)` and `hardImpute`. Remarkably, `primePCA` exhibits stronger performance than these other methods across each of the signal-to-noise ratio regimes and different missingness mechanisms. We also remark that `hardImpute` is inaccurate and unstable, because it might converge to the local optimum that is far from the truth.

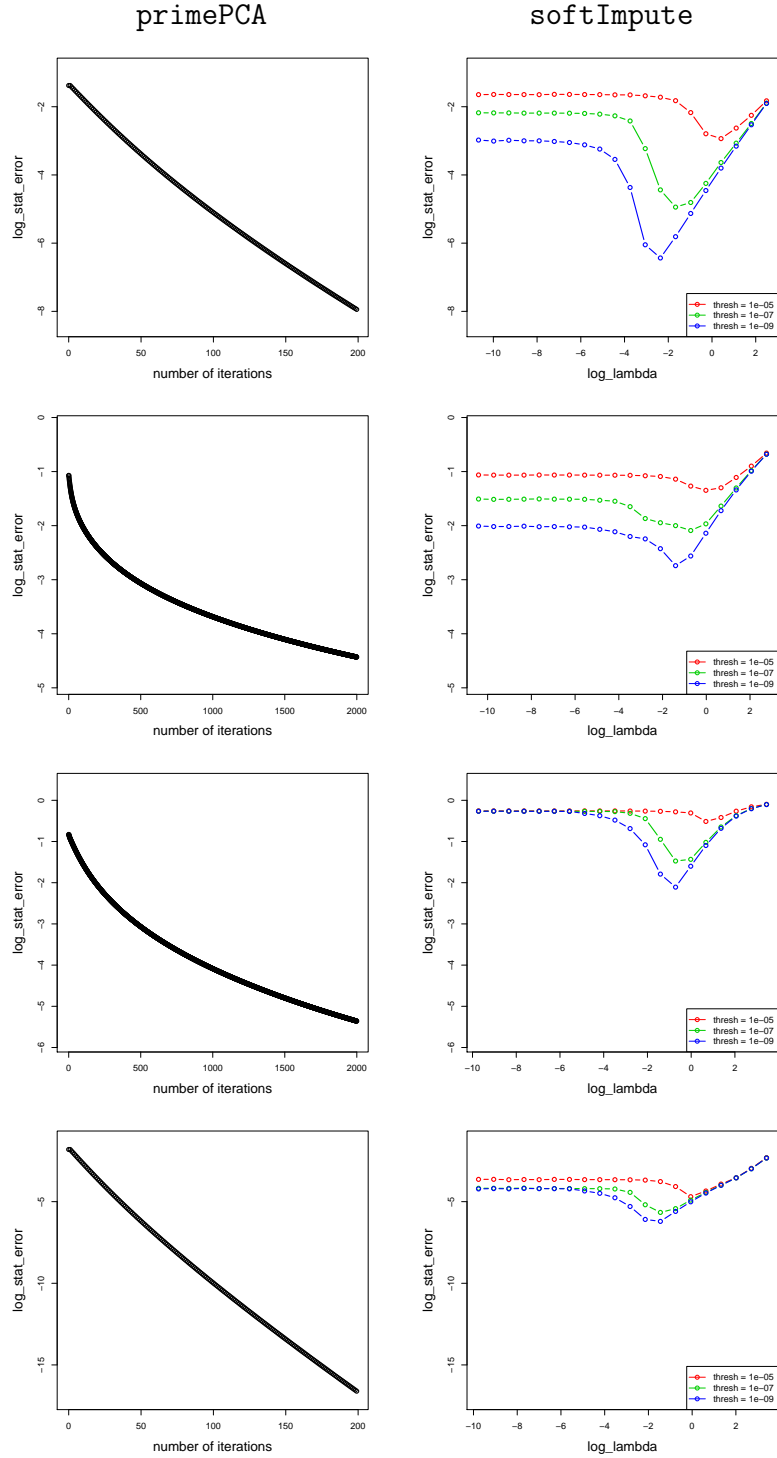


Figure 2: Logarithms of the average Frobenius norm $\sin \Theta$ error of **primePCA** and **softImpute** under various heterogeneity levels of missingness in absence of noise. The three rows of plots above, from the top to bottom, correspond to (H1), (H2), (H3) and (H4).

Table 1: Average losses (with standard errors in brackets) under (H1), (H2), (H3) and (H4).

		$\nu = 20$	$\nu = 40$	$\nu = 60$
(H1)	<code>hardImpute</code>	0.444 _(0.001)	0.251 _(0.001)	0.186 _(0.0005)
	<code>softImpute(oracle)</code>	0.186 _(0.0004)	0.095 _(0.0002)	0.064 _(0.0002)
	<code>primePCA_init</code>	0.306 _(0.001)	0.266 _(0.001)	0.259 _(0.001)
	<code>primePCA</code>	0.171 _(0.0004)	0.084 _(0.0002)	0.056 _(0.0001)
(H2)	<code>hardImpute</code>	0.473 _(0.001)	0.291 _(0.001)	0.236 _(0.001)
	<code>softImpute(oracle)</code>	0.308 _(0.001)	0.185 _(0.001)	0.141 _(0.001)
	<code>primePCA_init</code>	0.399 _(0.002)	0.357 _(0.001)	0.349 _(0.001)
	<code>primePCA</code>	0.232 _(0.001)	0.115 _(0.001)	0.077 _(0.0005)
(H3)	<code>hardImpute</code>	0.479 _(0.001)	0.385 _(0.001)	0.427 _(0.001)
	<code>softImpute(oracle)</code>	0.374 _(0.001)	0.222 _(0.001)	0.170 _(0.001)
	<code>primePCA_init</code>	0.486 _(0.001)	0.449 _(0.001)	0.442 _(0.001)
	<code>primePCA</code>	0.290 _(0.001)	0.145 _(0.001)	0.097 _(0.0004)
(H4)	<code>hardImpute</code>	0.174 _(0.0005)	0.089 _(0.0003)	0.062 _(0.0003)
	<code>softImpute(oracle)</code>	0.121 _(0.0002)	0.062 _(0.0001)	0.042 _(0.0001)
	<code>primePCA_init</code>	0.203 _(0.001)	0.175 _(0.0005)	0.169 _(0.0004)
	<code>primePCA</code>	0.116 _(0.0003)	0.058 _(0.0002)	0.038 _(0.0001)

4.3 Near low-rank case

Here, we set $n = 2000$, $d = 500$, $K = 10$, $\Sigma_{\mathbf{u}} = \text{diag}(2^{10}, 2^9, \dots, 2)$, and fixed \mathbf{V}_K once for all experiments to be the top K eigenvectors of one realisation¹ of the sample covariance matrix of n independent $N_d(\mathbf{0}, \mathbf{I}_d)$ random vectors. Here $\|\mathbf{V}_K\|_{\infty}/d^{1/2} < 3.63$, and we again generated the rows of \mathbf{Z} as independent $N_d(\mathbf{0}, \mathbf{I}_d)$ random vectors. Table 2 reports the average loss of estimating the top \hat{K} eigenvectors of $\Sigma_{\mathbf{y}}$, where \hat{K} varies from 1 to 5. Interestingly, even in this misspecified setting, `primePCA` is competitive with the oracle version of `softImpute`.

5 Real data analysis: Million Song Dataset

We apply `primePCA` to a subset of the Million Song Dataset² to analyse music preferences. The original data can be expressed as a matrix with 110,000 users (rows) and 163,206 songs (columns), with entries representing the number of times a song was played by a particular user. The proportion of non-missing entries in the matrix is 0.008%. Since the matrix is very sparse, and since most songs have very few listeners, we enhance the signal-to-noise ratio by restricting our attention to songs that have at least 100 listeners (1,777 songs in total). This improves the proportion of non-missing entries to 0.23%. Further summary information

¹In R, we set the random seed to be 2019 before generating \mathbf{V}_K .

²<https://www.kaggle.com/c/msdchallenge/data>

Table 2: Average losses (with standard errors in brackets) in the setting of Section 4.3 under (H1), (H2), (H3) and (H4).

		$\widehat{K} = 1$	$\widehat{K} = 2$	$\widehat{K} = 3$	$\widehat{K} = 4$	$\widehat{K} = 5$
(H1)	hardImpute	0.308 _(0.002)	0.507 _(0.002)	0.764 _(0.004)	1.199 _(0.006)	1.524 _(0.004)
	softImpute(oracle)	0.107 _(0.001)	0.182 _(0.001)	0.275 _(0.001)	0.401 _(0.001)	0.596 _(0.001)
	primePCA_init	0.203 _(0.001)	0.345 _(0.001)	0.554 _(0.003)	1.074 _(0.007)	1.427 _(0.006)
	primePCA	0.141 _(0.001)	0.200 _(0.001)	0.269 _(0.001)	0.374 _(0.001)	0.580 _(0.001)
(H2)	hardImpute	0.298 _(0.002)	0.466 _(0.002)	0.696 _(0.003)	1.124 _(0.006)	1.452 _(0.004)
	softImpute(oracle)	0.188 _(0.001)	0.283 _(0.001)	0.410 _(0.001)	0.562 _(0.001)	0.751 _(0.001)
	primePCA_init	0.285 _(0.001)	0.443 _(0.004)	0.757 _(0.013)	1.201 _(0.004)	1.533 _(0.003)
	primePCA	0.190 _(0.002)	0.267 _(0.002)	0.368 _(0.003)	0.543 _(0.008)	0.797 _(0.009)
(H3)	hardImpute	0.302 _(0.001)	0.482 _(0.002)	0.695 _(0.002)	1.004 _(0.006)	1.373 _(0.004)
	softImpute(oracle)	0.206 _(0.001)	0.338 _(0.001)	0.492 _(0.001)	0.664 _(0.002)	0.878 _(0.002)
	primePCA_init	0.341 _(0.001)	0.528 _(0.019)	1.097 _(0.008)	1.306 _(0.008)	1.597 _(0.004)
	primePCA	0.222 _(0.001)	0.330 _(0.002)	0.452 _(0.003)	0.641 _(0.008)	0.919 _(0.007)
(H4)	hardImpute	0.090 _(0.001)	0.148 _(0.001)	0.226 _(0.001)	0.346 _(0.002)	0.589 _(0.007)
	softImpute(oracle)	0.071 _(0.001)	0.112 _(0.001)	0.164 _(0.001)	0.233 _(0.001)	0.332 _(0.001)
	primePCA_init	0.139 _(0.001)	0.220 _(0.001)	0.325 _(0.001)	0.475 _(0.002)	0.805 _(0.012)
	primePCA	0.098 _(0.001)	0.135 _(0.001)	0.176 _(0.001)	0.236 _(0.001)	0.328 _(0.001)

about the filtered data is provided below:

1. Quantiles of non-missing matrix entry values:

0%	10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
1	1	1	1	1	1	2	3	5	8	500

90%	91%	92%	93%	94%	95%	96%	97%	98%	99%	100%
8	9	9	10	11	13	15	18	23	33	500

2. Quantiles of the number of listeners for each song:

0%	10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
100	108	117	126	139	154	178	214	272.8	455.6	5043

3. Quantiles of the total play counts of each user:

0%	10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
0	0	1	3	4	6	9	14	21	38	1114

90%	91%	92%	93%	94%	95%	96%	97%	98%	99%	100%
38	41	44	48	54	60	68	79	97	132	1114

Moreover, from the first numbered point above, we see that the distribution of play counts has an extremely heavy tail. To guard against excessive influence from the outliers, we discretise the play counts into five interest levels as follows:

Play count	1	2 – 3	4 – 6	7 – 10	≥ 11
Level of interest	1	2	3	4	5

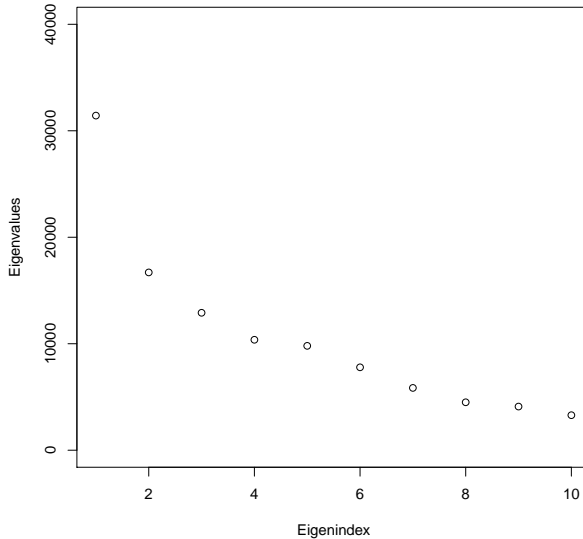


Figure 3: Leading eigenvalues of $\widehat{\Sigma}_{\mathbf{y}}$.

We are now in a position to analyse the data using `primePCA`. For $i = 1, \dots, n = 110,000$ and $j = 1, \dots, d = 1,777$, let $Y_{ij} \in \{1, \dots, 5\}$ denote the level of interest of user i in song j , let $\widehat{K} = 10$ and let $\mathcal{I} = \{i : \|\omega_i\|_1 > \widehat{K}\}$. Our initial goal is to assess the top \widehat{K} eigenvalues of $\Sigma_{\mathbf{y}}$ to see if there is low-rank signal in $\mathbf{Y} = (Y_{ij})$. To this end, we first apply Algorithm 2 to obtain $\mathbf{V}_{\widehat{K}}^{\text{prime}}$; next, for each $i \in \mathcal{I}$, we run Steps 2–5 of Algorithm 1 to obtain the estimated principal score $\widehat{\mathbf{u}}_i$, so that we can approximate \mathbf{y}_i by $\widehat{\mathbf{y}}_i = \widehat{\mathbf{V}}_{\widehat{K}}^{\text{prime}} \widehat{\mathbf{u}}_i$. This allows us to estimate $\Sigma_{\mathbf{y}}$ by $\widehat{\Sigma}_{\mathbf{y}} = n^{-1} \sum_{i \in \mathcal{I}} \widehat{\mathbf{y}}_i \widehat{\mathbf{y}}_i^\top$. Figure 3 displays the top \widehat{K} eigenvalues of $\widehat{\Sigma}_{\mathbf{y}}$, which exhibit a fairly rapid decay, thereby providing evidence for the existence of low-rank signal in \mathbf{Y} .

In the left panel of Figure 4, we present the estimate $\widehat{\mathbf{V}}_2^{\text{prime}}$ of the top two eigenvectors of the covariance matrix $\Sigma_{\mathbf{y}}$, with colours indicating the genre of the song. The outliers in the x -axis of this plot are particularly interesting: they reveal songs that polarise opinion among users (see Table 3) and that best capture variation in individuals’ preferences for types of music measured by the first principal component. It is notable that Rock songs are overrepresented among the outliers (see Table 4), relative to, say, Country songs. Users who express a preference for particular songs are also more likely to enjoy songs that are

nearby in the plot. Such information is therefore potentially commercially valuable, both as an efficient means of gauging users’ preferences, and for providing recommendations.

The right panel of Figure 4 presents the principal scores $\{\hat{\mathbf{u}}_i\}_{i=1}^n$ of the users, with frequent users (whose total song plays are in the top 10% of all users) in red and occasional users in blue. This plot reveals, for instance, that the second principal component is well aligned with general interest in the website. Returning to the left plot, we can now interpret a positive y -coordinate for a particular song (which is the case for the large majority of songs) as being associated with an overall interest in the music provided by the site.

