

## STATISTICAL AND COMPUTATIONAL TRADE-OFFS IN ESTIMATION OF SPARSE PRINCIPAL COMPONENTS

BY TENG YAO WANG<sup>\*,1</sup>, QUENTIN BERTHET<sup>\*,†,2</sup>  
AND RICHARD J. SAMWORTH<sup>\*,3</sup>

*University of Cambridge*<sup>\*</sup> and *California Institute of Technology*<sup>†</sup>

In recent years, sparse principal component analysis has emerged as an extremely popular dimension reduction technique for high-dimensional data. The theoretical challenge, in the simplest case, is to estimate the leading eigenvector of a population covariance matrix under the assumption that this eigenvector is sparse. An impressive range of estimators have been proposed; some of these are fast to compute, while others are known to achieve the minimax optimal rate over certain Gaussian or sub-Gaussian classes. In this paper, we show that, under a widely-believed assumption from computational complexity theory, there is a fundamental trade-off between statistical and computational performance in this problem. More precisely, working with new, larger classes satisfying a restricted covariance concentration condition, we show that there is an effective sample size regime in which no randomised polynomial time algorithm can achieve the minimax optimal rate. We also study the theoretical performance of a (polynomial time) variant of the well-known semidefinite relaxation estimator, revealing a subtle interplay between statistical and computational efficiency.

**1. Introduction.** Principal Component Analysis (PCA), which involves projecting a sample of multivariate data onto the space spanned by the leading eigenvectors of the sample covariance matrix, is one of the oldest and most widely-used dimension reduction devices in statistics. It has proved to be particularly effective when the dimension of the data is relatively small by comparison with the sample size. However, the work of [Johnstone and Lu \(2009\)](#) and [Paul \(2007\)](#) shows

---

**Tribute:** Peter was a remarkable person: not only a prolific and highly influential researcher, but also someone with a wonderful warmth and generosity of spirit. He was a great inspiration to so many statisticians around the world. We are deeply saddened that he is no longer with us, and dedicate this paper to his memory. Further personal reflections on Peter Hall's life and work from the third author can be found in [Samworth \(2016\)](#).

Received May 2015; revised July 2015.

<sup>1</sup>Supported by a Benefactors' scholarship from St. John's College, Cambridge.

<sup>2</sup>Supported by Air Force Office of Scientific Research (AFOSR) Grant FA9550-14-1-0098 at the Center for the Mathematics of Information at the California Institute of Technology.

<sup>3</sup>Supported by Engineering and Physical Sciences Research Council Early Career Fellowship EP/J017213/1 and Leverhulme Trust Grant PLP-2014-353.

*MSC2010 subject classifications.* 62H25, 68Q17.

*Key words and phrases.* Computational lower bounds, planted clique problem, polynomial time algorithm, sparse principal component analysis.

that PCA breaks down in the high-dimensional settings that are frequently encountered in many diverse modern application areas. For instance, consider the spiked covariance model where  $X_1, \dots, X_n$  are independent  $N_p(0, \Sigma)$  random vectors, with  $\Sigma = I_p + \theta v_1 v_1^\top$  for some  $\theta > 0$  and an arbitrary unit vector  $v_1 \in \mathbb{R}^p$ . In this case,  $v_1$  is the leading eigenvector (principal component) of  $\Sigma$ , and the classical PCA estimate would be  $\hat{v}_1$ , a unit-length leading eigenvector of the sample covariance matrix  $\hat{\Sigma} := n^{-1} \sum_{i=1}^n X_i X_i^\top$ . In the high-dimensional setting where  $p = p_n$  is such that  $p/n \rightarrow c \in (0, 1)$ , Paul (2007) showed that

$$|\hat{v}_1^\top v_1| \xrightarrow{\text{a.s.}} \begin{cases} \sqrt{\frac{1 - c/\theta^2}{1 + c/\theta}}, & \text{if } \theta > \sqrt{c}, \\ 0, & \text{if } \theta \leq \sqrt{c}. \end{cases}$$

In other words,  $\hat{v}_1$  is inconsistent as an estimator of  $v_1$  in this asymptotic regime. This phenomenon is related to the so-called ‘‘BBP’’ transition in random matrix theory [Baik, Ben Arous and P ech e (2005)].

Sparse principal component analysis was designed to remedy this inconsistency and to give additional interpretability to the projected data. In the simplest case, it is assumed that the leading eigenvector  $v_1$  of the population covariance matrix  $\Sigma$  belongs to the  $k$ -sparse unit Euclidean sphere in  $\mathbb{R}^p$ , given by

$$(1) \quad B_0(k) := \left\{ u = (u_1, \dots, u_p)^\top \in \mathbb{R}^p : \sum_{j=1}^p \mathbb{1}_{\{u_j \neq 0\}} \leq k, \|u\|_2 = 1 \right\}.$$

A remarkable number of recent papers have proposed estimators of  $v_1$  in this setting, including Jolliffe, Trendafilov and Uddin (2003), Zou, Hastie and Tibshirani (2006), d’Aspremont et al. (2007), Johnstone and Lu (2009), Witten, Tibshirani and Hastie (2009), Journ ee et al. (2010), Birnbaum et al. (2013), Cai, Ma and Wu (2013), Ma (2013), Shen, Shen and Marron (2013) and Vu and Lei (2013).

Sparse PCA methods have gained high popularity in many diverse applied fields where high-dimensional datasets are routinely handled. These include computer vision for online visual tracking [Wang, Lu and Yang (2013)] and pattern recognition [Naikal, Yang and Sastry (2011)], signal processing for image compression [Majumdar (2009)] and electrocardiography feature extraction [Johnstone and Lu (2009)], and biomedical research for gene expression analysis [Chan and Hall (2010), Chun and S und uz (2009), Parkhomenko, Tritchler and Beyene (2009), Zou, Hastie and Tibshirani (2006)], RNA-seq classification [Tan, Petersen and Witten (2014)] and metabolomics studies [Allen and Maletic-Savatic (2011)]. In these applications, sparse PCA is employed to identify a small number of interpretable directions that represent the data succinctly, typically as the first stage of a more involved procedure such as classification, clustering or regression.

The success of the ultimate inferential methods in the types of application described above depends critically on how well the particular sparse PCA technique

involved identifies the relevant meaningful directions in the underlying population. It therefore becomes important to understand the ways in which our ability to estimate these directions from data depends on the characteristics of the problem, including the sample size, dimensionality, sparsity level and signal-to-noise ratio. Such results form a key component of any theoretical analysis of an inference problem in which sparse PCA is employed as a first step.

In terms of the theoretical properties of existing methods for sparse PCA, Ma (2013) was able to show that his estimator attains the minimax rate of convergence over a certain Gaussian class of distributions, provided that  $k$  is treated as a fixed constant. Both Cai, Ma and Wu (2013) and Vu and Lei (2013) also study minimax properties, but treat  $k$  as a parameter of the problem that may vary with the sample size  $n$ . In particular, for a certain class  $\mathcal{P}_p(n, k)$  of sub-Gaussian distributions and in a particular asymptotic regime, Vu and Lei (2013) show<sup>4</sup> that

$$\inf_{\hat{v}} \sup_{P \in \mathcal{P}_p(n, k)} \mathbb{E}_P \{1 - (v_1^\top \hat{v})^2\} \asymp \frac{k \log p}{n},$$

where the infimum is taken over all estimators  $\hat{v}$ ; see also Birnbaum et al. (2013). Moreover, they show that the minimax rate is attained by a leading  $k$ -sparse eigenvector of  $\hat{\Sigma}$ , given by

$$(2) \quad \hat{v}_{\max}^k \in \operatorname{argmax}_{u \in B_0(k)} u^\top \hat{\Sigma} u.$$

The papers cited above would appear to settle the question of sparse principal component estimation (at least in a sub-Gaussian setting) from the perspective of statistical theory. However, there remains an unsettling feature, namely that neither the estimator of Cai, Ma and Wu (2013), nor that of Vu and Lei (2013), is computable in polynomial time.<sup>5</sup> For instance, computing the estimator (2) is an NP-hard problem, and the naive algorithm that searches through all  $\binom{p}{k}$  of the  $k \times k$  principal submatrices of  $\hat{\Sigma}$  quickly becomes infeasible for even moderately large  $p$  and  $k$ .

Given that sparse PCA methods are typically applied to massive high-dimensional datasets, it is crucial to understand the rates that can be achieved using only computationally efficient procedures. Specifically, in this paper, we address the question of whether it is possible to find an estimator of  $v_1$  that is computable in (randomised) polynomial time, and that attains the minimax optimal rate of convergence when the sparsity of  $v_1$  is allowed to vary with the sample size. Some progress in a related direction was made by Berthet and Rigollet (2013a, 2013b),

<sup>4</sup>Here and below,  $a_n \asymp b_n$  means  $0 < \liminf_{n \rightarrow \infty} |a_n/b_n| \leq \limsup_{n \rightarrow \infty} |a_n/b_n| < \infty$ .

<sup>5</sup>Since formal definitions of such notions from computational complexity theory may be unfamiliar to many statisticians, and to keep the paper as self-contained as possible, we provide a brief introduction to this topic in Section 2 of the online supplementary material [Wang, Berthet and Samworth (2015)].

who considered the problem of testing the null hypothesis  $H_0 : \Sigma = I_p$  against the alternative  $H_1 : v^\top \Sigma v \geq 1 + \theta$  for some  $v \in B_0(k)$  and  $\theta > 0$ . Of interest here is the minimal level  $\theta = \theta_{n,p,k}$  that ensures small asymptotic testing error. Under a hypothesis on the computational intractability of a certain well-known problem from theoretical computer science (the ‘‘Planted Clique’’ detection problem), Berthet and Rigollet showed that for certain classes of distributions, there is a gap between the minimal  $\theta$ -level permitting successful detection with a randomised polynomial time test, and the corresponding  $\theta$ -level when arbitrary tests are allowed.

The particular classes of distributions considered in Berthet and Rigollet (2013a, 2013b) were highly tailored to the testing problem, and do not provide sufficient structure to study principal component estimation. The thesis of this paper, however, is that from the point of view of both theory and applications, it is the estimation of sparse principal components, rather than testing for the existence of a distinguished direction, that is the more natural and fundamental (as well as more challenging) problem. Indeed, we observe subtle phase transition phenomena that are absent from the hypothesis testing problem; see Section 4.4 for further details. It is worth noting that different results for statistical and computational trade-offs for estimation and testing were also observed in the context of  $k$ -SAT formulas in Feldman, Perkins and Vempala (2015) and Berthet (2015), respectively.

Our first contribution, in Section 2, is to introduce a new Restricted Covariance Concentration (RCC) condition that underpins the classes of distributions  $\mathcal{P}_p(n, k, \theta)$  over which we perform the statistical and computational analyses [see (4) for a precise definition]. The RCC condition is satisfied by sub-Gaussian distributions, and moreover has the advantage of being more robust to certain mixture contaminations that turn out to be of key importance in the statistical analysis under the computational constraint. We show that subject to mild restrictions on the parameter values,

$$\inf_{\hat{v}} \sup_{P \in \mathcal{P}_p(n,k,\theta)} \mathbb{E}_P L(\hat{v}, v_1) \asymp \sqrt{\frac{k \log p}{n\theta^2}},$$

where  $L(u, v) := \{1 - (u^\top v)^2\}^{1/2}$ , and where no restrictions are placed on the class of estimators  $\hat{v}$ . By contrast, in Section 3, we show that a variant  $\hat{v}^{\text{SDP}}$  of the semidefinite relaxation estimator of d’Aspremont et al. (2007) and Bach, Ahipasaoglu and d’Aspremont (2010), which is computable in polynomial time, satisfies

$$\sup_{P \in \mathcal{P}_p(n,k,\theta)} \mathbb{E}_P L(\hat{v}^{\text{SDP}}, v_1) \leq (16\sqrt{2} + 2) \sqrt{\frac{k^2 \log p}{n\theta^2}}.$$

Our main result, in Section 4, is that, under a much weaker planted clique hypothesis than that in Berthet and Rigollet (2013a, 2013b), for any  $\alpha \in (0, 1)$ , there exists

a moderate effective sample size asymptotic regime in which every sequence  $(\hat{v}^{(n)})$  of randomised polynomial time estimators satisfies

$$\sqrt{\frac{n\theta^2}{k^{1+\alpha} \log p}} \sup_{P \in \mathcal{P}_p(n,k,\theta)} \mathbb{E}_P L(\hat{v}^{(n)}, v_1) \rightarrow \infty.$$

This result shows that there is a fundamental trade-off between statistical and computational efficiency in the estimation of sparse principal components, and that there is in general no consistent sequence of randomised polynomial time estimators in this regime. Interestingly, in a high effective sample size regime, where even randomised polynomial time estimators can be consistent, we are able to show in Theorem 7 that under additional distributional assumptions, a modified (but still polynomial time) version of  $\hat{v}^{\text{SDP}}$  attains the minimax optimal rate. Thus, the trade-off disappears for a sufficiently high effective sample size, at least over a subset of the parameter space.