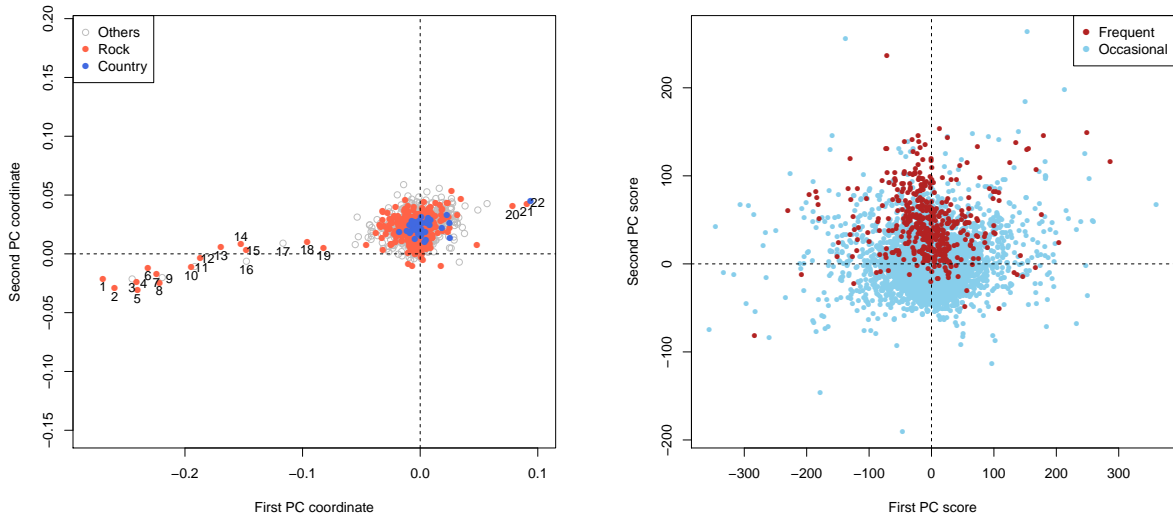


Figure 4: Plots of the first two principal components $\hat{\mathbf{V}}_2^{\text{prime}}$ (left) and the associated scores $\{\hat{\mathbf{u}}_i\}_{i=1}^n$ (right).

6 Discussion

Heterogeneous missingness is ubiquitous in contemporary, large-scale data sets, yet we currently understand very little about how existing procedures perform or should be adapted to cope with the challenges this presents. Here we attempt to extract the lessons learned from this study of high-dimensional PCA, in order to see how related ideas may be relevant in other statistical problems where one wishes to recover low-dimensional structure with data corrupted in a heterogeneous manner.

A key insight, as gleaned from Section 2.2, is that the way in which the heterogeneity interacts with the underlying structure of interest is crucial. In the worst case, the missingness may be constructed to conceal precisely the structure one seeks to uncover, thereby rendering the problem infeasible by any method. The only hope, then, in terms of providing

Table 3: Titles, artists and genres of the 22 outlier songs in Figure 4.

ID	Title	Artist	Genre
1	Your Hand In Mine	Explosions In The Sky	Rock
2	All These Things That I've Done	The Killers	Rock
3	Lady Marmalade	Christina Aguilera / Lil' Kim/ Mya / Pink	Pop
4	Here It Goes Again	Ok Go	Rock
5	I Hate Pretending (Album Version)	Secret Machines	Rock
6	No Rain	Blind Melon	Rock
7	Comatose (Comes Alive Version)	Skillet	Rock
8	Life In Technicolor	Coldplay	Rock
9	New Soul	Yael Nam	Pop
10	Blurry	Puddle Of Mudd	Rock
11	Give It Back	Polly Paulusma	Pop
12	Walking On The Moon	The Police	Rock
13	Face Down (Album Version)	The Red Jumpsuit Apparatus	Rock
14	Savior	Rise Against	Rock
15	Swing Swing	The All-American Rejects	Rock
16	Without Me	Eminem	Rap
17	Almaz	Randy Crawford	Pop
18	Hotel California	Eagles	Rock
19	Hey There Delilah	Plain White T's	Rock
20	Revelry	Kings Of Leon	Rock
21	Undo	Bjrk	Rock
22	You're The One	Dwight Yoakam	Country

theoretical guarantees, is to rule out such an adversarial interaction. This was achieved via our incoherence condition in Section 3, and we look forward to seeing how the relevant interactions between structure and heterogeneity can be controlled in other statistical problems such as those mentioned in the introduction. For instance, in sparse linear regression, one would anticipate that missingness of covariates with strong signal would be much more harmful than corresponding missingness for noise variables.

Our study also contributes to the broader understanding of the uses and limitations

Table 4: Genre distribution of the outliers (songs whose corresponding coordinate in the estimated leading principal component is of magnitude larger than 0.07).

	Rock	Pop	Electronic	Rap	Country	RnB	Latin	Others
Population	48.92%	18.53%	9.12%	7.15%	4.33%	2.35%	2.26%	7.34%
Outliers	72.73%	18.18%	0%	4.54%	4.54%	0%	0%	0%

of spectral methods for estimating hidden low-dimensional structures in high-dimensional problems. We have seen that the IPW estimator is both methodologically simple and achieves near-minimax optimality when the noise level is of constant order. Similar results have been obtained for spectral clustering for network community detection in stochastic block models (Rohe et al., 2011) and in low-rank-plus-sparse matrix estimation problems (Fan, Liao and Mincheva, 2013). On the other hand, the IPW estimator fails to provide exact recovery of the principal components in the noiseless setting. In these other aforementioned problems, it has also been observed that refinement of an initial spectral estimator can enhance performance, particularly in high signal-to-noise ratio regimes (Gao et al., 2016; Zhang, Cai and Wu, 2018), as we were able to show for our primePCA algorithm. This suggests that such a refinement has the potential to confer a sharper dependence of the statistical error rate on the signal-to-noise ratio compared with a vanilla spectral algorithm, and understanding this phenomenon in greater detail provides another interesting avenue for future research.

7 Proofs of main results

We define two linear maps $\mathcal{D}, \mathcal{F} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}$, such that for any $\mathbf{A} = (A_{ij}) \in \mathbb{R}^{d \times d}$, we have $[\mathcal{D}(\mathbf{A})]_{ij} := A_{ij} \mathbb{1}_{\{i=j\}}$ and $\mathcal{F}(\mathbf{A}) := \mathbf{A} - \mathcal{D}(\mathbf{A})$. In other words, $\mathcal{D}(\mathbf{A})$ and $\mathcal{F}(\mathbf{A})$ correspond to the diagonal and off-diagonal parts of \mathbf{A} respectively.

Proof of Theorem 1. Since $\mathbf{y}_i = \mathbf{V}_K \mathbf{u}_i + \mathbf{z}_i$, we have that

$$\|\mathbf{y}_i\|_{\psi_2} \leq \|\mathbf{V}_K \mathbf{u}_i\|_{\psi_2} + \|\mathbf{z}_i\|_{\psi_2} = \|\mathbf{u}_i\|_{\psi_2} + \|\mathbf{z}_i\|_{\psi_2} \leq (\lambda^{1/2} + 1)\tau. \quad (9)$$

Moreover, since $\max_{j \in [d]} \|y_{1j}\|_{\psi_2} \leq M^{1/2}$ by Lemma 1, it follows from van der Vaart and Wellner (1996, Lemma 2.2.2) that there exist a universal constant $C > 0$ such that³

$$\|\|\mathbf{y}_i\|_{\infty}\|_{\psi_2} \leq \{CM \log d\}^{1/2}. \quad (10)$$

Recall that $\tilde{\mathbf{y}}_i^\top = (\tilde{y}_{i1}, \dots, \tilde{y}_{id})$ denotes the i th row of \mathbf{Y}_Ω . Define $\mathbf{A}_i := \mathcal{F}(\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^\top)$ and

³In van der Vaart and Wellner (1996), the ψ_2 -norm of a random variable is defined slightly differently as $\|X\|_{\psi_2} := \inf\{a : \mathbb{E}e^{(X/a)^2} \leq 2\}$. It can be shown (Vershynin, 2012, Lemma 5.5) that these two norms are equivalent.

$\mathbf{B}_i := \mathcal{D}(\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^\top)$. We have the following decomposition of $\widehat{\mathbf{G}}$:

$$\begin{aligned}
\widehat{\mathbf{G}} &= \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{\widehat{p}^2} \mathbf{A}_i - \frac{1}{p^2} \mathbb{E} \mathbf{A}_i \right) + \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{\widehat{p}} \mathbf{B}_i - \frac{1}{p} \mathbb{E} \mathbf{B}_i \right) + \boldsymbol{\Sigma}_{\mathbf{y}} \\
&= \frac{1}{n \widehat{p}^2} \sum_{i=1}^n (\mathbf{A}_i - \mathbb{E} \mathbf{A}_i) + \frac{1}{n \widehat{p}} \sum_{i=1}^n (\mathbf{B}_i - \mathbb{E} \mathbf{B}_i) + \left(\frac{1}{\widehat{p}^2} - \frac{1}{p^2} \right) \mathbb{E} \mathbf{A}_1 + \left(\frac{1}{\widehat{p}} - \frac{1}{p} \right) \mathbb{E} \mathbf{B}_1 + \boldsymbol{\Sigma}_{\mathbf{y}} \\
&= \frac{1}{n \widehat{p}^2} \sum_{i=1}^n (\mathbf{A}_i - \mathbb{E} \mathbf{A}_i) + \frac{1}{n \widehat{p}} \sum_{i=1}^n (\mathbf{B}_i - \mathbb{E} \mathbf{B}_i) + \left(\frac{p^2}{\widehat{p}^2} - 1 \right) \mathcal{F}(\boldsymbol{\Sigma}_{\mathbf{y}}) + \left(\frac{p}{\widehat{p}} - 1 \right) \mathcal{D}(\boldsymbol{\Sigma}_{\mathbf{y}}) + \boldsymbol{\Sigma}_{\mathbf{y}} \\
&= \frac{1}{n \widehat{p}^2} \sum_{i=1}^n (\mathbf{A}_i - \mathbb{E} \mathbf{A}_i) + \frac{1}{n \widehat{p}} \sum_{i=1}^n (\mathbf{B}_i - \mathbb{E} \mathbf{B}_i) + \left(\frac{p}{\widehat{p}} - \frac{p^2}{\widehat{p}^2} \right) \mathcal{D}(\boldsymbol{\Sigma}_{\mathbf{y}}) + \frac{p^2}{\widehat{p}^2} \boldsymbol{\Sigma}_{\mathbf{y}}.
\end{aligned}$$

We regard $\widehat{\mathbf{G}}$ as a perturbed version of $(p^2/\widehat{p}^2)\boldsymbol{\Sigma}_{\mathbf{y}}$. Applying [Yu, Wang and Samworth \(2015, Theorem 2\)](#), we have

$$\begin{aligned}
L(\widehat{\mathbf{V}}_K, \mathbf{V}_K) &\leq \frac{2K^{1/2} \widehat{p}^2}{p^2 \lambda_K} \left\| \frac{1}{n \widehat{p}^2} \sum_{i=1}^n (\mathbf{A}_i - \mathbb{E} \mathbf{A}_i) + \frac{1}{n \widehat{p}} \sum_{i=1}^n (\mathbf{B}_i - \mathbb{E} \mathbf{B}_i) + \left(\frac{p}{\widehat{p}} - \frac{p^2}{\widehat{p}^2} \right) \mathcal{D}(\boldsymbol{\Sigma}_{\mathbf{y}}) \right\|_{\text{op}} \\
&\leq \frac{2K^{1/2}}{\lambda_K} \left(\left\| \frac{1}{n p^2} \sum_{i=1}^n (\mathbf{A}_i - \mathbb{E} \mathbf{A}_i) \right\|_{\text{op}} + \left\| \frac{\widehat{p}}{n p^2} \sum_{i=1}^n (\mathbf{B}_i - \mathbb{E} \mathbf{B}_i) \right\|_{\text{op}} + \left\| \left(\frac{\widehat{p}}{p} - 1 \right) \mathcal{D}(\boldsymbol{\Sigma}_{\mathbf{y}}) \right\|_{\text{op}} \right). \tag{11}
\end{aligned}$$

We will control the expectation of the three terms on the right-hand side of [\(11\)](#) separately. Define $\widehat{p}_i := d^{-1} \sum_{j=1}^d \omega_{ij}$. For notational simplicity, we write \mathbb{P}' and \mathbb{E}' respectively for the probability and expectation conditional on $(\widehat{p}_1, \dots, \widehat{p}_n)$. Also, let $\widehat{p}_i^{(2)} := \mathbb{E}'(\omega_{i1} \omega_{i2})$ and $\widehat{p}_i^{(3)} := \mathbb{E}'(\omega_{i1} \omega_{i2} \omega_{i3})$ (if $d = 2$, then $\widehat{p}_i^{(3)} := 0$). For the first term, we apply a symmetrisation argument. Let $\{\mathbf{A}_i^*\}_{i=1}^n$ denote copies of $\{\mathbf{A}_i\}_{i=1}^n$ that are independent of $\{\mathbf{u}_i, \mathbf{z}_i, \boldsymbol{\omega}_i\}_{i=1}^n$, let $\{\epsilon_i\}_{i=1}^n$ be independent Rademacher random variables that are independent of $\{\mathbf{u}_i, \mathbf{z}_i, \boldsymbol{\omega}_i, \mathbf{A}_i^*\}_{i=1}^n$ and write \mathbb{E}^* for expectation conditional on $\{\mathbf{u}_i, \mathbf{z}_i, \boldsymbol{\omega}_i\}_{i=1}^n$. Then by Jensen's inequality,

$$\begin{aligned}
\mathbb{E} \left\| \frac{1}{n p^2} \sum_{i=1}^n (\mathbf{A}_i - \mathbb{E} \mathbf{A}_i) \right\|_{\text{op}} &= \mathbb{E} \left\| \frac{1}{n p^2} \sum_{i=1}^n (\mathbf{A}_i - \mathbb{E}^* \mathbf{A}_i^*) \right\|_{\text{op}} \leq \mathbb{E} \left\| \frac{1}{n p^2} \sum_{i=1}^n (\mathbf{A}_i - \mathbf{A}_i^*) \right\|_{\text{op}} \\
&= \mathbb{E} \left\| \frac{1}{n p^2} \sum_{i=1}^n \epsilon_i (\mathbf{A}_i - \mathbf{A}_i^*) \right\|_{\text{op}} \leq 2 \mathbb{E} \left\| \frac{1}{n p^2} \sum_{i=1}^n \epsilon_i \mathbf{A}_i \right\|_{\text{op}}. \tag{12}
\end{aligned}$$

Since $\mathbf{A}_i = \tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^\top - \mathcal{D}(\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^\top)$, we have that

$$\mathbb{E}' \{ (\mathbf{A}_i^2)_{jk} \mid \mathbf{y}_i \} = \begin{cases} \mathbb{E}' \{ \widehat{y}_{ij}^2 \|\tilde{\mathbf{y}}_i\|_2^2 - \widehat{y}_{ij}^4 \mid \mathbf{y}_i \} = \widehat{p}_i^{(2)} y_{ij}^2 \sum_{t \neq j} y_{it}^2, & \text{if } j = k, \\ \sum_{t \notin \{j, k\}} \mathbb{E}' \{ \widehat{y}_{ij} \widehat{y}_{ik} \widehat{y}_{it}^2 \mid \mathbf{y}_i \} = \widehat{p}_i^{(3)} y_{ij} y_{ik} \sum_{t \notin \{j, k\}} y_{it}^2, & \text{if } j \neq k. \end{cases}$$

Writing $\mathbf{y}_{i,-t} := \mathbf{y}_i - y_{it}\mathbf{e}_t$, we then have

$$\begin{aligned}\mathbb{E}'(\mathbf{A}_i^2 \mid \mathbf{y}_i) &= \widehat{p}_i^{(3)} \sum_{t=1}^d y_{it}^2 \mathbf{y}_{i,-t} \mathbf{y}_{i,-t}^\top + (\widehat{p}_i^{(2)} - \widehat{p}_i^{(3)}) \mathcal{D} \left(\sum_{t=1}^d y_{it}^2 \mathbf{y}_{i,-t} \mathbf{y}_{i,-t}^\top \right) \\ &\preceq \widehat{p}_i^{(3)} \|\mathbf{y}_i\|_\infty^2 \sum_{t=1}^d \mathbf{y}_{i,-t} \mathbf{y}_{i,-t}^\top + (\widehat{p}_i^{(2)} - \widehat{p}_i^{(3)}) \|\mathbf{y}_i\|_\infty^2 \mathcal{D} \left(\sum_{t=1}^d \mathbf{y}_{i,-t} \mathbf{y}_{i,-t}^\top \right).\end{aligned}$$