Statistical and computational trade-offs have also recently been studied in the context of convex relaxation algorithms [Chandrasekaran and Jordan (2013)], submatrix signal detection [Chen and Xu (2014), Ma and Wu (2015)], sparse linear regression [Zhang, Wainwright and Jordan (2014)], community detection [Hajek, Wu and Xu (2014)] and sparse canonical correlation analysis [Gao, Ma and Zhou (2014)]. Given the importance of computationally feasible algorithms with good statistical performance in today's era of big data, it seems clear that understanding the extent of this phenomenon in different settings will represent a key challenge for theoreticians in the coming years.

Proofs of our main results are given in the [Appendix](#), while several ancillary results are deferred to the online supplementary material [Wang, Berthet and Samworth (2015)]. We end this section by introducing some notation used throughout the paper. For a vector  $u = (u_1, \dots, u_M)^\top \in \mathbb{R}^M$ , a matrix  $A = (A_{ij}) \in \mathbb{R}^{M \times N}$  and for  $q \in [1, \infty)$ , we write  $\|u\|_q := (\sum_{i=1}^M |u_i|^q)^{1/q}$  and  $\|A\|_q := (\sum_{i=1}^M \sum_{j=1}^N |A_{ij}|^q)^{1/q}$  for their (entrywise)  $\ell_q$ -norms. We also write  $\|u\|_0 := \sum_{i=1}^M \mathbb{1}_{\{u_i \neq 0\}}$ ,  $\text{supp}(u) := \{i : u_i \neq 0\}$ ,  $\|A\|_0 := \sum_{i=1}^M \sum_{j=1}^N \mathbb{1}_{\{A_{ij} \neq 0\}}$  and  $\text{supp}(A) := \{(i, j) : A_{ij} \neq 0\}$ . For  $S \subseteq \{1, \dots, M\}$  and  $T \subseteq \{1, \dots, N\}$ , we write  $u_S := (u_i : i \in S)^\top$  and write  $M_{S,T}$  for the  $|S| \times |T|$  submatrix of  $M$  obtained by extracting the rows and columns with indices in  $S$  and  $T$ , respectively. For positive sequences  $(a_n)$  and  $(b_n)$ , we write  $a_n \ll b_n$  to mean  $a_n/b_n \rightarrow 0$ .

## 2. Restricted covariance concentration and minimax rate of estimation.

Let  $p \geq 2$  and let  $\mathcal{P}$  denote the class of probability distributions  $P$  on  $\mathbb{R}^p$  with  $\int_{\mathbb{R}^p} x dP(x) = 0$  and such that the entries of  $\Sigma(P) := \int_{\mathbb{R}^p} xx^\top dP(x)$  are finite. For  $P \in \mathcal{P}$ , write  $\lambda_1(P), \dots, \lambda_p(P)$  for the eigenvalues of  $\Sigma(P)$ , arranged in decreasing order. When  $\lambda_1(P) - \lambda_2(P) > 0$ , the first principal component  $v_1(P)$ , that is, a unit-length eigenvector of  $\Sigma$  corresponding to the eigenvalue  $\lambda_1(P)$ , is well defined up to sign. In some places below, and where it is clear from the

context, we suppress the dependence of these quantities on  $P$ , or write the eigenvalues and eigenvectors as  $\lambda_1(\Sigma), \dots, \lambda_p(\Sigma)$  and  $v_1(\Sigma), \dots, v_p(\Sigma)$ , respectively. Let  $X_1, \dots, X_n$  be independent and identically distributed random vectors with distribution  $P$ , and form the  $n \times p$  matrix  $\mathbf{X} := (X_1, \dots, X_n)^\top$ . An estimator of  $v_1$  is a measurable function from  $\mathbb{R}^{n \times p}$  to  $\mathbb{R}^p$ , and we write  $\mathcal{V}_{n,p}$  for the class of all such estimators.

Given unit vectors  $u, v \in \mathbb{R}^p$ , let  $\Theta(u, v) := \cos^{-1}(|u^\top v|)$  denote the acute angle between  $u$  and  $v$ , and define the loss function

$$L(u, v) := \sin \Theta(u, v) = \{1 - (u^\top v)^2\}^{1/2} = \frac{1}{\sqrt{2}} \|uu^\top - vv^\top\|_2.$$

Note that  $L(\cdot, \cdot)$  is invariant to sign changes of either of its arguments. The *directional variance* of  $P$  along a unit vector  $u \in \mathbb{R}^p$  is defined to be  $V(u) := \mathbb{E}\{(u^\top X_1)^2\} = u^\top \Sigma u$ . Its empirical counterpart is  $\hat{V}(u) := n^{-1} \sum_{i=1}^n (u^\top X_i)^2 = u^\top \hat{\Sigma} u$ , where  $\hat{\Sigma} := n^{-1} \sum_{i=1}^n X_i X_i^\top$  denotes the sample covariance matrix.

Recall the definition of the  $k$ -sparse unit ball  $B_0(k)$  from (1). Given  $\ell \in \{1, \dots, p\}$  and  $C \in (0, \infty)$ , we say  $P$  satisfies a *Restricted Covariance Concentration* (RCC) condition with parameters  $p, n, \ell$  and  $C$ , and write  $P \in \text{RCC}_p(n, \ell, C)$ , if

$$(3) \quad \mathbb{P} \left\{ \sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| \geq C \max \left( \sqrt{\frac{\ell \log(p/\delta)}{n}}, \frac{\ell \log(p/\delta)}{n} \right) \right\} \leq \delta$$

for all  $\delta > 0$ . It is also convenient to define

$$\text{RCC}_p(\ell, C) := \bigcap_{n=1}^\infty \text{RCC}_p(n, \ell, C) \quad \text{and} \quad \text{RCC}_p(C) := \bigcap_{\ell=1}^p \text{RCC}_p(\ell, C).$$

The RCC conditions amount to uniform Bernstein-type concentration properties of the directional variance around its expectation along all sparse directions. This condition turns out to be particularly convenient in the study of convergence rates in sparse PCA, and moreover, as we show in Proposition 1 below, sub-Gaussian distributions satisfy an RCC condition for all sample sizes  $n$  and all sparsity levels  $\ell$ . Recall that a mean-zero distribution  $Q$  on  $\mathbb{R}^p$  is *sub-Gaussian* with parameter<sup>6</sup>  $\sigma^2 \in (0, \infty)$ , written

$$Q \in \text{sub-Gaussian}_p(\sigma^2),$$

if whenever  $Y \sim Q$ , we have  $\mathbb{E}(e^{u^\top Y}) \leq e^{\sigma^2 \|u\|^2/2}$  for all  $u \in \mathbb{R}^p$ .

---

<sup>6</sup>Note that some authors say that distributions satisfying this condition are sub-Gaussian with parameter  $\sigma$ , rather than  $\sigma^2$ .

PROPOSITION 1. (i) For every  $\sigma > 0$ , we have

$$\text{sub-Gaussian}_p(\sigma^2) \subseteq \text{RCC}_p\left(16\sigma^2\left(1 + \frac{9}{\log p}\right)\right).$$

(ii) In the special case where  $P = N_p(0, \Sigma)$ , we have  $P \in \text{RCC}_p(8\lambda_1(P)(1 + \frac{9}{\log p}))$ .

Our convergence rate results for sparse principal component estimation will be proved over the following classes of distributions. For  $\theta > 0$ , let

$$(4) \quad \mathcal{P}_p(n, k, \theta) := \{P \in \text{RCC}_p(n, 2, 1) \cap \text{RCC}_p(n, 2k, 1) : v_1(P) \in B_0(k), \lambda_1(P) - \lambda_2(P) \geq \theta\}.$$

Observe that RCC classes have the scaling property that if the distribution of a random vector  $Y$  belongs to  $\text{RCC}_p(n, \ell, C)$  and if  $r > 0$ , then the distribution of  $rY$  belongs to  $\text{RCC}_p(n, \ell, r^2C)$ . It is therefore convenient to fix  $C = 1$  in both RCC classes in (4), so that  $\theta$  becomes a measure of the signal-to-noise level.

For a symmetric  $A \in \mathbb{R}^{p \times p}$ , define  $\hat{v}_{\max}^k(A) := \text{sargmax}_{u \in B_0(k)} u^\top Au$  to be the  $k$ -sparse maximum eigenvector of  $A$ , where  $\text{sargmax}$  denotes the smallest element of the argmax in the lexicographic ordering. [This choice ensures that  $\hat{v}_{\max}^k(A)$  is a measurable function of  $A$ .] Theorem 2 below gives a finite-sample minimax upper bound for estimating  $v_1(P)$  over  $\mathcal{P}_p(n, k, \theta)$ . For similar bounds over Gaussian or sub-Gaussian classes, see Cai, Ma and Wu (2013) and Vu and Lei (2013), who consider the more general problem of principal subspace estimation. As well as working with a larger class of distributions, our different proof techniques facilitate an explicit constant.

THEOREM 2. For  $2k \log p \leq n$ , the  $k$ -sparse empirical maximum eigenvector,  $\hat{v}_{\max}^k(\hat{\Sigma})$ , satisfies

$$\sup_{P \in \mathcal{P}_p(n, k, \theta)} \mathbb{E}_P L(\hat{v}_{\max}^k(\hat{\Sigma}), v_1(P)) \leq 2\sqrt{2}\left(1 + \frac{1}{\log p}\right) \sqrt{\frac{k \log p}{n\theta^2}} \leq 7\sqrt{\frac{k \log p}{n\theta^2}}.$$

A matching minimax lower bound of the same order in all parameters  $k, p, n$  and  $\theta$  is given below. The proof techniques are adapted from Vu and Lei (2013).

THEOREM 3. Suppose that  $7 \leq k \leq p^{1/2}$  and  $0 < \theta \leq \frac{1}{16(1+9/\log p)}$ . Then

$$\inf_{\hat{v} \in \mathcal{V}_{n,p}} \sup_{P \in \mathcal{P}_p(n, k, \theta)} \mathbb{E}_P L(\hat{v}, v_1(P)) \geq \min\left\{\frac{1}{1660} \sqrt{\frac{k \log p}{n\theta^2}}, \frac{5}{18\sqrt{3}}\right\}.$$

We remark that the conditions in the statement of Theorem 3 can be strengthened or weakened, with a corresponding weakening or strengthening of the con-



stands in the bound. For instance, a bound of the same order in  $k, p, n$  and  $\theta$  could be obtained assuming only that  $k \leq p^{1-\delta}$  for some  $\delta > 0$ . The upper bound on  $\theta$  is also not particularly restrictive. For example, if  $P = N_p(0, \sigma^2 I_p + \theta e_1 e_1^\top)$ , where  $e_1$  is the first standard basis vector in  $\mathbb{R}^p$ , then it can be shown that the condition  $P \in \mathcal{P}_p(n, k, \theta)$  requires that  $\theta \leq 1 - \sigma^2$ .

**3. Computationally efficient estimation.** As was mentioned in the [Introduction](#), the trouble with the estimator  $\hat{v}_{\max}^k(\hat{\Sigma})$  of Section 2, as well as the estimator of [Cai, Ma and Wu \(2013\)](#), is that there are no known polynomial time algorithms for their computation. In this section, we therefore study the (polynomial time) semidefinite relaxation estimator  $\hat{v}^{\text{SDP}}$  defined by Algorithm 1 below. This estimator is a variant of one proposed by [d’Aspremont et al. \(2007\)](#), whose support recovery properties were studied for a particular class of Gaussian distributions and a known sparsity level by [Amini and Wainwright \(2009\)](#).

To motivate the main step (Step 2) of Algorithm 1, it is convenient to let  $\mathcal{M}$  denote the class of  $p \times p$  nonnegative definite real, symmetric matrices, and let  $\mathcal{M}_1 := \{M \in \mathcal{M} : \text{tr}(M) = 1\}$ . Let  $\mathcal{M}_{1,1}(k^2) := \{M \in \mathcal{M}_1 : \text{rank}(M) = 1, \|M\|_0 = k^2\}$  and observe that

$$\max_{u \in B_0(k)} u^\top \hat{\Sigma} u = \max_{u \in B_0(k)} \text{tr}(\hat{\Sigma} u u^\top) = \max_{M \in \mathcal{M}_{1,1}(k^2)} \text{tr}(\hat{\Sigma} M).$$

In the final expression, the rank and sparsity constraints are nonconvex. We therefore adopt the standard semidefinite relaxation approach of dropping the rank constraint and replacing the sparsity ( $\ell_0$ ) constraint with an  $\ell_1$  penalty to obtain the convex optimisation problem

$$(5) \quad \max_{M \in \mathcal{M}_1} \{\text{tr}(\hat{\Sigma} M) - \lambda \|M\|_1\}.$$

We now discuss the complexity of computing  $\hat{v}^{\text{SDP}}$  in detail. One possible way of implementing Step 2 is to use a generic interior-point method. However,

**Algorithm 1:** Pseudo-code for computing the semidefinite relaxation estimator  $\hat{v}^{\text{SDP}}$

**Input:**  $\mathbf{X} = (X_1, \dots, X_n)^\top \in \mathbb{R}^{n \times p}$ ,  $\lambda > 0$ ,  $\varepsilon > 0$

**begin**

**Step 1:** Set  $\hat{\Sigma} \leftarrow n^{-1} \mathbf{X}^\top \mathbf{X}$ .

**Step 2:** For  $f(M) := \text{tr}(\hat{\Sigma} M) - \lambda \|M\|_1$ , let  $\hat{M}^\varepsilon$  be an  $\varepsilon$ -maximiser of  $f$  in  $\mathcal{M}_1$ . In other words,  $\hat{M}^\varepsilon$  satisfies  $f(\hat{M}^\varepsilon) \geq \max_{M \in \mathcal{M}_1} f(M) - \varepsilon$ .