Notice that

$$\sum_{t=1}^d \mathbf{y}_{i,-t} \mathbf{y}_{i,-t}^\top = \sum_{t=1}^d (\mathbf{y}_i \mathbf{y}_i^\top - y_{it} \mathbf{e}_t \mathbf{y}_i^\top - y_{it} \mathbf{y}_i \mathbf{e}_t^\top + y_{it}^2 \mathbf{e}_t \mathbf{e}_t^\top) = (d-2) \mathbf{y}_i \mathbf{y}_i^\top + \mathcal{D}(\mathbf{y}_i \mathbf{y}_i^\top).$$

Therefore,

$$\begin{aligned}\mathbb{E}'(\mathbf{A}_i^2 \mid \mathbf{y}_i) &\preceq \|\mathbf{y}_i\|_\infty^2 \{ \widehat{p}_i^{(3)} (d-2) \mathbf{y}_i \mathbf{y}_i^\top + ((d-1)\widehat{p}_i^{(2)} - (d-2)\widehat{p}_i^{(3)}) \mathcal{D}(\mathbf{y}_i \mathbf{y}_i^\top) \} \\ &\preceq d \|\mathbf{y}_i\|_\infty^2 \{ \widehat{p}_i^{(3)} \mathbf{y}_i \mathbf{y}_i^\top + \widehat{p}_i^{(2)} \mathcal{D}(\mathbf{y}_i \mathbf{y}_i^\top) \}.\end{aligned}$$

Now, observe that $\|\mathbf{A}_i\|_{\text{op}} \leq d\widehat{p}_i \|\mathbf{y}_i\|_\infty^2$, so for $q \geq 2$,

$$\mathbb{E}'(\mathbf{A}_i^q) \preceq \mathbb{E}' \{ (d\widehat{p}_i \|\mathbf{y}_i\|_\infty^2)^{q-2} \mathbb{E}'(\mathbf{A}_i^2 \mid \mathbf{y}_i) \} \preceq d^{q-1} \widehat{p}_i^{q-2} \mathbb{E}' [\|\mathbf{y}_i\|_\infty^{2q-2} \{ \widehat{p}_i^{(3)} \mathbf{y}_i \mathbf{y}_i^\top + \widehat{p}_i^{(2)} \mathcal{D}(\mathbf{y}_i \mathbf{y}_i^\top) \}].$$

By the Cauchy–Schwarz inequality, we therefore have that

$$\begin{aligned}\|\mathbb{E}'(\epsilon_i^q \mathbf{A}_i^q)\|_{\text{op}} &\leq d^{q-1} \widehat{p}_i^{q-2} \widehat{p}_i^{(3)} \left[\mathbb{E}(\|\mathbf{y}_i\|_\infty^{4q-4}) \sup_{\mathbf{v} \in \mathcal{S}^{d-1}} \mathbb{E}\{(\mathbf{v}^\top \mathbf{y}_i)^4\} \right]^{1/2} + d^{q-1} \widehat{p}_i^{q-2} \widehat{p}_i^{(2)} \mathbb{E} \|\mathbf{y}_i\|_\infty^{2q} \\ &\leq d^{q-1} \widehat{p}_i^{q-2} \left\{ \widehat{p}_i^{(3)} (4q-4)^{q-1} (CM \log d)^{q-1} 8R\tau^2 + \widehat{p}_i^{(2)} (2q)^q (CM \log d)^q \right\} \\ &\leq \frac{q!}{2} \left\{ 32eCMR\tau^2 \widehat{p}_i^{(3)} d \log d + e^2 \widehat{p}_i^{(2)} d (CM \log d)^2 \right\} (4eCM\widehat{p}_i d \log d)^{q-2} \\ &\leq \frac{q!}{2} C' M d \log d \left\{ R\tau^2 \widehat{p}_i^{(3)} + \widehat{p}_i^{(2)} M \log d \right\} (4eCM\widehat{p}_i d \log d)^{q-2},\end{aligned}$$

where $C' > 0$ is a universal constant, the second inequality uses (9) and (10) and the penultimate bound uses Stirling's inequality.

Let $\rho := 4eCMd(\max_i \widehat{p}_i) \log d$ and $\sigma^2 := C'Mn^{-1}d \log d \sum_{i=1}^n \{ R\tau^2 \widehat{p}_i^{(3)} + \widehat{p}_i^{(2)} M \log d \}$. Then by Tropp (2012, Theorem 6.2), we obtain that

$$\mathbb{P}' \left(\left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbf{A}_i \right\|_{\text{op}} \geq t \right) \leq 2d \exp \left(\frac{-nt^2/2}{\sigma^2 + \rho t} \right).$$

Consequently, for $t_0 := 2\sigma n^{-1/2} \log^{1/2} d + 4\rho n^{-1} \log d$, we have

$$\mathbb{E}' \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbf{A}_i \right\|_{\text{op}} \leq t_0 + \int_{t_0}^{\infty} 2d \{ e^{-nt^2/(4\sigma^2)} + e^{-nt/(4\rho)} \} dt \leq 4t_0.$$

Given (12), integrating the left-hand side of the above inequality over $(\widehat{p}_i)_{i=1}^n$ yields

$$\begin{aligned} \mathbb{E} \left\| \frac{1}{np^2} \sum_{i=1}^n (\mathbf{A}_i - \mathbb{E}\mathbf{A}_i) \right\|_{\text{op}} &\lesssim \frac{(\mathbb{E}\sigma^2)^{1/2} \log^{1/2} d}{n^{1/2} p^2} + \frac{\mathbb{E}\rho \log d}{np^2} \\ &\lesssim \sqrt{\frac{Md\{R\tau^2 p + M \log d\} \log^2 d}{np^2}} + \frac{Md \log^2 d \log n}{np}, \end{aligned} \quad (13)$$

where the first inequality uses Jensen's inequality and the second inequality uses Lemma 4.

For the second sum on the right-hand side of (11), we have by van der Vaart and Wellner (1996, Lemma 2.2.2) again that

$$\begin{aligned} \left\| \left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{B}_i - \mathbb{E}\mathbf{B}_i) \right\|_{\text{op}} \right\|_{\psi_1} &= \left\| \max_{j \in [d]} \left\| \frac{1}{n} \sum_{i=1}^n (\widetilde{y}_{ij}^2 - \mathbb{E}\widetilde{y}_{ij}^2) \right\| \right\|_{\psi_1} \\ &\lesssim \frac{\log d}{n} \left\| \sum_{i=1}^n (\widetilde{y}_{i1}^2 - \mathbb{E}\widetilde{y}_{i1}^2) \right\|_{\psi_1} \lesssim \frac{M \log d}{\sqrt{n}}, \end{aligned}$$

where the final inequality uses Lemma 2 and the fact that $\|\widetilde{y}_{i1}^2 - \mathbb{E}\widetilde{y}_{i1}^2\|_{\psi_1} \leq \|\widetilde{y}_{i1}^2\|_{\psi_1} + \mathbb{E}\widetilde{y}_{i1}^2 \leq 2M$. Now by the Cauchy–Schwarz inequality,

$$\begin{aligned} \mathbb{E} \left\| \frac{\widehat{p}}{np^2} \sum_{i=1}^n (\mathbf{B}_i - \mathbb{E}\mathbf{B}_i) \right\|_{\text{op}} &\leq \left\{ \mathbb{E} \left(\frac{\widehat{p}^2}{p^4} \right) \mathbb{E} \left(\left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{B}_i - \mathbb{E}\mathbf{B}_i) \right\|_{\text{op}}^2 \right) \right\}^{1/2} \\ &\lesssim \left\{ \left(\frac{1}{p^2} + \frac{1}{ndp^3} \right) \frac{M^2 \log^2 d}{n} \right\}^{1/2} \lesssim \frac{M \log d}{p\sqrt{n}}, \end{aligned} \quad (14)$$

which is dominated by the bound in (13).

Finally, for the third term on the right-hand side of (11), we have by the Cauchy–Schwarz inequality again that

$$\mathbb{E} \left\| \left(\frac{\widehat{p}}{p} - 1 \right) \mathcal{D}(\boldsymbol{\Sigma}_y) \right\|_{\text{op}} \lesssim \frac{M}{\sqrt{ndp}}, \quad (15)$$

which is also dominated by the bound in (13). Substituting (13), (14) and (15) into (11) establishes (5). If we regard M and τ as constants and if $n \geq d \log^2 d \log^2 n / (\lambda_1 p + \log d)$, then the second term in the bracket of the right-hand side of (5) is dominated by the first term, and claim (6) follows immediately. \square

Proof of Theorem 2. Without loss of generality, we may assume that $d \geq 50$ and that d is even, and write $d = 2h$ for some $h \in \mathbb{N}$. By the Gilbert–Varshamov lemma (see, e.g. Massart, 2007, Lemma 4.7), there exist $W \subseteq \{0, 1\}^h$ such that $\log |W| \geq h/16$ and $d_{\text{H}}(\mathbf{w}, \mathbf{w}') \geq h/4$ for any distinct pair of vectors $\mathbf{w}, \mathbf{w}' \in W$. Let $\gamma \in [0, \pi/2]$ be a real number to be specified later. To each $\mathbf{w} \in W$, we can associate a distribution $P_{\mathbf{w}} \in \mathcal{P}_{n,d}(\lambda_1, p)$ such that \mathbf{U} is

a random vector ($n \times 1$ random matrix) with independent $N(0, \lambda_1)$ entries, \mathbf{Z} is an $n \times d$ random matrix with independent $N(0, 1)$ entries, and

$$\mathbf{V}_1 = \mathbf{V}_{1, \mathbf{w}} := \frac{1}{\sqrt{h}} \left\{ \mathbf{w} \otimes \begin{pmatrix} \cos \gamma \\ \sin \gamma \end{pmatrix} + (\mathbf{1}_h - \mathbf{w}) \otimes \begin{pmatrix} \cos \gamma \\ -\sin \gamma \end{pmatrix} \right\} \in \mathcal{S}^{d-1}.$$

Fixing distinct $\mathbf{w}, \mathbf{w}' \in W$, we write $\mathbf{v} = (v_j)_{j \in [d]} := \mathbf{V}_{1, \mathbf{w}}$ and $\mathbf{v}' = (v'_j)_{j \in [d]} := \mathbf{V}_{1, \mathbf{w}'}$ and let $Q_{\mathbf{w}}$ and $Q_{\mathbf{w}'}$ denote respectively the marginal distribution of $(\tilde{\mathbf{y}}_1, \boldsymbol{\omega}_1)$ under $P_{\mathbf{w}}$ and $P_{\mathbf{w}'}$. Define $S := \{j \in [d] : \omega_{1j} = 1\}$ and also set $\bar{\mathbf{v}}_S := (v_j \mathbf{1}_{\{j \in S\}})_{j \in [d]} \in \mathbb{R}^d$ and $\bar{\mathbf{v}}'_S := (v'_j \mathbf{1}_{\{j \in S\}})_{j \in [d]} \in \mathbb{R}^d$. Then

$$\begin{aligned} \text{KL}(P_{\mathbf{w}}, P_{\mathbf{w}'}) &= \text{KL}(Q_{\mathbf{w}}^{\otimes n}, Q_{\mathbf{w}'}^{\otimes n}) = n \text{KL}(Q_{\mathbf{w}}, Q_{\mathbf{w}'}) = n \mathbb{E}_{Q_{\mathbf{w}}} \left\{ \mathbb{E}_{Q_{\mathbf{w}}} \left(\log \frac{dQ_{\mathbf{w}}}{dQ_{\mathbf{w}'}} \mid \boldsymbol{\omega}_1 \right) \right\} \\ &= n \mathbb{E} \text{KL}(N_d(\mathbf{0}, \mathbf{I}_d + \lambda_1 \bar{\mathbf{v}}_S \bar{\mathbf{v}}_S^\top), N_d(\mathbf{0}, \mathbf{I}_d + \lambda_1 \bar{\mathbf{v}}'_S \bar{\mathbf{v}}'^\top)), \end{aligned} \quad (16)$$

where the final expectation is over the marginal distribution of S under $P_{\mathbf{w}}$. We partition $S = S_0 \sqcup S_{1+} \sqcup S_{1-}$, where $S_0 := \{j \in S : j \text{ is odd}\}$, $S_{1+} := \{j \in S : j \text{ is even and } v_j = v'_j\}$ and $S_{1-} := \{j \in S : j \text{ is even and } v_j \neq v'_j\}$. Since by construction we always have $\|\bar{\mathbf{v}}_S\|_2 = \|\bar{\mathbf{v}}'_S\|_2$, we can apply Lemma 5 to obtain

$$\begin{aligned} \text{KL}(N(\mathbf{0}, \mathbf{I}_d + \lambda_1 \bar{\mathbf{v}}_S \bar{\mathbf{v}}_S^\top), N(\mathbf{0}, \mathbf{I}_d + \lambda_1 \bar{\mathbf{v}}'_S (\bar{\mathbf{v}}'_S)^\top)) &= \frac{\lambda_1^2 (\|\bar{\mathbf{v}}_S\|_2^4 - \langle \bar{\mathbf{v}}_S, \bar{\mathbf{v}}'_S \rangle^2)}{2(1 + \lambda_1 \|\bar{\mathbf{v}}_S\|_2^2)} \\ &\leq \frac{\lambda_1^2 \langle \bar{\mathbf{v}}_S, \bar{\mathbf{v}}_S + \bar{\mathbf{v}}'_S \rangle \langle \bar{\mathbf{v}}_S, \bar{\mathbf{v}}_S - \bar{\mathbf{v}}'_S \rangle}{2 \max\{1, \lambda_1 \|\bar{\mathbf{v}}_S\|_2^2\}} = \frac{\lambda_1^2 (\sum_{j \in S_0 \cup S_{1+}} 2v_j^2) (\sum_{j \in S_{1-}} 2v_j^2)}{2 \max\{1, \lambda_1 \sum_{j \in S} v_j^2\}} \\ &\leq \min \left\{ \frac{2\lambda_1^2}{h^2} (|S_0 \times S_{1-}| \sin^2 \gamma \cos^2 \gamma + |S_{1+} \times S_{1-}| \sin^4 \gamma), \frac{2\lambda_1 |S_{1-}| \sin^2 \gamma}{h} \right\}. \end{aligned}$$

Substituting the above bound into (16), we have

$$\text{KL}(P_{\mathbf{w}}, P_{\mathbf{w}'}) \leq 2n\lambda_1 p \min\{1, \lambda_1 p\} \sin^2 \gamma. \quad (17)$$

On the other hand, since $d_{\text{H}}(\mathbf{w}, \mathbf{w}') \geq h/4$, we also have

$$\sin^2 \Theta(\mathbf{v}, \mathbf{v}') = 1 - (\mathbf{v}^\top \mathbf{v}')^2 = 1 - \left(1 - \frac{2d_{\text{H}}(\mathbf{w}, \mathbf{w}') \sin^2 \gamma}{h} \right)^2 \geq \frac{1}{2} \sin^2 \gamma. \quad (18)$$

By (17), (18) and Fano's inequality (Yu, 1997, Lemma 3),

$$\begin{aligned} \inf_{\hat{\mathbf{v}}} \sup_{P \in \mathcal{P}_{n, d, 1}(\lambda_1, p)} \mathbb{E}_P L(\hat{\mathbf{v}}, \mathbf{v}) &\geq \inf_{\hat{\mathbf{v}}} \max_{\mathbf{w} \in W} \mathbb{E}_{P_{\mathbf{w}}} L(\hat{\mathbf{v}}, \mathbf{v}) \\ &\geq \frac{1}{2\sqrt{2}} \sin \gamma \left(1 - \frac{\log 2 + 2n\lambda_1 p \min\{1, \lambda_1 p\} \sin^2 \gamma}{\log |W|} \right). \end{aligned}$$