**Step 3:** Let  $\hat{v}^{\text{SDP}} := \hat{v}_{\lambda, \varepsilon}^{\text{SDP}} \in \text{argmax}_{u: \|u\|_2=1} u^\top \hat{M}^\varepsilon u$ .

**end**

**Output:**  $\hat{v}^{\text{SDP}}$



---

**Algorithm 2:** A possible implementation of Step 2 of Algorithm 1

---

**Input:**  $\hat{\Sigma} \in \mathcal{M}, \lambda > 0, \varepsilon > 0.$

**begin**

Set  $M_0 \leftarrow I_p/p, U_0 \leftarrow 0 \in \mathbb{R}^{p \times p}$  and  $N \leftarrow \lceil \frac{\lambda^2 p^2 + 1}{\sqrt{2\varepsilon}} \rceil.$

**for**  $t \leftarrow 1$  **to**  $N$  **do**

$U'_t \leftarrow \Pi_{\mathcal{U}}(U_{t-1} - \frac{1}{\sqrt{2}}M_{t-1}), M'_t \leftarrow \Pi_{\mathcal{M}_1}(M_{t-1} + \frac{1}{\sqrt{2}}\hat{\Sigma} + \frac{1}{\sqrt{2}}U_{t-1}).$

$U_t \leftarrow \Pi_{\mathcal{U}}(U_{t-1} - \frac{1}{\sqrt{2}}M'_t), M_t \leftarrow \Pi_{\mathcal{M}_1}(M_{t-1} + \frac{1}{\sqrt{2}}\hat{\Sigma} + \frac{1}{\sqrt{2}}U'_t).$

**end**

Set  $\hat{M}^\varepsilon \leftarrow \frac{1}{N} \sum_{t=1}^N M'_t.$

**end**

**Output:**  $\hat{M}^\varepsilon$

---

as shown in [Nesterov \(2005\)](#), [Nemirovski \(2004\)](#) and [Bach, Ahipaşaoğlu and d’Aspremont \(2010\)](#), certain first-order algorithms [i.e., methods requiring  $O(1/\varepsilon)$  steps to find a feasible point achieving an  $\varepsilon$ -approximation of the optimal objective function value] can significantly outperform such generic interior-point solvers. The key idea in both [Nesterov \(2005\)](#) and [Nemirovski \(2004\)](#) is that the optimisation problem in Step 2 can be rewritten in a saddlepoint formulation:

$$\max_{M \in \mathcal{M}_1} \text{tr}(\hat{\Sigma}M) - \lambda \|M\|_1 = \max_{M \in \mathcal{M}_1} \min_{U \in \mathcal{U}} \text{tr}((\hat{\Sigma} + U)M),$$

where  $\mathcal{U} := \{U \in \mathbb{R}^{p \times p} : U^\top = U, \|U\|_\infty \leq \lambda\}$ . The fact that  $\text{tr}((\hat{\Sigma} + U)M)$  is linear in both  $M$  and  $U$  makes the problem amenable to proximal methods. In Algorithm 2 above, we state a possible implementation of Step 2 of Algorithm 1, derived from the “basic implementation” in [Nemirovski \(2004\)](#). In the algorithm, the  $\|\cdot\|_2$ -norm projection  $\Pi_{\mathcal{U}}(A)$  of a symmetric matrix  $A = (A_{ij}) \in \mathbb{R}^{p \times p}$  onto  $\mathcal{U}$  is given by

$$(\Pi_{\mathcal{U}}(A))_{ij} := \text{sign}(A_{ij}) \min(|A_{ij}|, \lambda).$$

For the projection  $\Pi_{\mathcal{M}_1}(A)$ , first decompose  $A =: PDP^\top$  for some orthogonal  $P$  and diagonal  $D = \text{diag}(d)$ , where  $d = (d_1, \dots, d_p)^\top \in \mathbb{R}^p$ . Now let  $\Pi_{\mathcal{W}}(d)$  be the projection image of  $d$  on the unit  $(p - 1)$ -simplex  $\mathcal{W} := \{(w_1, \dots, w_p) : w_j \geq 0, \sum_{j=1}^p w_j = 1\}$ . Finally, transform back to obtain  $\Pi_{\mathcal{M}_1}(A) := P \text{diag}(\Pi_{\mathcal{W}}(d))P^\top$ . The fact that Algorithm 2 outputs an  $\varepsilon$ -maximiser of the optimisation problem in Step 2 of Algorithm 1 follows from [Nemirovski \[\(2004\), Theorem 3.2\]](#), which implies in our particular case that after  $N$  iterations,

$$\max_{M \in \mathcal{M}_1} \min_{U \in \mathcal{U}} \text{tr}((\hat{\Sigma} + U)M) - \min_{U \in \mathcal{U}} \text{tr}((\hat{\Sigma} + U)\hat{M}^\varepsilon) \leq \frac{\lambda^2 p^2 + 1}{\sqrt{2N}}.$$

In Algorithm 1, Step 1 takes  $O(np^2)$  floating point operations; Step 3 takes  $O(p^3)$  operations in the worst case, though other methods such as the Lanczos method [Golub and Van Loan (1996), Lanczos (1950)] require only  $O(p^2)$  operations under certain conditions. Our particular implementation (Algorithm 2) for Step 2 requires  $O(\frac{\lambda^2 p^2 + 1}{\varepsilon})$  iterations in the worst case, though this number may often be considerably reduced by terminating the **for** loop if the primal-dual gap

$$\lambda_1(\hat{U}_t + \hat{\Sigma}) - \{\text{tr}(\hat{M}_t \hat{\Sigma}) - \lambda \|\hat{M}_t\|_1\}$$

falls below  $\varepsilon$ , where  $\hat{U}_t := t^{-1} \sum_{s=1}^t U'_s$  and  $\hat{M}_t := t^{-1} \sum_{s=1}^t M'_s$ . The most costly step within the **for** loop is the eigen-decomposition used to compute the projection  $\Pi_{\mathcal{M}_1}$ , which takes  $O(p^3)$  operations. Taking  $\lambda := 4\sqrt{\frac{\log p}{n}}$  and  $\varepsilon := \frac{\log p}{4n}$  as in Theorem 5 below, we find an overall complexity for the algorithm of  $O(\max(p^5, \frac{np^3}{\log p}))$  operations in the worst case.