We now choose $\gamma \in [0, \pi/2]$ such that $\sin^2 \gamma = \min\left\{\frac{\log |W|}{8n\lambda_1 p \min\{1, \lambda_1 p\}}, 1\right\}$. Since $d \geq 50$, we obtain $\log |W| \geq d/32 \geq 2 \log 2$. Therefore,

$$\inf_{\widehat{\mathbf{v}}} \sup_{P \in \mathcal{P}_{n,d,1}(\lambda_1, p)} \mathbb{E}_P L(\widehat{\mathbf{v}}, \mathbf{v}) \geq \frac{1}{8\sqrt{2}} \sin \gamma \geq \min\left\{\frac{1}{200\lambda_1} \sqrt{\frac{d \max(1, \lambda_1 p)}{np^2}}, \frac{1}{8\sqrt{2}}\right\},$$

as desired. \square

Proof of Proposition 1. For notational simplicity, we write $\widehat{\mathbf{V}}_K := \widehat{\mathbf{V}}_K^{(\text{in})}$ and $\widehat{\mathbf{V}}_{S,K} := (\widehat{\mathbf{V}}_K)_S$ for any $S \subseteq [d]$. For $i \in \mathcal{I}$, let $\ell_i^\top \in \mathbb{R}^K$ denote the i th row of \mathbf{L} . For any $i \in \mathcal{I}$, we have $\widehat{\mathbf{y}}_{i, \mathcal{J}_i} = \mathbf{y}_{i, \mathcal{J}_i}$ and

$$\begin{aligned} \widehat{\mathbf{y}}_{i, \mathcal{J}_i^c} - \mathbf{y}_{i, \mathcal{J}_i^c} &= \widehat{\mathbf{V}}_{\mathcal{J}_i^c, K} (\widehat{\mathbf{V}}_{\mathcal{J}_i, K}^\top \widehat{\mathbf{V}}_{\mathcal{J}_i, K})^{-1} \widehat{\mathbf{V}}_{\mathcal{J}_i, K}^\top \mathbf{y}_{i, \mathcal{J}_i} - \mathbf{y}_{i, \mathcal{J}_i^c} \\ &= \widehat{\mathbf{V}}_{\mathcal{J}_i^c, K} (\widehat{\mathbf{V}}_{\mathcal{J}_i, K}^\top \widehat{\mathbf{V}}_{\mathcal{J}_i, K})^{-1} \widehat{\mathbf{V}}_{\mathcal{J}_i, K}^\top \mathbf{R}_{\mathcal{J}_i} \mathbf{W}_{\widehat{\mathbf{V}}_K, \mathbf{R}} \mathbf{W}_{\widehat{\mathbf{V}}_K, \mathbf{R}}^{-1} \boldsymbol{\Gamma} \ell_i - \mathbf{R}_{\mathcal{J}_i^c} \boldsymbol{\Gamma} \ell_i \\ &= \widehat{\mathbf{V}}_{\mathcal{J}_i^c, K} (\widehat{\mathbf{V}}_{\mathcal{J}_i, K}^\top \widehat{\mathbf{V}}_{\mathcal{J}_i, K})^{-1} \widehat{\mathbf{V}}_{\mathcal{J}_i, K}^\top (\mathbf{R}_{\mathcal{J}_i} \mathbf{W}_{\widehat{\mathbf{V}}_K, \mathbf{R}} - \widehat{\mathbf{V}}_{\mathcal{J}_i, K}) \mathbf{W}_{\widehat{\mathbf{V}}_K, \mathbf{R}}^{-1} \boldsymbol{\Gamma} \ell_i \\ &\quad + (\widehat{\mathbf{V}}_{\mathcal{J}_i^c, K} - \mathbf{R}_{\mathcal{J}_i^c} \mathbf{W}_{\widehat{\mathbf{V}}_K, \mathbf{R}}) \mathbf{W}_{\widehat{\mathbf{V}}_K, \mathbf{R}}^{-1} \boldsymbol{\Gamma} \ell_i. \end{aligned}$$

Thus

$$\begin{aligned} \|\widehat{\mathbf{y}}_{i, \mathcal{J}_i^c} - \mathbf{y}_{i, \mathcal{J}_i^c}\|_\infty &\leq \sigma_* \sqrt{d} \|\widehat{\mathbf{V}}_{\mathcal{J}_i^c, K}\|_{2 \rightarrow \infty} \|\mathbf{R}_{\mathcal{J}_i} \mathbf{W}_{\widehat{\mathbf{V}}_K, \mathbf{R}} - \widehat{\mathbf{V}}_{\mathcal{J}_i, K}\|_{2 \rightarrow \infty} \|\boldsymbol{\Gamma} \ell_i\|_2 \\ &\quad + \|\widehat{\mathbf{V}}_{\mathcal{J}_i^c, K} - \mathbf{R}_{\mathcal{J}_i^c} \mathbf{W}_{\widehat{\mathbf{V}}_K, \mathbf{R}}\|_{2 \rightarrow \infty} \|\boldsymbol{\Gamma} \ell_i\|_2 \\ &\leq \Delta \|\boldsymbol{\Gamma} \ell_i\|_2 (1 + \sigma_* \sqrt{d} \|\widehat{\mathbf{V}}_K\|_{2 \rightarrow \infty}) \\ &\leq \Delta \sigma_1(\boldsymbol{\Gamma}) \mu_1 \left(\frac{K}{n}\right)^{1/2} \{1 + \sigma_* (\mu_2 \sqrt{K} + \Delta \sqrt{d})\} \\ &\leq \frac{C'}{n^{1/2}} \Delta \sigma_1(\boldsymbol{\Gamma}) \mu_1 \mu_2 K =: m, \end{aligned}$$

say, where $C' > 0$ depends only on σ_* and c_1 . Note that the inequality above holds for all $i \in \mathcal{I}$. Writing $\mathbf{E} := \widehat{\mathbf{Y}} - \mathbf{Y}$ for convenience, we have found that $\|\mathbf{E}\|_\infty \leq m$. Let $\mathbf{L}_\perp \in \mathbb{O}^{n \times (n-K)}$, $\mathbf{R}_\perp \in \mathbb{O}^{d \times (d-K)}$ be the orthogonal complements of $\mathbf{L} \in \mathbb{O}^{n \times K}$ and $\mathbf{R} \in \mathbb{O}^{d \times K}$ respectively, so that $(\mathbf{L}, \mathbf{L}_\perp) \in \mathbb{O}^{n \times n}$ and $(\mathbf{R}, \mathbf{R}_\perp) \in \mathbb{O}^{d \times d}$. We wish to apply [Cai and Zhang \(2018a, Theorem 1\)](#). To this end, note that

$$\|\mathbf{L}^\top \mathbf{E} \mathbf{R}\|_{\text{op}} = \sup_{\mathbf{s}, \mathbf{t} \in \mathcal{S}^{K-1}} (\mathbf{L} \mathbf{s})^\top \mathbf{E} (\mathbf{R} \mathbf{t}) \leq \|\mathbf{L}\|_{2 \rightarrow \infty} \|\mathbf{R}\|_{2 \rightarrow \infty} \|\mathbf{E}\|_1 \leq \frac{K \mu_1 \mu_2 m \|\boldsymbol{\Omega}^c\|_1}{\sqrt{nd}}.$$

Hence, writing $\alpha := \sigma_K(\boldsymbol{\Gamma} + \mathbf{L}^\top \mathbf{E} \mathbf{R})$, we have by Weyl's inequality that

$$\sigma_K(\boldsymbol{\Gamma}) - \frac{K \mu_1 \mu_2 m \|\boldsymbol{\Omega}^c\|_1}{\sqrt{nd}} \leq \alpha \leq \sigma_K(\boldsymbol{\Gamma}) + \frac{K \mu_1 \mu_2 m \|\boldsymbol{\Omega}^c\|_1}{\sqrt{nd}}.$$

Now, writing $\beta := \|\mathbf{L}_\perp^\top \widehat{\mathbf{Y}} \mathbf{R}_\perp\|_{\text{op}} = \|\mathbf{L}_\perp^\top \mathbf{E} \mathbf{R}_\perp\|_{\text{op}}$, we have

$$\beta \leq \|\mathbf{E}\|_{\text{op}} \leq \|\mathbf{E}\|_{\text{F}} \leq m \sqrt{\|\boldsymbol{\Omega}^c\|_1}.$$

In addition, by Cauchy–Schwarz and Jensen’s inequality,

$$\begin{aligned} \|\mathbf{L}^\top \mathbf{E}\|_{\text{op}} &= \sup_{\substack{\mathbf{s} \in \mathcal{S}^{K-1} \\ \mathbf{t} \in \mathcal{S}^{d-1}}} (\mathbf{L}\mathbf{s})^\top \mathbf{E}\mathbf{t} \leq \|\mathbf{L}\|_{2 \rightarrow \infty} \sup_{\mathbf{t} \in \mathcal{S}^{K-1}} \|\mathbf{E}\mathbf{t}\|_1 \\ &\leq \mu_1 (Kn)^{1/2} \frac{1}{n} \sum_{i=1}^n m \sqrt{\|\boldsymbol{\omega}_i^c\|_1} \leq \mu_1 m (K \|\boldsymbol{\Omega}^c\|_1)^{1/2}. \end{aligned}$$

Similarly,

$$\|\mathbf{E} \mathbf{R}\|_{\text{op}} \leq \mu_2 m (K \|\boldsymbol{\Omega}^c\|_1)^{1/2}.$$

Hence there exists $c_1 > 0$, depending only on μ_1, μ_2 and σ_* , such that whenever $\Delta \leq \frac{c_1 \sigma_K(\boldsymbol{\Gamma})}{K^2 \sigma_1(\boldsymbol{\Gamma}) \sqrt{d}}$, we have

$$\alpha^2 - \beta^2 - \min(\|\mathbf{L}^\top \mathbf{E}\|_{\text{op}}^2, \|\mathbf{E} \mathbf{R}\|_{\text{op}}^2) \geq \frac{\sigma_K^2(\boldsymbol{\Gamma})}{2} \quad \text{and} \quad \alpha, \beta \leq 2\sigma_K(\boldsymbol{\Gamma}).$$

Now let $\widehat{\mathbf{Y}}_{\mathcal{I}} = \widehat{\mathbf{L}} \widehat{\boldsymbol{\Gamma}} \widehat{\mathbf{R}}^\top$ be an SVD of $\widehat{\mathbf{Y}}$. We can now apply [Cai and Zhang \(2018a, Theorem 1\)](#) to deduce that for such c_2 ,

$$\begin{aligned} \|\sin \Theta(\widehat{\mathbf{R}}, \mathbf{R})\|_{\text{op}} &\leq \frac{\alpha \|\mathbf{L}^\top \mathbf{E}\|_{\text{op}} + \beta \|\mathbf{E} \mathbf{R}\|_{\text{op}}}{\alpha^2 - \beta^2 - \min(\|\mathbf{L}^\top \mathbf{E}\|_{\text{op}}^2, \|\mathbf{E} \mathbf{R}\|_{\text{op}}^2)} \leq \frac{4m(\mu_1 + \mu_2)(K \|\boldsymbol{\Omega}^c\|_1)^{1/2}}{\sigma_K(\boldsymbol{\Gamma})} \\ &\leq \frac{4C' K^{3/2} \sigma_1(\boldsymbol{\Gamma})(\mu_1 + \mu_2) \mu_1 \mu_2}{\sigma_K(\boldsymbol{\Gamma})} \left(\frac{\|\boldsymbol{\Omega}^c\|_1}{n}\right)^{1/2} \Delta =: \kappa \Delta, \end{aligned}$$

say. Similarly,

$$\|\sin \Theta(\widehat{\mathbf{L}}, \mathbf{L})\|_{\text{op}} \leq \frac{\alpha \|\mathbf{E} \mathbf{R}\|_{\text{op}} + \beta \|\mathbf{L}^\top \mathbf{E}\|_{\text{op}}}{\alpha^2 - \beta^2 - \min(\|\mathbf{L}^\top \mathbf{E}\|_{\text{op}}^2, \|\mathbf{E} \mathbf{R}\|_{\text{op}}^2)} \leq \kappa \Delta.$$

We are now in a position to show contraction in terms of two-to-infinity norm. By [Cape, Tang and Priebe \(2018, Theorem 3.7\)](#),

$$\begin{aligned} \mathcal{T}(\widehat{\mathbf{R}}, \mathbf{R}) &\leq \frac{2\|\mathbf{R}_\perp \mathbf{R}_\perp^\top \mathbf{E}^\top \mathbf{L} \mathbf{L}^\top\|_{2 \rightarrow \infty}}{\sigma_K(\boldsymbol{\Gamma})} + \frac{2\|\mathbf{R}_\perp \mathbf{R}_\perp^\top \mathbf{E}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top\|_{2 \rightarrow \infty}}{\sigma_K(\boldsymbol{\Gamma})} \|\sin \Theta(\widehat{\mathbf{L}}, \mathbf{L})\|_{\text{op}} \\ &\quad + \|\sin \Theta(\widehat{\mathbf{R}}, \mathbf{R})\|_{\text{op}}^2 \|\mathbf{R}\|_{2 \rightarrow \infty} =: T_1 + T_2 + T_3, \end{aligned} \quad (19)$$

say. Note that

$$\begin{aligned} \|\mathbf{R}_\perp \mathbf{R}_\perp^\top\|_{\infty \rightarrow \infty} &\leq \|\mathbf{I}_d\|_{\infty \rightarrow \infty} + \|\mathbf{R} \mathbf{R}^\top\|_{\infty \rightarrow \infty} = 1 + \sup_{\|\mathbf{v}\|_\infty \leq 1} \|\mathbf{R} \mathbf{R}^\top \mathbf{v}\|_\infty \\ &\leq 1 + \sup_{\|\mathbf{v}\|_2 \leq \sqrt{d}} \|\mathbf{R}\|_{2 \rightarrow \infty} \|\mathbf{R}^\top \mathbf{v}\|_2 \leq 1 + \sqrt{K} \mu_2. \end{aligned}$$

Hence,

$$\begin{aligned} T_1 &\leq \frac{2(1 + \sqrt{K}\mu_2)\|\mathbf{E}^\top \mathbf{L}\mathbf{L}^\top\|_{2 \rightarrow \infty}}{\sigma_K(\mathbf{\Gamma})} \leq \frac{2(1 + \sqrt{K}\mu_2)\|\mathbf{E}^\top \mathbf{L}\|_{2 \rightarrow \infty}}{\sigma_K(\mathbf{\Gamma})} \\ &\leq \frac{2(1 + \sqrt{K}\mu_2)\mu_1\sqrt{K}m\|\mathbf{\Omega}^c\|_{1 \rightarrow 1}}{\sqrt{n}\sigma_K(\mathbf{\Gamma})} \lesssim_{\mu_1, \mu_2} \frac{K^2\sigma_1(\mathbf{\Gamma})\|\mathbf{\Omega}^c\|_{1 \rightarrow 1}\Delta}{n\sigma_K(\mathbf{\Gamma})}. \end{aligned}$$

Moreover,

$$\begin{aligned} T_2 &\leq \frac{2(1 + \sqrt{K}\mu_2)\|\mathbf{E}^\top\|_{2 \rightarrow \infty}\kappa\Delta}{\sigma_K(\mathbf{\Gamma})} \leq \frac{2(1 + \sqrt{K}\mu_2)m\|\mathbf{\Omega}^c\|_{1 \rightarrow 1}^{1/2}\kappa\Delta}{\sigma_K(\mathbf{\Gamma})} \\ &\lesssim_{\mu_1, \mu_2} \frac{K^{3/2}\sigma_1(\mathbf{\Gamma})\|\mathbf{\Omega}^c\|_{1 \rightarrow 1}^{1/2}\kappa\Delta^2}{\sqrt{n}\sigma_K(\mathbf{\Gamma})}. \end{aligned}$$