We now turn to the theoretical properties of the estimator  $\hat{v}^{\text{SDP}}$  computed using Algorithm 1. Lemma 4 below is stated in a general, deterministic fashion, but will be used in Theorem 5 below to bound the loss incurred by the estimator on the event that the sample and population covariance matrices are close in  $\ell_\infty$ -norm. See also Vu et al. [(2013), Theorem 3.1] for a closely related result in the context of a projection matrix estimation problem. Recall that  $\mathcal{M}$  denotes the class of  $p \times p$  nonnegative definite real, symmetric matrices.

LEMMA 4. *Let  $\Sigma \in \mathcal{M}$  be such that  $\theta := \lambda_1(\Sigma) - \lambda_2(\Sigma) > 0$ . Let  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and  $\hat{\Sigma} := n^{-1} \mathbf{X}^\top \mathbf{X}$ . For arbitrary  $\lambda > 0$  and  $\varepsilon > 0$ , if  $\|\hat{\Sigma} - \Sigma\|_\infty \leq \lambda$ , then the semidefinite relaxation estimator  $\hat{v}^{\text{SDP}}$  in Algorithm 1 with inputs  $\mathbf{X}, \lambda, \varepsilon$  satisfies*

$$L(\hat{v}^{\text{SDP}}, v_1(\Sigma)) \leq \frac{4\sqrt{2}\lambda k}{\theta} + 2\sqrt{\frac{\varepsilon}{\theta}}.$$

Theorem 5 below describes the statistical properties of the estimator  $\hat{v}^{\text{SDP}}$  over  $\mathcal{P}_p(n, k, \theta)$  classes. It reveals in particular that we incur a loss of statistical efficiency of a factor of  $\sqrt{k}$  compared with the minimax upper bound in Theorem 2 in Section 2 above. As well as applying Lemma 4 on the event  $\{\|\hat{\Sigma} - \Sigma\|_\infty \leq \lambda\}$ , the proof relies on Lemma 5 in the online supplementary material [Wang, Berthet and Samworth (2015)], which relates the event  $\{\|\hat{\Sigma} - \Sigma\|_\infty > \lambda\}$  to the  $\text{RCC}_p(n, 2, 1)$  condition. Indeed, this explains why we incorporated this condition into the definition of the  $\mathcal{P}_p(n, k, \theta)$  classes.

THEOREM 5. *For an arbitrary  $P \in \mathcal{P}_p(n, k, \theta)$  and  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} P$ , we write  $\hat{v}^{\text{SDP}}(\mathbf{X})$  for the output of Algorithm 1 with input  $\mathbf{X} := (X_1, \dots, X_n)^\top$ ,  $\lambda := 4\sqrt{\frac{\log p}{n}}$  and  $\varepsilon := \frac{\log p}{4n}$ . If  $4 \log p \leq n \leq k^2 p^2 \theta^{-2} \log p$  and  $\theta \in (0, k]$ , then*

$$(6) \quad \sup_{P \in \mathcal{P}_p(n, k, \theta)} \mathbb{E}_P L(\hat{v}^{\text{SDP}}(\mathbf{X}), v_1(P)) \leq \min \left\{ (16\sqrt{2} + 2) \sqrt{\frac{k^2 \log p}{n\theta^2}}, 1 \right\}.$$

We remark that  $\hat{v}^{\text{SDP}}$  has the attractive property of being fully adaptive in the sense that it can be computed without knowledge of the sparsity level  $k$ . On the other hand,  $\hat{v}^{\text{SDP}}$  is not necessarily  $k$ -sparse. If a specific sparsity level is desired in a particular application, Algorithm 1 can be modified to obtain a (nonadaptive)  $k$ -sparse estimator having similar estimation risk. Specifically, we can find

$$\hat{v}_0^{\text{SDP}} \in \underset{u \in B_0(k)}{\operatorname{argmin}} L(\hat{v}^{\text{SDP}}, u).$$

Since  $L(\hat{v}^{\text{SDP}}, u)^2 = 1 - (u^\top \hat{v}^{\text{SDP}})^2$ , we can compute  $\hat{v}_0^{\text{SDP}}$  by setting all but the top  $k$  coordinates of  $\hat{v}^{\text{SDP}}$  in absolute value to zero and renormalising the vector. In particular,  $\hat{v}_0^{\text{SDP}}$  is computable in polynomial time. We deduce that under the same conditions as in Theorem 5, for any  $P \in \mathcal{P}_p(n, k, \theta)$ ,

$$\begin{aligned} \mathbb{E}L(\hat{v}_0^{\text{SDP}}, v_1) &\leq \mathbb{E}[\{L(\hat{v}_0^{\text{SDP}}, \hat{v}^{\text{SDP}}) + L(\hat{v}^{\text{SDP}}, v_1)\} \mathbb{1}_{\{\|\hat{\Sigma} - \Sigma\|_\infty \leq \lambda\}}] + \mathbb{P}(\|\hat{\Sigma} - \Sigma\|_\infty > \lambda) \\ &\leq 2\mathbb{E}\{L(\hat{v}_0^{\text{SDP}}, v_1) \mathbb{1}_{\{\|\hat{\Sigma} - \Sigma\|_\infty \leq \lambda\}}\} + \mathbb{P}(\|\hat{\Sigma} - \Sigma\|_\infty > \lambda) \\ &\leq (32\sqrt{2} + 3)\sqrt{\frac{k^2 \log p}{n\theta^2}}, \end{aligned}$$

where the final inequality follows from the proof of Theorem 5.

**4. Computational lower bounds in sparse principal component estimation.**

Theorems 5 and 2 reveal a gap between the provable performance of our semidefinite relaxation estimator  $\hat{v}^{\text{SDP}}$  and the minimax optimal rate. It is natural to ask whether there exists a computationally efficient algorithm that achieves the statistically optimal rate of convergence. In fact, as we will see in Theorem 6 below, the effective sample size region over which  $\hat{v}^{\text{SDP}}$  is consistent is essentially tight among the class of all *randomised polynomial time algorithms*.<sup>7</sup> Indeed, any randomised polynomial time algorithm with a faster rate of convergence could otherwise be adapted to solve instances of the planted clique problem that are believed to be hard; see Section 4.1 below for formal definitions and discussion. In this sense, the extra factor of  $\sqrt{k}$  is an intrinsic price in statistical efficiency that we have to pay for computational efficiency, and the estimator  $\hat{v}^{\text{SDP}}$  studied in Section 3 has essentially the best possible rate of convergence among computable estimators.

---

<sup>7</sup>In this section, terms from computational complexity theory defined Section 2 of the online supplementary material [Wang, Berthet and Samworth (2015)] are written in italics at their first occurrence.

4.1. *The planted clique problem.* A graph  $G := (V(G), E(G))$  is an ordered pair in which  $V(G)$  is a countable set, and  $E(G)$  is a subset of  $\{\{x, y\} : x, y \in V(G), x \neq y\}$ . For  $x, y \in V(G)$ , we say  $x$  and  $y$  are *adjacent*, and write  $x \sim y$ , if  $\{x, y\} \in E(G)$ . A *clique*  $C$  is a subset of  $V(G)$  such that  $\{x, y\} \in E(G)$  for all distinct  $x, y \in C$ . The problem of finding a clique of maximum size in a given graph  $G$  is known to be *NP-complete* [Karp (1972)]. It is therefore natural to consider randomly generated input graphs with a clique “planted” in, where the signal is much less confounded by the noise. Such problems were first suggested by Jerrum (1992) and Kučera (1995) as a potentially easier variant of the classical clique problem.

Let  $\mathbb{G}_m$  denote the collection of all graphs with  $m$  vertices. Define  $\mathcal{G}_m$  to be the distribution on  $\mathbb{G}_m$  associated with the standard Erdős–Rényi random graph. In other words, under  $\mathcal{G}_m$ , each pair of vertices is adjacent independently with probability  $1/2$ . For any  $\kappa \in \{1, \dots, m\}$ , let  $\mathcal{G}_{m,\kappa}$  be a distribution on  $\mathbb{G}_m$  constructed by first picking  $\kappa$  distinct vertices uniformly at random and connecting all edges (the “planted clique”), then joining each remaining pair of distinct vertices by an edge independently with probability  $1/2$ . The planted clique problem has input graphs randomly sampled from the distribution  $\mathcal{G}_{m,\kappa}$ . Due to the random nature of the problem, the goal of the planted clique problem is to find (possibly randomised) algorithms that can locate a maximum clique  $K_m$  with high probability.

It is well known that, for a standard Erdős–Rényi graph,  $\frac{|K_m|}{2 \log_2 m} \xrightarrow{\text{a.s.}} 1$  [e.g., Grimmett and McDiarmid (1975)]. In fact, if  $\kappa = \kappa_m$  is such that

$$\liminf_{m \rightarrow \infty} \frac{\kappa}{2 \log_2 m} > 1,$$

it can be shown that the planted clique is asymptotically almost surely also the unique maximum clique in the input graph. As observed in Kučera (1995), there exists  $C > 0$  such that, if  $\kappa > C \sqrt{m \log m}$ , then asymptotically almost surely, vertices in the planted clique have larger degrees than all other vertices, in which case they can be located in  $O(m^2)$  operations. Alon, Krivelevich and Sudakov (1998) improved the above result by exhibiting a spectral method that, given any  $c > 0$ , identifies planted cliques of size  $\kappa \geq c \sqrt{m}$  asymptotically almost surely.

Although several other polynomial time algorithms have subsequently been discovered for the  $\kappa \geq c \sqrt{m}$  case [e.g., Ames and Vavasis (2011), Feige and Krauthgamer (2000), Feige and Ron (2010)], there is no known randomised polynomial time algorithm that can detect below this threshold. Jerrum (1992) hinted at the hardness of this problem by showing that a specific Markov chain approach fails to work when  $\kappa = O(m^{1/2-\delta})$  for some  $\delta > 0$ . Feige and Krauthgamer (2003) showed that Lovàcz–Schrijver semidefinite programming relaxation methods also fail in this regime. Feldman et al. (2013) recently presented further evidence of the hardness of this problem by showing that a broad class of algorithms, which they refer to as “statistical algorithms”, cannot solve the planted clique problem with  $\kappa = O(m^{1/2-\delta})$  in randomised polynomial time, for any  $\delta > 0$ . It is now widely

accepted in theoretical computer science that the planted clique problem is hard, in the sense that the following assumption holds with  $\tau = 0$ :

(A1)( $\tau$ ) For any sequence  $\kappa = \kappa_m$  such that  $\kappa \leq m^\beta$  for some  $0 < \beta < 1/2 - \tau$ , there is no randomised polynomial time algorithm that can correctly identify the planted clique with probability tending to 1 as  $m \rightarrow \infty$ .

We state the assumption in terms of a general parameter  $\tau \in [0, 1/2)$ , because it will turn out below that even if only (A1)( $\tau$ ) holds for some  $\tau \in (0, 1/6)$ , there are still regimes of  $(n, p, k, \theta)$  in which no randomised polynomial time algorithm can attain the minimax optimal rate.

Researchers have used the hardness of the planted clique problem as an assumption to prove various impossibility results in other problems. Examples include cryptographic applications [Applebaum, Barak and Wigderson (2010), Juels and Peinado (2000)], testing  $k$ -wise independence [Alon et al. (2007)] and approximating Nash equilibria [Hazan and Krauthgamer (2011)]. Recent works by Berthet and Rigollet (2013a, 2013b) and Ma and Wu (2015) used a stronger hypothesis on the hardness of detecting the presence of a planted clique to establish computational lower bounds in sparse principal component detection and sparse submatrix detection problems, respectively. Our assumption (A1)(0) assumes only the computational intractability of identifying the entire planted clique, so in particular, is implied by hypothesis  $A_{PC}$  of Berthet and Rigollet (2013b) and Hypothesis 1 of Ma and Wu (2015).

4.2. *Computational lower bounds.* In this section, we use a reduction argument to show that, under assumption (A1)( $\tau$ ), it is impossible to achieve the statistically optimal rate of sparse principal component estimation using randomised polynomial time algorithms. For  $\rho \in \mathbb{N}$ , and for  $x \in \mathbb{R}$ , we let  $[x]_\rho$  denote  $x$  in its binary representation, rounded to  $\rho$  significant figures. Let  $[\mathbb{R}]_\rho := \{[x]_\rho : x \in \mathbb{R}\}$ . We say  $(\hat{v}^{(n)})$  is a *sequence of randomised polynomial time estimators* of  $v_1 \in \mathbb{R}^{p_n}$  if  $\hat{v}^{(n)}$  is a measurable function from  $\mathbb{R}^{n \times p_n}$  to  $\mathbb{R}^{p_n}$  and if, for every  $\rho \in \mathbb{N}$ , there exists a randomised polynomial time algorithm  $M_{pr}$  such that for any  $\mathbf{x} \in ([\mathbb{R}]_\rho)^{n \times p_n}$  we have  $[\hat{v}^{(n)}(\mathbf{x})]_\rho = [M_{pr}(\mathbf{x})]_\rho$ . The sequence of semidefinite programming estimators  $(\hat{v}^{SDP})$  defined in Section 3 is an example of a sequence of randomised polynomial time estimators of  $v_1(P)$ .