Finally,

$$T_3 \leq \mu_2\kappa^2\Delta^2\left(\frac{K}{d}\right)^{1/2}.$$

Write

$$\eta := \frac{K^2\sigma_1(\mathbf{\Gamma})\|\mathbf{\Omega}^c\|_{1 \rightarrow 1}^{1/2}}{\sqrt{n}\sigma_K(\mathbf{\Gamma})}$$

for simplicity, so that $\kappa \lesssim_{\mu_1, \mu_2, \sigma_*} (d/K)^{1/2}\eta$. Given that $\mathcal{T}(\widehat{\mathbf{V}}_K^{(\text{out})}, \mathbf{V}_K) = \mathcal{T}(\widehat{\mathbf{R}}, \mathbf{R})$, substituting the bounds for T_1, T_2, T_3 into (19) yields that

$$\begin{aligned} \mathcal{T}(\widehat{\mathbf{V}}_K^{(\text{out})}, \mathbf{V}_K) &\lesssim_{\mu_1, \mu_2, \sigma_*} \left\{ \eta \left(\frac{\|\mathbf{\Omega}^c\|_{1 \rightarrow 1}}{n} \right)^{1/2} + \frac{\sqrt{d}}{K}\eta^2\Delta + \left(\frac{d}{K} \right)^{1/2} \eta^2\Delta \right\} \Delta \\ &\leq \eta^2 \left\{ \frac{\sigma_K(\mathbf{\Gamma})}{K^2\sigma_1(\mathbf{\Gamma})} + 2 \left(\frac{d}{K} \right)^{1/2} \Delta \right\} \Delta \lesssim_{\mu_1, \mu_2} \frac{\eta^2\sigma_K(\mathbf{\Gamma})}{K^2\sigma_1(\mathbf{\Gamma})} \Delta = \frac{K^2\sigma_1(\mathbf{\Gamma})\|\mathbf{\Omega}^c\|_{1 \rightarrow 1}\Delta}{\sigma_K(\mathbf{\Gamma})n}, \end{aligned}$$

as desired. \square

Proof of Theorem 3. We prove this result by induction on t . The case $t = 0$ is true by definition of Δ , so suppose that the conclusion holds for some $t \in \{0\} \cup [n_{\text{iter}} - 1]$. We make the following two claims:

(a) $\mathcal{I}^{(t)} = \mathcal{I}$;

(b) The error is further contracted by refinement, i.e., $\mathcal{T}(\widehat{\mathbf{V}}_K^{(t+1)}, \mathbf{V}_K) \leq \rho\mathcal{T}(\widehat{\mathbf{V}}_K^{(t)}, \mathbf{V}_K)$.

To prove claim (a), notice that for each $i \in [n]$, by Weyl's inequality and the inductive hypothesis,

$$\begin{aligned} |\sigma_K((\widehat{\mathbf{V}}_K^{(t)})_{\mathcal{J}_i}) - \sigma_K((\mathbf{V}_K)_{\mathcal{J}_i})| &= |\sigma_K((\widehat{\mathbf{V}}_K^{(t)})_{\mathcal{J}_i}) - \sigma_K((\mathbf{V}_K)_{\mathcal{J}_i} \mathbf{W}_{\widehat{\mathbf{V}}_K^{(t)}, \mathbf{V}_K})| \\ &\leq \|(\widehat{\mathbf{V}}_K^{(t)})_{\mathcal{J}_i} - (\mathbf{V}_K)_{\mathcal{J}_i} \mathbf{W}_{\widehat{\mathbf{V}}_K^{(t)}, \mathbf{V}_K}\|_{\text{op}} \\ &\leq |\mathcal{J}_i|^{1/2} \mathcal{T}(\widehat{\mathbf{V}}_K^{(t)}, \mathbf{V}_K) \leq |\mathcal{J}_i|^{1/2} \rho^t \Delta. \end{aligned}$$

Now, for $i \in \mathcal{I}$,

$$\begin{aligned}\sigma_K((\widehat{\mathbf{V}}_K^{(t)})_{\mathcal{J}_i}) &\geq \sigma_K((\mathbf{V}_K)_{\mathcal{J}_i}) - |\sigma_K((\widehat{\mathbf{V}}_K^{(t)})_{\mathcal{J}_i}) - \sigma_K((\mathbf{V}_K)_{\mathcal{J}_i})| \\ &\geq (\sigma_*^{-1} + \epsilon - \sqrt{d}\Delta)(|\mathcal{J}_i|/d)^{1/2}.\end{aligned}$$

On the other hand, if $i \in \mathcal{I}^c$ and $\|\boldsymbol{\omega}_i\|_1 > K$, then

$$\begin{aligned}\sigma_K((\widehat{\mathbf{V}}_K^{(t)})_{\mathcal{J}_i}) &\leq \sigma_K((\mathbf{V}_K)_{\mathcal{J}_i}) + |\sigma_K((\widehat{\mathbf{V}}_K^{(t)})_{\mathcal{J}_i}) - \sigma_K((\mathbf{V}_K)_{\mathcal{J}_i})| \\ &\leq (\sigma_*^{-1} - \epsilon + \sqrt{d}\Delta)(|\mathcal{J}_i|/d)^{1/2}.\end{aligned}$$

Hence, if we choose $c_1 \leq \epsilon$, then $\sqrt{d}\Delta < \epsilon$, so for $i \in \mathcal{I}$,

$$\sigma_K((\widehat{\mathbf{V}}_K^{(t)})_{\mathcal{J}_i}) > \left(\frac{|\mathcal{J}_i|}{d\sigma_*}\right)^{1/2};$$

moreover, for $i \in \mathcal{I}^c$,

$$\sigma_K((\widehat{\mathbf{V}}_K^{(t)})_{\mathcal{J}_i}) < \left(\frac{|\mathcal{J}_i|}{d\sigma_*}\right)^{1/2}.$$

Claim (a) follows. As for claim (b), note that $\widehat{\mathbf{V}}_K^{(t+1)} = \text{refine}(K, \widehat{\mathbf{V}}_K^{(t)}, \boldsymbol{\Omega}_{\mathcal{I}^{(t)}}, (\mathbf{Y}_\Omega)_{\mathcal{I}^{(t)}})$. Taking $c_1, C > 0$ from Proposition 1, and reducing c_1 if necessary so that $c_1 \leq \epsilon$, we may apply this proposition to deduce that whenever

- (i) $\mathcal{T}(\widehat{\mathbf{V}}_K^{(t)}, \mathbf{V}_K) \leq \frac{c_1 \sigma_K(\mathbf{Y}_{\mathcal{I}})}{K^2 \sigma_1(\mathbf{Y}_{\mathcal{I}}) \sqrt{d}}$;
- (ii) $\rho := \frac{CK^2 \sigma_1(\mathbf{Y}_{\mathcal{I}}) \|\boldsymbol{\Omega}_{\mathcal{I}}^\sharp\|_{1 \rightarrow 1}}{\sigma_K(\mathbf{Y}_{\mathcal{I}}) |\mathcal{I}|} < 1$,

we have $\mathcal{T}(\widehat{\mathbf{V}}_K^{(t+1)}, \mathbf{V}_K) \leq \rho \mathcal{T}(\widehat{\mathbf{V}}_K^{(t)}, \mathbf{V}_K)$. But the conditions (i) and (ii) are ensured by the inductive hypothesis and our assumptions, so the conclusion follows. \square

Proof of Theorem 4. Let $\mathbf{E} := \widetilde{\mathbf{G}} - \mathbb{E}^\Omega \widetilde{\mathbf{G}} = \widetilde{\mathbf{G}} - \boldsymbol{\Sigma}_y$. By Cape, Tang and Priebe (2018, Theorem 3.7), when $\lambda_K \geq 2\|\mathbf{E}\|_{\text{op}}$, we have that

$$\begin{aligned}\mathcal{T}(\widetilde{\mathbf{V}}_K, \mathbf{V}_K) &\leq 2\lambda_K^{-1} \|\mathbf{V}_{-K} \mathbf{V}_{-K}^\top \mathbf{E} \mathbf{V}_K \mathbf{V}_K^\top\|_{2 \rightarrow \infty} \\ &\quad + 2\lambda_K^{-1} \|\mathbf{V}_{-K} \mathbf{V}_{-K}^\top \mathbf{E} \mathbf{V}_{-K} \mathbf{V}_{-K}^\top\|_{2 \rightarrow \infty} \|\sin \Theta(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)\|_{\text{op}} \\ &\quad + 2\lambda_K^{-1} \|\mathbf{V}_{-K} \mathbf{V}_{-K}^\top \boldsymbol{\Sigma}_y \mathbf{V}_{-K} \mathbf{V}_{-K}^\top\|_{2 \rightarrow \infty} \|\sin \Theta(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)\|_{\text{op}} \\ &\quad + \|\sin \Theta(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)\|_{\text{op}}^2 \|\mathbf{V}_K\|_{2 \rightarrow \infty} \\ &=: T_1 + T_2 + T_3 + T_4.\end{aligned}$$

Note that if λ_K satisfies (27) for some $\xi > 1$, then $\mathbb{P}^\Omega(\|\mathbf{E}\|_{\text{op}} \geq \lambda_K/4) \leq 4d^{-(\xi-1)}$. In fact, since $\|\widetilde{\mathbf{W}}\|_{\text{op}} \leq \|\widetilde{\mathbf{W}}\|_{\text{F}}$, there exists $c_{M,\tau,\mu} > 0$ such that (7) implies (27), which, together with (26) ensures that $\mathbb{P}^\Omega(\|\mathbf{E}\|_{\text{op}} \geq \lambda_K/2) \leq 4d^{-(\xi-1)}$.

To bound T_1 , we have

$$\begin{aligned} \|\mathbf{V}_{-K}\mathbf{V}_{-K}^\top\mathbf{E}\mathbf{V}_K\mathbf{V}_K^\top\|_{2\rightarrow\infty} &\leq \|\mathbf{V}_{-K}\mathbf{V}_{-K}^\top\|_{\infty\rightarrow\infty}\|\mathbf{E}\mathbf{V}_K\mathbf{V}_K^\top\|_{2\rightarrow\infty} \\ &\leq (1+K\mu^2)\max_{j\in[d]}\sup_{\mathbf{u}\in\mathcal{S}^{K-1}}\mathbf{e}_j^\top\mathbf{E}\mathbf{V}_K\mathbf{u}, \end{aligned} \quad (20)$$

where the second inequality is due to the fact that

$$\|\mathbf{V}_{-K}\mathbf{V}_{-K}^\top\|_{\infty\rightarrow\infty} \leq \|\mathbf{I}_d\|_{\infty\rightarrow\infty} + \|\mathbf{V}_K\mathbf{V}_K^\top\|_{\infty\rightarrow\infty} \leq 1 + K\mu^2.$$

We use a covering argument to bound the supremum term. Let $\mathcal{N}_K(1/2)$ be a $1/2$ -net of the Euclidean sphere \mathcal{S}^{K-1} , i.e., for any $\mathbf{u} \in \mathcal{S}^{K-1}$, there exists a point $\pi(\mathbf{u}) \in \mathcal{N}_K(1/2)$ such that $\|\mathbf{u} - \pi(\mathbf{u})\|_2 \leq 1/2$. Note that for any $\mathbf{u} \in \mathcal{S}^{K-1}$,

$$\mathbf{e}_j^\top\mathbf{E}\mathbf{V}_K\mathbf{u} = \mathbf{e}_j^\top\mathbf{E}\mathbf{V}_K\pi(\mathbf{u}) + \mathbf{e}_j^\top\mathbf{E}\mathbf{V}_K(\mathbf{u} - \pi(\mathbf{u})) \leq \max_{\mathbf{v}\in\mathcal{N}_K(1/2)}\mathbf{e}_j^\top\mathbf{E}\mathbf{V}_K\mathbf{v} + \frac{1}{2}\sup_{\mathbf{v}\in\mathcal{S}^{K-1}}\mathbf{e}_j^\top\mathbf{E}\mathbf{V}_K\mathbf{v},$$

which further implies that

$$\sup_{\mathbf{u}\in\mathcal{S}^{K-1}}\mathbf{e}_j^\top\mathbf{E}\mathbf{V}_K\mathbf{u} \leq 2\max_{\mathbf{u}\in\mathcal{N}_K(1/2)}\mathbf{e}_j^\top\mathbf{E}\mathbf{V}_K\mathbf{u}. \quad (21)$$

We then argue similarly as in (23), with $\mathbf{V}_K\mathbf{u}$ taking the role of \mathbf{v}_k there (since $\|\mathbf{V}_K\mathbf{u}\|_\infty \leq \mu(K/d)^{1/2}$, we correspondingly have $\sqrt{K}\mu$ taking the role of the incoherence parameter μ there), to obtain that for any $\xi > 0$,

$$\mathbb{P}^\Omega\left\{|\mathbf{e}_j^\top\mathbf{E}\mathbf{V}_K\mathbf{u}| \geq 2e\tau\mu\left(\frac{KMR}{d}\right)^{1/2}\left(\frac{\xi^{1/2}\|\widetilde{\mathbf{W}}_j\|_1^{1/2}}{n^{1/2}} + \frac{\xi\|\widetilde{\mathbf{W}}_j\|_2}{n}\right)\right\} \leq 2e^{-\xi}.$$

By Vershynin (2012, Lemma 5.2), $|\mathcal{N}_K(1/2)| \leq 5^K$. Hence, by (20), (21) and a union bound, we have for any $\xi > \log 5$ that

$$\mathbb{P}^\Omega\left\{T_1 \geq \frac{8\tau\mu(1+K\mu^2)}{\lambda_K}\left(\frac{KMR}{d}\right)^{1/2}\left(\frac{\xi^{1/2}\|\widetilde{\mathbf{W}}\|_{\infty\rightarrow\infty}^{1/2}}{n^{1/2}} + \frac{\xi\|\widetilde{\mathbf{W}}\|_{2\rightarrow\infty}}{n}\right)\right\} \leq 2de^{K\log 5 - \xi}.$$

Next we bound T_2 . Note that

$$\|\mathbf{V}_{-K}\mathbf{V}_{-K}^\top\mathbf{E}\mathbf{V}_{-K}\mathbf{V}_{-K}^\top\|_{2\rightarrow\infty} \leq \|\mathbf{V}_{-K}\mathbf{V}_{-K}^\top\|_{\infty\rightarrow\infty}\|\mathbf{E}\|_{2\rightarrow\infty} \leq (1+K\mu^2)\|\mathbf{E}\|_{2\rightarrow\infty}.$$

For $j, k \in [d]$, let $\mathcal{I}_{jk} := \{i : \omega_{ij}\omega_{ik} = 1\}$ and $n_{jk} := |\mathcal{I}_{jk}| = n/\widetilde{W}_{jk}$. Then

$$E_{jk} = \frac{1}{n}\sum_{i=1}^n \widetilde{y}_{ij}\widetilde{y}_{ik}\widetilde{W}_{jk} - [\mathbb{E}^\Omega\widetilde{\mathbf{G}}]_{jk} = \frac{1}{n_{jk}}\sum_{i\in\mathcal{I}_{jk}} y_{ij}y_{ik} - [\mathbb{E}^\Omega\widetilde{\mathbf{G}}]_{jk}.$$

By applying both parts of Lemma 1, for any $i \in [n]$ and $j, k \in [d]$, we have $\|y_{ij}y_{ik}\|_{\psi_1} \leq 2\|y_{ij}\|_{\psi_2}\|y_{ik}\|_{\psi_2} \leq 2M$. Applying Bernstein's inequality (Boucheron, Lugosi and Massart, 2013, Theorem 2.10) yields that for any $\xi > 0$,

$$\mathbb{P}^\Omega\left\{|E_{jk}| \geq 2eM\left(\left(\frac{2\xi\widetilde{W}_{jk}}{n}\right)^{1/2} + \frac{\xi\widetilde{W}_{jk}}{n}\right)\right\} \leq 2e^{-\xi}.$$

Therefore, a union bound with $(j, k) \in [d] \times [d]$ yields that

$$\mathbb{P}^\Omega \left\{ T_2 \geq \frac{4\sqrt{2}eM(1 + K\mu^2)}{\lambda_K} \left(\left(\frac{2\xi \|\widetilde{\mathbf{W}}\|_{\infty \rightarrow \infty}}{n} \right)^{1/2} + \frac{\xi \|\widetilde{\mathbf{W}}\|_{2 \rightarrow \infty}}{n} \right) \|\sin \Theta(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)\|_{\text{op}} \right\} \leq 2d^2 e^{-\xi}.$$

Now we bound T_3 . We have that

$$T_3 = \frac{2\|\mathbf{V}_{-K}\mathbf{V}_{-K}^\top\|_{2 \rightarrow \infty}}{\lambda_K} \|\sin \Theta(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)\|_{\text{op}} \leq \frac{2\{1 + \mu(K/d)^{1/2}\}}{\lambda_K} \|\sin \Theta(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)\|_{\text{op}}.$$

Finally, T_4 satisfies

$$T_4 \leq \frac{\mu K^{1/2}}{d^{1/2}} \|\sin \Theta(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)\|_{\text{op}}^2.$$

Since $\|\sin \Theta(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)\|_{\text{op}} \leq \min\{L(\widetilde{\mathbf{V}}_K, \mathbf{V}_K), 1\}$, combining our bounds for $\{T_j\}_{j=1}^4$ yields that there exists $C_{M,\tau,\mu} > 0$ such that for any $\xi > 2$,

$$\begin{aligned} \mathbb{P}^\Omega \left\{ \mathcal{T}(\widetilde{\mathbf{V}}_K, \mathbf{V}_K) \geq \frac{KC_{M,\tau,\mu}}{\lambda_K} \left\{ L(\widetilde{\mathbf{V}}_K, \mathbf{V}_K) + \left(\frac{KR}{d} \right)^{1/2} \right\} \left(\frac{\xi^{1/2} \|\widetilde{\mathbf{W}}\|_{\infty \rightarrow \infty}^{1/2}}{n^{1/2}} + \frac{\xi \|\widetilde{\mathbf{W}}\|_{2 \rightarrow \infty}}{n} \right) \right. \\ \left. + \mu \left(\frac{K^{1/2}}{d^{1/2}} + \frac{4}{\lambda_K} \right) L(\widetilde{\mathbf{V}}_K, \mathbf{V}_K) \right\} \leq 2de^{K \log 5 - \xi} + 2d^2 e^{-\xi} + 4d^{-(\xi-1)}. \end{aligned}$$

It therefore follows from Proposition 2, which applies because condition (7) for a suitable $c_{M,\tau,\mu}$ implies (8), together with the facts that $\|\widetilde{\mathbf{W}}\|_1 \leq d\|\widetilde{\mathbf{W}}\|_{\infty \rightarrow \infty}$ and $\|\widetilde{\mathbf{W}}\|_{\text{F}} \leq d^{1/2}\|\widetilde{\mathbf{W}}\|_{2 \rightarrow \infty}$ that the first conclusion of the theorem holds. The second conclusion then follows immediately. \square

Proof of Proposition 2. In this proof, we use the shorthand $\mathbf{D}_{\mathbf{u}} := \text{diag}(\mathbf{u})$ for $\mathbf{u} \in \mathbb{R}^d$. We represent $\widetilde{\mathbf{G}}$ under the orthonormal basis $(\mathbf{V}_K, \mathbf{V}_{-K})$ as follows:

$$\widetilde{\mathbf{G}} = (\mathbf{V}_K, \mathbf{V}_{-K}) \begin{pmatrix} \mathbf{V}_K^\top \widetilde{\mathbf{G}} \mathbf{V}_K & \mathbf{V}_K^\top \widetilde{\mathbf{G}} \mathbf{V}_{-K} \\ \mathbf{V}_{-K}^\top \widetilde{\mathbf{G}} \mathbf{V}_K & \mathbf{V}_{-K}^\top \widetilde{\mathbf{G}} \mathbf{V}_{-K} \end{pmatrix} \begin{pmatrix} \mathbf{V}_K^\top \\ \mathbf{V}_{-K}^\top \end{pmatrix}.$$

Define

$$\mathbf{G}^* := (\mathbf{V}_K, \mathbf{V}_{-K}) \begin{pmatrix} \mathbf{V}_K^\top \widetilde{\mathbf{G}} \mathbf{V}_K & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_{-K}^\top \widetilde{\mathbf{G}} \mathbf{V}_{-K} \end{pmatrix} \begin{pmatrix} \mathbf{V}_K^\top \\ \mathbf{V}_{-K}^\top \end{pmatrix}.$$

In the sequel, we regard $\widetilde{\mathbf{G}}$ as a corrupted version of \mathbf{G}^* with the off-diagonal blocks $\mathbf{V}_K^\top \widetilde{\mathbf{G}} \mathbf{V}_{-K}$ and $\mathbf{V}_{-K}^\top \widetilde{\mathbf{G}} \mathbf{V}_K$ as perturbations. We have

$$\|\mathbf{V}_K^\top \widetilde{\mathbf{G}} \mathbf{V}_{-K}\|_{\text{F}} = \|\mathbf{V}_K^\top (\widetilde{\mathbf{G}} - \Sigma_{\mathbf{y}}) \mathbf{V}_{-K}\|_{\text{F}} \leq \|\mathbf{V}_K^\top (\widetilde{\mathbf{G}} - \mathbb{E}^\Omega \widetilde{\mathbf{G}})\|_{\text{F}}$$

We control the right-hand side through a concentration inequality, and for $k \in [K]$ let \mathbf{v}_k denote the k th column of \mathbf{V}_K . For any $j \in [d]$ and $k \in [K]$,

$$\begin{aligned} \mathbf{v}_k^\top (\tilde{\mathbf{G}} - \mathbb{E}^\Omega \tilde{\mathbf{G}}) \mathbf{e}_j &= \frac{1}{n} \sum_{i=1}^n \mathbf{v}_k^\top \{ \tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^\top \circ \tilde{\mathbf{W}} - \mathbb{E}^\Omega (\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^\top \circ \tilde{\mathbf{W}}) \} \mathbf{e}_j \\ &= \frac{1}{n} \sum_{i=1}^n \{ \tilde{\mathbf{y}}_i^\top \mathbf{D}_{\mathbf{v}_k} \tilde{\mathbf{W}} \mathbf{D}_{\mathbf{e}_j} \tilde{\mathbf{y}}_i - \mathbb{E}^\Omega (\tilde{\mathbf{y}}_i^\top \mathbf{D}_{\mathbf{v}_k} \tilde{\mathbf{W}} \mathbf{D}_{\mathbf{e}_j} \tilde{\mathbf{y}}_i) \} \\ &= \frac{1}{n} \sum_{i=1}^n \{ \tilde{y}_{ij} \tilde{\mathbf{y}}_i^\top \mathbf{D}_{\mathbf{v}_k} \tilde{\mathbf{W}}_j - \mathbb{E}^\Omega (\tilde{y}_{ij} \tilde{\mathbf{y}}_i^\top \mathbf{D}_{\mathbf{v}_k} \tilde{\mathbf{W}}_j) \}, \end{aligned} \quad (22)$$

where $\tilde{\mathbf{W}}_j$ denotes the j th column of $\tilde{\mathbf{W}}$. Note that

$$\|\mathbf{y}_i\|_{\psi_2^*} \leq \sup_{\mathbf{v} \in \mathcal{S}^{d-1}} \frac{\|\mathbf{v}^\top \mathbf{V}_K \mathbf{u}_i\|_{\psi_2} + \|\mathbf{v}^\top \mathbf{z}_i\|_{\psi_2}}{\sqrt{\mathbf{v}^\top \mathbf{V}_K \Sigma_{\mathbf{u}} \mathbf{V}_K^\top \mathbf{v} + 1}} \leq 2\tau.$$

Thus for any vector $\mathbf{a} \in \mathbb{R}^d$, we have by Lemma 1 that

$$\|y_{ij}(\mathbf{a}^\top \mathbf{y}_i)\|_{\psi_1} \leq 2\|y_{ij}\|_{\psi_2} \|\mathbf{a}^\top \mathbf{y}_i\|_{\psi_2} \leq 4\tau (M \mathbf{a}^\top \Sigma_{\mathbf{y}} \mathbf{a})^{1/2}.$$

For $i \in [n]$, let $\mathbf{a}_i := \omega_{ij} \tilde{\mathbf{W}}_j \circ \mathbf{v}_k \circ \omega_i$. Now for any $q \geq 2$,

$$\begin{aligned} \mathbb{E}^\Omega |\tilde{y}_{ij}(\tilde{\mathbf{W}}_j^\top \mathbf{D}_{\mathbf{v}_k} \tilde{\mathbf{y}}_i)|^q &= \mathbb{E}^\Omega |y_{ij} \mathbf{a}_i^\top \mathbf{y}_i|^q \leq \left(4q\tau \sqrt{M \mathbf{a}_i^\top \Sigma_{\mathbf{y}} \mathbf{a}_i}\right)^q \\ &\leq \frac{16q^q \tau^2 \mu^2 MR}{d} \left(4\tau\mu \sqrt{MR \|\tilde{\mathbf{W}}_j\|_2^2/d}\right)^{q-2} \sum_{t=1}^d \tilde{W}_{tj}^2 \omega_{it} \omega_{ij} \\ &\leq \frac{8e^2 q! \tau^2 \mu^2 MR}{d} \left(4e\tau\mu \sqrt{MR \|\tilde{\mathbf{W}}_j\|_2^2/d}\right)^{q-2} \sum_{t=1}^d \tilde{W}_{tj}^2 \omega_{it} \omega_{ij}, \end{aligned}$$

where the penultimate inequality uses the fact that $\|\mathbf{a}_i\|_2^2 \leq \mu^2 \|\tilde{\mathbf{W}}_j\|_2^2/d$, and the last inequality is due to Stirling's approximation. Hence,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{E}^\Omega |\tilde{y}_{ij} \tilde{\mathbf{W}}_j^\top \mathbf{D}_{\mathbf{v}_k} \tilde{\mathbf{y}}_i|^q &\leq \frac{8e^2 q! \tau^2 \mu^2 MR}{d} \left(4e\tau\mu \sqrt{MR \|\tilde{\mathbf{W}}_j\|_2^2/d}\right)^{q-2} \sum_{t=1}^d \sum_{i=1}^n \frac{\tilde{W}_{jt}^2 \omega_{it} \omega_{ij}}{n} \\ &= \frac{8e^2 q! \tau^2 \mu^2 MR}{d} \left(4e\tau\mu \sqrt{MR \|\tilde{\mathbf{W}}_j\|_2^2/d}\right)^{q-2} \|\tilde{\mathbf{W}}_j\|_1. \end{aligned}$$

Thus by (22) and Bernstein's inequality (Boucheron, Lugosi and Massart, 2013, Theorem 2.10), we have that for any $\xi > 0$,

$$\mathbb{P}^\Omega \left\{ |\mathbf{v}_k^\top (\tilde{\mathbf{G}} - \mathbb{E}^\Omega \tilde{\mathbf{G}}) \mathbf{e}_j| \geq 2^{5/2} e\tau\mu \left(\frac{MR}{d}\right)^{1/2} \left(\left(\frac{\xi \|\tilde{\mathbf{W}}_j\|_1}{n}\right)^{1/2} + \frac{\xi \|\tilde{\mathbf{W}}_j\|_2}{n} \right) \right\} \leq 2e^{-\xi}. \quad (23)$$

By a union bound over $(j, k) \in [d] \times [K]$, for any $\xi > 1$,

$$\mathbb{P}^\Omega \left\{ \|\mathbf{V}_K^\top \tilde{\mathbf{G}} \mathbf{V}_{-K}\|_F \geq 8e\tau\mu \left(\frac{KMR}{d} \right)^{1/2} \left(\frac{\xi^{1/2} \|\tilde{\mathbf{W}}\|_1^{1/2} \log^{1/2} d}{n^{1/2}} + \frac{\xi \|\tilde{\mathbf{W}}\|_F \log d}{n} \right) \right\} \leq 2Kd^{-(\xi-1)}. \quad (24)$$

Now we provide a condition under which $\lambda_{\min}(\mathbf{V}_K^\top \tilde{\mathbf{G}} \mathbf{V}_K) > \|\mathbf{V}_{-K}^\top \tilde{\mathbf{G}} \mathbf{V}_{-K}\|_{\text{op}}$, which ensures that \mathbf{V}_K is the top K eigenspace of \mathbf{G}^* . Note that

$$\lambda_{\min}(\mathbf{V}_K^\top \tilde{\mathbf{G}} \mathbf{V}_K) \geq \lambda_K + 1 - \|\mathbf{V}_K^\top (\tilde{\mathbf{G}} - \Sigma_y) \mathbf{V}_K\|_{\text{op}} \geq \lambda_K + 1 - \|\tilde{\mathbf{G}} - \Sigma_y\|_{\text{op}}$$

and

$$\|\mathbf{V}_{-K}^\top \tilde{\mathbf{G}} \mathbf{V}_{-K}\|_{\text{op}} \leq 1 + \|\tilde{\mathbf{G}} - \Sigma_y\|_{\text{op}}.$$

This implies that if $\lambda_K > 4\|\tilde{\mathbf{G}} - \Sigma_y\|_{\text{op}}$, then

$$\lambda_{\min}(\mathbf{V}_K^\top \tilde{\mathbf{G}} \mathbf{V}_K) - \|\mathbf{V}_{-K}^\top \tilde{\mathbf{G}} \mathbf{V}_{-K}\|_{\text{op}} > \lambda_K/2. \quad (25)$$

In the following, we derive an exponential tail bound for $\|\tilde{\mathbf{G}} - \Sigma_y\|_{\text{op}} = \|\tilde{\mathbf{G}} - \mathbb{E}^\Omega \tilde{\mathbf{G}}\|_{\text{op}}$. Let $\mathbf{A}_i := \tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^\top \circ \tilde{\mathbf{W}}$ and note that $\|\mathbf{A}_i\|_{\text{op}} \leq \|\mathbf{y}_i\|_\infty^2 \|\tilde{\mathbf{W}}\|_{\text{op}}$. Thus, for any $\mathbf{v} = (v_1, \dots, v_d)^\top \in \mathcal{S}^{d-1}$ and any integer $q \geq 2$,

$$\begin{aligned} \mathbb{E}^\Omega(\mathbf{v}^\top |\mathbf{A}_i|^{q\mathbf{v}}) &\leq \mathbb{E}^\Omega(\|\mathbf{A}_i\|_{\text{op}}^{q-2} \mathbf{v}^\top \mathbf{A}_i^2 \mathbf{v}) \leq \mathbb{E}^\Omega\{(\|\tilde{\mathbf{W}}\|_{\text{op}} \|\mathbf{y}_i\|_\infty^2)^{q-2} \mathbf{v}^\top (\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^\top \circ \tilde{\mathbf{W}})^2 \mathbf{v}\} \\ &= \|\tilde{\mathbf{W}}\|_{\text{op}}^{q-2} \mathbb{E}^\Omega\{\|\mathbf{y}_i\|_\infty^{2(q-2)} \mathbf{v}^\top \mathbf{D}_{\tilde{\mathbf{y}}_i} \tilde{\mathbf{W}} \mathbf{D}_{\tilde{\mathbf{y}}_i} \mathbf{D}_{\tilde{\mathbf{y}}_i} \tilde{\mathbf{W}} \mathbf{D}_{\tilde{\mathbf{y}}_i} \mathbf{v}\} \\ &= \|\tilde{\mathbf{W}}\|_{\text{op}}^{q-2} \mathbb{E}^\Omega\{\|\mathbf{y}_i\|_\infty^{2(q-2)} \text{tr}(\mathbf{D}_{\tilde{\mathbf{y}}_i}^2 \tilde{\mathbf{W}} \mathbf{D}_{\tilde{\mathbf{y}}_i} \mathbf{v} \mathbf{v}^\top \mathbf{D}_{\tilde{\mathbf{y}}_i} \tilde{\mathbf{W}})\} \\ &= \|\tilde{\mathbf{W}}\|_{\text{op}}^{q-2} \sum_{j=1}^d \omega_{ij} \mathbb{E}^\Omega\{y_{ij}^2 \|\mathbf{y}_i\|_\infty^{2(q-2)} (\tilde{\mathbf{W}}_j^\top \mathbf{D}_{\mathbf{v}} \tilde{\mathbf{y}}_i)^2\}. \end{aligned}$$

Now, for each $j \in [d]$, and $q \geq 2$,

$$\begin{aligned} \mathbb{E}^\Omega\{y_{ij}^2 \|\mathbf{y}_i\|_\infty^{2(q-2)} (\tilde{\mathbf{W}}_j^\top \mathbf{D}_{\mathbf{v}} \tilde{\mathbf{y}}_i)^2\} &= \mathbb{E}^\Omega[y_{ij}^2 \|\mathbf{y}_i\|_\infty^{2(q-2)} \{(\tilde{\mathbf{W}}_j \circ \mathbf{v} \circ \boldsymbol{\omega}_i)^\top \mathbf{y}_i\}^2] \\ &\leq (\mathbb{E} y_{ij}^8)^{1/4} \{\mathbb{E}(\|\mathbf{y}_i\|_\infty^{8(q-2)})\}^{1/4} 8R\tau^2 \|\tilde{\mathbf{W}}_j \circ \mathbf{v} \circ \boldsymbol{\omega}_i\|_2^2 \\ &\lesssim MR\tau^2 \{8(q-2)CM \log d\}^{q-2} \sum_{t=1}^d (v_t \tilde{W}_{tj} \omega_{it})^2, \end{aligned}$$

where the last inequality is due to the fact that $\|\|\mathbf{y}_i\|_\infty\|_{\psi_2} \leq (CM \log d)^{1/2}$ by (10). Therefore,

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}^\Omega(\mathbf{v}^\top |\mathbf{A}_i|^{q\mathbf{v}}) &\lesssim MR\tau^2 \{8(q-2)CM \|\tilde{\mathbf{W}}\|_{\text{op}} \log d\}^{q-2} \sum_{j,t=1}^d \sum_{i=1}^n \omega_{ij} \omega_{it} v_t^2 \tilde{W}_{tj}^2 \\ &= nMR\tau^2 \{8(q-2)CM \|\tilde{\mathbf{W}}\|_{\text{op}} \log d\}^{q-2} \sum_{j,t=1}^d v_t^2 \tilde{W}_{tj} \\ &\lesssim q! nMR\tau^2 \|\tilde{\mathbf{W}}^\top\|_{1 \rightarrow 1} (8eCM \|\tilde{\mathbf{W}}\|_{\text{op}} \log d)^{q-2}, \end{aligned}$$

where $\|\widetilde{\mathbf{W}}^\top\|_{1 \rightarrow 1} = \sup_{\|\mathbf{u}\|_1=1} \|\widetilde{\mathbf{W}}^\top \mathbf{u}\|_1 = \|\widetilde{\mathbf{W}}\|_{1 \rightarrow 1}$. Since the above inequality holds for all $\mathbf{v} \in \mathcal{S}^{d-1}$, we have

$$\left\| \sum_{i=1}^n \mathbb{E}^\Omega(|\mathbf{A}_i|^q) \right\|_{\text{op}} \lesssim q! n M R \tau^2 \|\widetilde{\mathbf{W}}\|_{1 \rightarrow 1} (8eCM \|\widetilde{\mathbf{W}}\|_{\text{op}} \log d)^{q-2}.$$

By a version of the Matrix Bernstein inequality for non-central absolute moments, which we give as Lemma 3, there exists a universal constant $C_1 > 0$ such that for any $\xi > 1$,

$$\mathbb{P}^\Omega \left\{ \|\widetilde{\mathbf{G}} - \mathbb{E}^\Omega \widetilde{\mathbf{G}}\|_{\text{op}} \geq C_1 \left(\left(\frac{MR\tau^2 \|\widetilde{\mathbf{W}}\|_{1 \rightarrow 1} \xi \log d}{n} \right)^{1/2} + \frac{M \|\widetilde{\mathbf{W}}\|_{\text{op}} \xi \log^2 d}{n} \right) \right\} \leq 4d^{-(\xi-1)}. \quad (26)$$

Now let

$$\mathcal{A} := \left\{ \lambda_{\min}(\mathbf{V}_K^\top \widetilde{\mathbf{G}} \mathbf{V}_K) - \|\mathbf{V}_{-K}^\top \widetilde{\mathbf{G}} \mathbf{V}_{-K}\|_{\text{op}} > \frac{\lambda_K}{2} \right\}.$$

From (25) and (26), we deduce that for any $\xi > 1$, if

$$\lambda_K \geq 4C_1 \left\{ \left(\frac{MR\tau^2 \|\widetilde{\mathbf{W}}\|_{1 \rightarrow 1} \xi \log d}{n} \right)^{1/2} + \frac{M \|\widetilde{\mathbf{W}}\|_{\text{op}} \xi \log^2 d}{n} \right\}, \quad (27)$$

then $\mathbb{P}^\Omega(\mathcal{A}^c) \leq \mathbb{P}^\Omega\{\|\widetilde{\mathbf{G}} - \Sigma_{\mathbf{y}}\|_{\text{op}} \geq \lambda_K/4\} \leq 4d^{-(\xi-1)}$. The desired result follows immediately by combining this with (24) and applying Yu, Wang and Samworth (2015, Theorem 2). \square

8 Auxiliary results

Lemma 1. *Let X and Y be two sub-Gaussian random variables. Then we have $\|X\|_{\psi_2}^2 \leq \|X^2\|_{\psi_1}$ and $\|XY\|_{\psi_1} \leq 2\|X\|_{\psi_2}\|Y\|_{\psi_2}$.*

Proof. For any $x \geq 0$, let $\lceil x \rceil := \inf\{z \in \mathbb{N} : z \geq x\}$. According to the definitions of the ψ_1 -norm and ψ_2 -norm, we have that

$$\|X\|_{\psi_2}^2 = \sup_{p \in \mathbb{N}} \frac{\mathbb{E}(|X|^p)^{2/p}}{p} \leq \sup_{p \in \mathbb{N}} \frac{\{\mathbb{E}(X^{2\lceil p/2 \rceil})\}^{1/\lceil p/2 \rceil}}{p} \leq \|X^2\|_{\psi_1},$$

where the penultimate inequality is due to Jensen's inequality and the last inequality is due to the fact that $p \geq \lceil p/2 \rceil$. For the second inequality,

$$\begin{aligned} \|XY\|_{\psi_1} &= \sup_{p \in \mathbb{N}} \frac{(\mathbb{E}|XY|^p)^{1/p}}{p} \leq 2 \sup_{p \in \mathbb{N}} \frac{(\mathbb{E}|X|^{2p})^{1/(2p)} (\mathbb{E}|Y|^{2p})^{1/(2p)}}{\sqrt{2p}} \\ &\leq 2 \sup_{p \in \mathbb{N}} \frac{(\mathbb{E}|X|^{2p})^{1/(2p)}}{\sqrt{2p}} \sup_{q \in \mathbb{N}} \frac{(\mathbb{E}|Y|^{2q})^{1/(2q)}}{\sqrt{2q}} \leq 2\|X\|_{\psi_2}\|Y\|_{\psi_2}, \end{aligned}$$

as required. \square

Lemma 2. *If X_1, \dots, X_n are independent centred random variables with $\max_{i \in [n]} \|X_i\|_{\psi_1} < \infty$, then there exists a universal constant $C > 0$ such that*

$$\left\| \sum_{i=1}^n X_i \right\|_{\psi_1} \leq C \left(\sum_{i=1}^n \|X_i\|_{\psi_1}^2 \right)^{1/2}.$$

Proof. Write $K_i := \|X_i\|_{\psi_1}$ and $\mathbf{K} := (K_1, \dots, K_n)^\top$. From [Vershynin \(2012, Lemma 5.15\)](#), there exist universal constants $c_1, C_1 > 0$ such that for $|t| \leq c_1/\|\mathbf{K}\|_\infty$,

$$\mathbb{E} \exp \left\{ t \sum_{i=1}^n X_i \right\} = \prod_{i=1}^n \mathbb{E} \exp \{ t X_i \} \leq \exp \{ C_1 t^2 \|\mathbf{K}\|_2^2 \}.$$

Setting $t = \min\{C_1^{-1/2}\|\mathbf{K}\|_2^{-1}, c_1\|\mathbf{K}\|_\infty^{-1}\}$ in the above expression, the right-hand side is bounded above by e . The desired result follows from the fact that (5.15) and (5.16) in [Vershynin \(2012\)](#) are two definitions that yield equivalent ψ_1 -norms. \square

The following lemma provides a variant of the existing matrix Bernstein inequality ([Tropp, 2012, Theorem 6.2](#)). The primary difference is that we impose non-central absolute moment inequalities, as opposed to central moment inequalities. We believe that this inequality may be of independent interest, with applications beyond the scope of this paper. Recall that if $\mathbf{A} \in \mathbb{S}^{d \times d}$, with eigendecomposition $\mathbf{A} = \mathbf{Q} \text{diag}(\mu_1, \dots, \mu_d) \mathbf{Q}^\top$ for some orthogonal $\mathbf{Q} \in \mathbb{R}^{d \times d}$, then we write $|\mathbf{A}| := \mathbf{Q} \text{diag}(|\mu_1|, \dots, |\mu_d|) \mathbf{Q}^\top$.

Lemma 3 (Matrix Bernstein inequality with non-central absolute moment conditions). *Let $\{\mathbf{X}_i\}_{i \in [n]}$ be independent symmetric $d \times d$ random matrices. Assume that*

$$\mathbb{E}(|\mathbf{X}_i|^q) \preceq \frac{q!}{2} R^{q-2} \mathbf{A}_i^2 \quad \text{for } q = 2, 3, 4, \dots$$

for some $R > 0$ and deterministic d -dimensional symmetric matrices $\{\mathbf{A}_i\}_{i \in [n]}$. Define the variance parameter

$$\sigma^2 := \left\| \sum_{i=1}^n \mathbf{A}_i^2 \right\|_{\text{op}}.$$

Then for each $t > 0$,

$$\mathbb{P} \left[\lambda_{\max} \left\{ \sum_{i=1}^n (\mathbf{X}_i - \mathbb{E} \mathbf{X}_i) \right\} \geq t \right] \leq 4d \exp \left(\frac{-t^2/32}{\sigma^2 + Rt} \right).$$

Proof. Let $\tilde{\mathbf{X}}_1, \dots, \tilde{\mathbf{X}}_n, \epsilon_1, \dots, \epsilon_n$ be independent random matrices and variables, independent of $(\mathbf{X}_1, \dots, \mathbf{X}_n)$, satisfying $\tilde{\mathbf{X}}_i \stackrel{d}{=} \mathbf{X}_i$ and $\epsilon_i \sim U(\{-1, 1\})$ for $i \in [n]$. Write $\mathbf{S}_n := \sum_{i=1}^n (\mathbf{X}_i - \mathbb{E} \mathbf{X}_i)$ and $\tilde{\mathbf{S}}_n := \sum_{i=1}^n (\tilde{\mathbf{X}}_i - \mathbb{E} \mathbf{X}_i)$. Given $\mathbf{X}_1, \dots, \mathbf{X}_n$, let $\mathbf{v}_* = \mathbf{v}_*(\mathbf{X}_1, \dots, \mathbf{X}_n)$ be a leading unit-length eigenvector of \mathbf{S}_n . Let $\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_d$ denote orthonormal eigenvectors of

$\tilde{\mathbf{X}}_1$ with corresponding eigenvalues $\tilde{\mu}_1, \dots, \tilde{\mu}_d$; fix $\mathbf{v} \in \mathcal{S}^{d-1}$, and let $w_j := (\tilde{\mathbf{v}}_j^\top \mathbf{v})^2$ for $j \in [d]$. Since $\sum_{j=1}^d w_j = 1$, we have by Jensen's inequality that for $q \in \{2, 3, \dots\}$,

$$|\mathbf{v}^\top \tilde{\mathbf{X}}_1 \mathbf{v}|^q = \left| \sum_{j=1}^d w_j \tilde{\mu}_j \right|^q \leq \sum_{j=1}^d w_j |\tilde{\mu}_j|^q = \mathbf{v}^\top |\tilde{\mathbf{X}}_1|^q \mathbf{v}.$$

We deduce that $\mathbb{E}\{(\mathbf{v}^\top \tilde{\mathbf{X}}_i \mathbf{v})_+^q\} \leq \mathbb{E}\{|\mathbf{v}^\top \tilde{\mathbf{X}}_i \mathbf{v}|^q\} \leq \frac{q!}{2} R^{q-2} \mathbf{v}^\top \mathbf{A}_i^2 \mathbf{v}$ for $i \in [n]$, so by Bernstein's inequality (Boucheron, Lugosi and Massart, 2013, Corollary 2.11),

$$\mathbb{P}(\mathbf{v}_*^\top \tilde{\mathbf{S}}_n \mathbf{v}_* > t/2 \mid \mathbf{X}_1, \dots, \mathbf{X}_n) \leq \exp\left(\frac{-t^2/8}{\mathbf{v}_*^\top \sum_{i=1}^n \mathbf{A}_i^2 \mathbf{v}_* + Rt}\right) \leq \exp\left(\frac{-t^2/8}{\sigma^2 + Rt}\right).$$

We may assume that the right-hand side of the above inequality is at most $1/2$, since otherwise the lemma is trivially true. Therefore,

$$\begin{aligned} \mathbb{P}\{\lambda_{\max}(\mathbf{S}_n) \geq t\} &= \mathbb{P}(\mathbf{v}_*^\top \mathbf{S}_n \mathbf{v}_* \geq t) \leq 2\mathbb{E}\left\{\mathbb{P}(\mathbf{v}_*^\top \tilde{\mathbf{S}}_n \mathbf{v}_* \leq t/2 \mid \mathbf{X}_1, \dots, \mathbf{X}_n) \mathbf{1}_{\{\mathbf{v}_*^\top \mathbf{S}_n \mathbf{v}_* \geq t\}}\right\} \\ &= 2\mathbb{P}(\mathbf{v}_*^\top \tilde{\mathbf{S}}_n \mathbf{v}_* \leq t/2 \text{ and } \mathbf{v}_*^\top \mathbf{S}_n \mathbf{v}_* \geq t) \leq 2\mathbb{P}(\mathbf{v}_*^\top (\mathbf{S}_n - \tilde{\mathbf{S}}_n) \mathbf{v}_* \geq t/2) \\ &\leq 2\mathbb{P}\left[\lambda_{\max}\left\{\sum_{i=1}^n \epsilon_i (\mathbf{X}_i - \tilde{\mathbf{X}}_i)\right\} \geq t/2\right] \leq 4\mathbb{P}\left\{\lambda_{\max}\left(\sum_{i=1}^n \epsilon_i \mathbf{X}_i\right) \geq t/4\right\}, \end{aligned} \tag{28}$$

where we have used the fact that $\epsilon_i (\mathbf{X}_i - \tilde{\mathbf{X}}_i) \stackrel{d}{=} \mathbf{X}_i - \tilde{\mathbf{X}}_i$ for all i in the penultimate inequality.

Since $\mathbb{E}(\epsilon_i \mathbf{X}_i) = \mathbf{0}$ and $\mathbb{E}\{(\epsilon_i \mathbf{X}_i)^q\} \preceq \mathbb{E}\{|\mathbf{X}_i|^q\} \preceq \frac{q!}{2} R^{q-2} \mathbf{A}_i^2$ for $q \in \{2, 3, \dots\}$, applying the matrix Bernstein inequality (Tropp, 2012, Theorem 6.2) to the sequence $\{\epsilon_i \mathbf{X}_i\}_{i \in [n]}$ yields

$$\mathbb{P}\left\{\lambda_{\max}\left(\sum_{i=1}^n \epsilon_i \mathbf{X}_i\right) \geq t/4\right\} \leq d \exp\left(\frac{-t^2/32}{\sigma^2 + Rt}\right).$$

We attain the conclusion by combining the above inequality with (28). \square

Lemma 4. *Let X_1, \dots, X_n be independent $\text{Bin}(d, p)$ random variables and let $\hat{p}_i := X_i/d$. When $dp \geq 1$ and $n \geq 2$, we have*

$$\mathbb{E} \max_{i \in [n]} \hat{p}_i \leq 10p \log n.$$

Proof. By Bernstein's inequality (van der Vaart and Wellner, 1996, Lemma 2.2.9) and a union bound,

$$\mathbb{P}\left(\max_{i \in [n]} \hat{p}_i \geq p + t\right) \leq n \exp\left(-\frac{dt^2}{2(p + t/3)}\right).$$

Setting $t_0 := 2\sqrt{pd^{-1} \log n} + \frac{4}{3d} \log n$, we have

$$\mathbb{E} \max_{i \in [n]} \hat{p}_i = p + t_0 + \int_{t_0}^{\infty} n \{e^{-dt^2/(4p)} + e^{-3dt/4}\} dt \leq p + t_0 + \sqrt{\frac{\pi p}{d}} + \frac{4}{3d} \leq 10p \log n,$$

where we have used $\log n \geq \log 2$ and $1/d \leq p$ in the final inequality. \square

Lemma 5. Suppose that $\boldsymbol{\beta}, \boldsymbol{\eta} \in \mathbb{R}^d$ and $\|\boldsymbol{\eta}\|_2 = \|\boldsymbol{\beta}\|_2$. Let $\boldsymbol{\Sigma}_1 := \mathbf{I}_d + \boldsymbol{\beta}\boldsymbol{\beta}^\top$ and $\boldsymbol{\Sigma}_2 := \mathbf{I}_d + \boldsymbol{\eta}\boldsymbol{\eta}^\top$. Then

$$\text{KL}(N_d(\mathbf{0}, \boldsymbol{\Sigma}_1), N_d(\mathbf{0}, \boldsymbol{\Sigma}_2)) = \frac{\|\boldsymbol{\eta}\|_2^4 - (\boldsymbol{\eta}^\top \boldsymbol{\beta})^2}{2(1 + \|\boldsymbol{\eta}\|_2^2)}.$$

Proof. Since $\|\boldsymbol{\eta}\|_2 = \|\boldsymbol{\beta}\|_2$, the matrices $\boldsymbol{\Sigma}_1$ and $\boldsymbol{\Sigma}_2$ share the same set of eigenvalues. Hence $|\boldsymbol{\Sigma}_1| = |\boldsymbol{\Sigma}_2|$ and we have

$$\text{KL}(N_d(\mathbf{0}, \boldsymbol{\Sigma}_1), N_d(\mathbf{0}, \boldsymbol{\Sigma}_2)) = \frac{1}{2} \{ \text{tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) - d \} = \frac{1}{2} \{ \text{tr}((\mathbf{I}_d + \boldsymbol{\eta}\boldsymbol{\eta}^\top)^{-1}(\mathbf{I}_d + \boldsymbol{\beta}\boldsymbol{\beta}^\top)) - d \}.$$

Now, by the Sherman–Morrison formula,

$$(\mathbf{I}_d + \boldsymbol{\eta}\boldsymbol{\eta}^\top)^{-1} = \mathbf{I}_d - \frac{\boldsymbol{\eta}\boldsymbol{\eta}^\top}{1 + \|\boldsymbol{\eta}\|_2^2}$$

and thus we have

$$\begin{aligned} \text{KL}(N(\mathbf{0}, \boldsymbol{\Sigma}_1), N(\mathbf{0}, \boldsymbol{\Sigma}_2)) &= \frac{1}{2} \left[\text{tr} \left(\left(\mathbf{I}_d - \frac{\boldsymbol{\eta}\boldsymbol{\eta}^\top}{1 + \|\boldsymbol{\eta}\|_2^2} \right) (\mathbf{I}_d + \boldsymbol{\beta}\boldsymbol{\beta}^\top) \right) - d \right] \\ &= \frac{1}{2} \left(\|\boldsymbol{\beta}\|_2^2 - \frac{\|\boldsymbol{\eta}\|_2^2}{1 + \|\boldsymbol{\eta}\|_2^2} - \frac{(\boldsymbol{\eta}^\top \boldsymbol{\beta})^2}{1 + \|\boldsymbol{\eta}\|_2^2} \right) = \frac{\|\boldsymbol{\eta}\|_2^4 - (\boldsymbol{\eta}^\top \boldsymbol{\beta})^2}{2(1 + \|\boldsymbol{\eta}\|_2^2)}, \end{aligned}$$

as required. \square

Theorem 4 and Proposition 2 exhibit bounds on $\mathcal{T}(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)$ and $L(\widetilde{\mathbf{V}}_K, \mathbf{V}_K)$ given a deterministic observation scheme. To provide some intuition on the size of these bounds under the p -homogeneous missingness setting described in Section 2.1, the following lemma derives probabilistic bounds for various norms of $\widetilde{\mathbf{W}}$.

Lemma 6. Assume (A5). Then there exists an event of probability at least $1 - d^2 e^{-3np^2/32}$ on which each of the following inequalities hold:

- (i) $\|\widetilde{\mathbf{W}}\|_{\text{op}} \leq 2dp^{-2}$;
- (ii) $\|\widetilde{\mathbf{W}}\|_{1 \rightarrow 1} = \|\widetilde{\mathbf{W}}\|_{\infty \rightarrow \infty} \leq 2dp^{-2}$;
- (iii) $\|\widetilde{\mathbf{W}}\|_1 \leq 2d^2 p^{-2}$;
- (iv) $\|\widetilde{\mathbf{W}}\|_{\text{F}} \leq 2dp^{-2}$;
- (v) $\|\widetilde{\mathbf{W}}\|_{2 \rightarrow \infty} \leq 2d^{1/2} p^{-2}$.

Proof. Define an event

$$\mathcal{A} := \{ \|\widetilde{\mathbf{W}} - p^{-2} \mathbf{1}_d \mathbf{1}_d^\top\|_\infty \leq p^{-2} \}.$$

For $j, k \in [d]$, write $\widehat{P}_{jk} := n^{-1} \sum_{i=1}^n \omega_{ij} \omega_{ik}$. Then by a union bound and Bernstein's inequality (Wainwright, 2019, Proposition 2.14), we have

$$\mathbb{P}(\mathcal{A}^c) \leq \sum_{j, k \in [d]} \mathbb{P}(\widehat{P}_{jk} < p^2/2) \leq d^2 e^{-3np^2/32}.$$

Note that on \mathcal{A} , we have $\|\widetilde{\mathbf{W}}\|_\infty \leq 2p^{-2}$. The desired bounds then follow respectively from the following inequalities: $\|\widetilde{\mathbf{W}}\|_{\text{op}} \leq d\|\widetilde{\mathbf{W}}\|_\infty$, $\|\widetilde{\mathbf{W}}\|_{1 \rightarrow 1} = \|\widetilde{\mathbf{W}}\|_{\infty \rightarrow \infty} \leq d\|\widetilde{\mathbf{W}}\|_\infty$, $\|\widetilde{\mathbf{W}}\|_1 \leq d^2\|\widetilde{\mathbf{W}}\|_\infty$, $\|\widetilde{\mathbf{W}}\|_{\text{F}} \leq d\|\widetilde{\mathbf{W}}\|_\infty$ and $\|\widetilde{\mathbf{W}}\|_{2 \rightarrow \infty} \leq d^{1/2}\|\widetilde{\mathbf{W}}\|_\infty$. \square

References

- Anderson, T. W. (1957) Maximum likelihood estimates for a multivariate normal distribution when some observations are missing. *J. Amer. Statist. Assoc.*, **52**, 200–203.
- Belloni, A., Rosenbaum, M. and Tsybakov, A. B. (2017) Linear and conic programming estimators in high dimensional errors-in-variables models. *J. Roy. Statist. Soc., Ser. B*, **79**, 939–956.
- Boucheron, S., Lugosi, G. and Massart, P. (2013) *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, Oxford.
- Cai, T. T. and Zhang, A. (2018a) Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics. *Ann. Statist.*, **46**, 60–89.
- Cai, T. T. and Zhang, L. (2018b) High-dimensional linear discriminant analysis: optimality, adaptive algorithm, and missing data. [arXiv:1804.03018](https://arxiv.org/abs/1804.03018).
- Candès, E. J., Li, X., Ma, Y. and Wright, J. (2011) Robust principal component analysis? *J. ACM*, **58**, 11:1–11:37.
- Candès, E. J. and Plan, Y. (2010) Matrix completion with noise. *Proc. IEEE*, **98**, 925–936.
- Candès, E. J. and Recht, B. (2009) Exact matrix completion via convex optimization. *Found. Comput. Math.*, **9**, 717–772.
- Cape, J., Tang, M. and Priebe, C. E. (2018) The two-to-infinity norm and singular subspace geometry with applications to high-dimensional statistics. *Ann. Statist.*, to appear.
- Chi Y., Lu Y. and Chen Y. (2018) Nonconvex optimization meets low-rank matrix factorization: An overview. [arXiv preprint arXiv:1809.09573](https://arxiv.org/abs/1809.09573).
- Cho, J., Kim, D. and Rohe, K. (2017) Asymptotic theory for estimating the singular vectors and values of a partially-observed low rank matrix with noise. *Statist. Sinica*, **27**, 1921–1948.
- Dempster, A. P., Laird, N. M. and Rubin, D. B. (1977) Maximum likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc., Ser. B*, **39**, 1–38.
- Dray, S. and Josse J. (2015) Principal component analysis with missing values: a comparative survey of methods. *Plant Ecol.*, **216**, 657–667.
- Elsener, A and van de Geer, S. (2018) Sparse spectral estimation with missing and corrupted measurements. [arXiv:1811.10443](https://arxiv.org/abs/1811.10443).
- Large covariance estimation by thresholding principal orthogonal complements. *J. Roy. Statist. Soc., Ser. B*, **75**, 603–680

- Ford, B. L. (1983) An overview of hot-deck procedures. In W. G. Madow, I. Olkin and D. B. Rubin (Eds.) *Incomplete Data in Sample Surveys, Vol. 2: Theory and Bibliographies*, 185–207. Academic Press, New York.
- Gao, C., Ma, Z., Zhang, A. Y. and Zhou, H. H. (2016) Achieving optimal misclassification proportion in stochastic block models. *J. Mach. Learn. Res.*, **18**, 1–45.
- Hastie T., Mazumder R., Lee J. D. and Zadeh R. (2015) Matrix completion and low-rank SVD via fast alternating least squares. *J. Mach. Learn. Res.*, **16**, 3367–3402.
- Johnstone, I. M. and Lu, A. Y. (2009) On consistency and sparsity for principal components analysis in high dimensions. *J. Amer. Statist. Assoc.*, **104**, 682–693.
- Josse J. and Husson F. (2012) Handling missing values in exploratory multivariate data analysis methods. *J. de la Société Française de Statistique*, **153**, 1–21.
- Josse J., Pagès J. and Husson F. (2009) Gestion des données manquantes en analyse en composantes principales. *J. de la Société Française de Statistique*, **150**, 28–51.
- Keshavan, R. H., Montanari, A. and Oh, S. (2010) Matrix completion from a few entries. *IEEE Trans. Inform. Theory*, **56**, 2980–2998.
- Kiers H. A. L. (1997) Weighted least squares fitting using ordinary least squares algorithms. *Psychometrika*, **62**, 251–266.
- Koltchinskii, V., Lounici, K. and Tsybakov, A. B. (2011) Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *Ann. Statist.*, **39**, 2302–2329.
- Little, R. J. and Rubin, D. B. (2014) *Statistical Analysis with Missing Data*, John Wiley & Sons, Hoboken.
- Loh, P.-L. and Tan, X. L. (2018) High-dimensional robust precision matrix estimation: Cell-wise corruption under ϵ -contamination. *Electron. J. Statist.*, **12**, 1429–1467.
- Loh, P.-L. and Wainwright, M. J. (2012) High-dimensional regression with noisy and missing data: provable guarantees with nonconvexity. *Ann. Statist.*, **40**, 1637–1664.
- Lounici, K. (2013) Sparse principal component analysis with missing observations. In C. Houdré et al. (Eds.) *High Dimensional Probability VI*, 327–356. Birkhäuser, Basel.
- Lounici, K. (2014) High-dimensional covariance matrix estimation with missing observations. *Bernoulli*, **20**, 1029–1058.
- Massart, P. (2007) *Concentration Inequalities and Model Selection*, Springer, Berlin.
- Mazumder, R., Hastie, T. and Tibshirani R. (2010) Spectral regularization algorithms for learning large incomplete matrices. *J. Mach. Learn. Res.*, **11**, 2287–2322.

- Negahban, S. and Wainwright, M. J. (2012) Restricted strong convexity and weighted matrix completion: optimal bounds with noise. *J. Mach. Learn. Res.*, **13**, 1665–1697.
- Rohe, K., Chatterjee, S. and Yu, B. (2011) Spectral clustering and the high-dimensional stochastic blockmodel. *Ann. Statist.*, **39**, 1878–1915.
- Rubin, D. B. (2004) *Multiple Imputation for Nonresponse in Surveys*. John Wiley & Sons, Hoboken.
- Shen, D., Shen, H., Zhu, H. and Marron, J. (2016) The statistics and mathematics of high dimension low sample size asymptotics. *Statist. Sinica*, **26**, 1747–1770.
- Tropp, J. A. (2012) User-friendly tail bounds for sums of random matrices. *Found. Comput. Math.*, **12**, 389–434.
- van der Vaart, A. W. and Wellner, J. A. (1996) *Weak Convergence and Empirical Processes*. Springer, New York.
- Vershynin, R. (2012) Introduction to the non-asymptotic analysis of random matrices. In Y. Eldar and G. Kutyniok (Eds.) *Compressed Sensing, Theory and Applications*. Cambridge University Press, Cambridge. 210–268.
- Wainwright, M. J. (2019) *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press, Cambridge.
- Wang, W. and Fan, J. (2017) Asymptotics of empirical eigen-structure for ultra-high dimensional spiked covariance model. *Ann. Statist.*, **45**, 1342–1374.
- Wold, H. and Lyttkens, E. (1969) Nonlinear iterative partial least squares (NIPALS) estimation procedures. *Bull. Int. Stat. Inst.*, **43**, 29–51.
- Yu, B. (1997) Assouad, Fano and Le Cam. In *Festschrift for Lucien Le Cam: Research Papers in Probability and Statistics*, Pollard, D., Torgersen, E. and Yang G. L. (Eds.), 423–435. Springer, New York.
- Yu, Y., Wang, T. and Samworth, R. J. (2015) A useful variant of the Davis–Kahan theorem for statisticians. *Biometrika*, **102**, 315–323.
- Zhang, A., Cai, T. T. and Wu, Y. (2018) Heteroskedastic PCA: Algorithm, optimality, and applications. [arXiv:1810.08316](https://arxiv.org/abs/1810.08316).
- Zhu, Z., Wang, T. and Samworth, R. J. (2019) `primePCA`: projected refinement for imputation of missing entries in principal component analysis. R package version 1.0. <https://cran.r-project.org/web/packages/primePCA/>.