**THEOREM 6.** Fix  $\tau \in [0, 1/6)$ , assume (A1)( $\tau$ ), and let  $\alpha \in (0, \frac{1-6\tau}{1-2\tau})$ . For any  $n \in \mathbb{N}$ , let  $(p, k, \theta) = (p_n, k_n, \theta_n)$  be parameters indexed by  $n$  such that  $k = O(p^{1/2-\tau-\delta})$  for some  $\delta \in (0, 1/2 - \tau)$ ,  $n = o(p \log p)$  and  $\theta \leq k^2/(1000p)$ . Suppose further that

$$\frac{k^{1+\alpha} \log p}{n\theta^2} \rightarrow 0$$

as  $n \rightarrow \infty$ . Let  $\mathbf{X}$  be an  $n \times p$  matrix with independent rows, each having distribution  $P$ . Then every sequence  $(\hat{v}^{(n)})$  of randomised polynomial time estimators of  $v_1(P)$  satisfies

$$\sqrt{\frac{n\theta^2}{k^{1+\alpha} \log p}} \sup_{P \in \mathcal{P}_p(n,k,\theta)} \mathbb{E}_P L(\hat{v}^{(n)}(\mathbf{X}), v_1(P)) \rightarrow \infty$$

as  $n \rightarrow \infty$ .

We note that the choices of parameters in the theorem imply that

$$(7) \quad \liminf_{n \rightarrow \infty} \frac{k^2 \log p}{n\theta^2} \geq \liminf_{n \rightarrow \infty} \frac{p}{k^2} = \infty.$$

As remarked in Section 4.1 above, the main interest in this theorem comes from the case  $\tau = 0$ . Here, our result reveals not only that no randomised polynomial time algorithm can attain the minimax optimal rate, but also that in the effective sample size regime described by (7), and provided the other side conditions of Theorem 6 hold, there is in general no consistent sequence of randomised polynomial time estimators. This is in contrast to Theorem 2, where we saw that consistent estimation with a computationally inefficient procedure is possible in the asymptotic regime (7). A further consequence of Theorem 6 is that, since any sequence  $(p, k, \theta) = (p_n, k_n, \theta_n)$  satisfying the conditions of Theorem 6 also satisfies the conditions of Theorem 5 for large  $n$ , the conclusion of Theorem 5 cannot be improved in terms of the exponent of  $k$  (at least, not uniformly over the parameter range given there). As mentioned in the Introduction, for a sufficiently large effective sample size, where even randomised polynomial time estimators can be consistent, the statistical and computational trade-off revealed by Theorems 2 and 6 may disappear. See Section 4.4 below for further details, and Gao, Ma and Zhou (2014) for recent extensions of these results to different classes of distributions.

Even though assumption (A1)(0) is widely believed, we also present results under the weaker family of conditions (A1)( $\tau$ ) for  $\tau \in (0, 1/6)$  to show that a statistical and computational trade-off still remains for certain parameter regimes even in these settings. The reason for assuming  $\tau < 1/6$  is to guarantee that there is a regime of parameters  $(n, p, k, \theta)$  satisfying the conditions of the theorem. Indeed, if  $\tau \in [0, 1/6)$  and  $\alpha \in (0, \frac{1-6\tau}{1-2\tau})$ , we can set  $p = n$ ,  $k = n^{1/2-\tau-\delta}$  for some  $\delta \in (0, \frac{1}{2} - \tau - \frac{1}{3-\alpha})$ ,  $\theta = k^2/(1000n)$ , and in that case,

$$\frac{k^{1+\alpha} \log p}{n\theta^2} = \frac{10^6 n \log n}{k^{3-\alpha}} \rightarrow 0,$$

as required.

4.3. *Sketch of the proof of Theorem 6.* The proof of Theorem 6 relies on a randomised polynomial time reduction from the planted clique problem to the sparse principal component estimation problem. The reduction is adapted from the “bottom-left transformation” of Berthet and Rigollet (2013b), and requires a rather different and delicate analysis.

In greater detail, suppose for a contradiction that we were given a randomised polynomial time algorithm  $\hat{v}$  for the sparse PCA problem with a rate  $\sup_{P \in \mathcal{P}_p(n, k, \theta)} \mathbb{E}_P L(\hat{v}, v_1) \leq \sqrt{\frac{k^{1+\alpha} \log p}{n\theta^2}}$  for some  $\alpha < 1$ . Set  $m \approx p \log p$  and  $\kappa \approx k \log p$ , so we are in the regime where (A1)( $\tau$ ) holds. Given any graph  $G \sim \mathcal{G}_{m, \kappa}$  with planted clique  $K \subseteq V(G)$ , we draw  $n + p$  vertices  $u_1, \dots, u_n, w_1, \dots, w_p$  uniformly at random without replacement from  $V(G)$ . On average there are about  $\kappa / \log \kappa$  clique vertices in  $\{w_1, \dots, w_p\}$ , and our initial aim is to identify a large fraction of these vertices. To do this, we form an  $n \times p$  matrix  $\mathbf{A} := (\mathbb{1}_{u_i \sim w_j})_{i, j}$ , which is an off-diagonal block of the adjacency matrix of  $G$ . We then replace each 0 in  $\mathbf{A}$  with  $-1$  and flip the signs of each row independently with probability  $1/2$  to obtain a new matrix  $\mathbf{X}$ . Each component of the  $i$ th row of  $\mathbf{X}$  has a marginal Rademacher distribution, but if  $u_i$  is a clique vertex, then the components  $\{j : w_j \in K\}$  are perfectly correlated. Writing  $\boldsymbol{\gamma}' := (\mathbb{1}_{\{w_j \in K\}})_{j=1, \dots, p}$ , the leading eigenvector of  $\mathbb{E}\{\mathbf{X}^\top \mathbf{X} / n | \boldsymbol{\gamma}'\}$  is proportional to  $\boldsymbol{\gamma}'$ , which suggests that a spectral method might be able to find  $\{w_1, \dots, w_p\} \cap K$  with high probability. Unfortunately, the joint distribution of the rows of  $\mathbf{X}$  is difficult to deal with directly, but since  $n$  and  $p$  are small relative to  $m$ , we can approximate  $\boldsymbol{\gamma}'$  by a random vector  $\boldsymbol{\gamma}$  having independent  $\text{Bern}(\kappa/m)$  components. We can then approximate  $\mathbf{X}$  by a matrix  $\mathbf{Y}$ , whose rows are independent conditional on  $\boldsymbol{\gamma}$  and have the same marginal distribution conditional on  $\boldsymbol{\gamma} = g$  as the rows of  $\mathbf{X}$  conditional on  $\boldsymbol{\gamma}' = g$ .

It turns out that the distribution of an appropriately scaled version of an arbitrary row of  $\mathbf{Y}$ , conditional on  $\boldsymbol{\gamma} = g$ , belongs to  $\mathcal{P}_p(n, k, \theta)$  for  $g$  belonging to a set of high probability. We could therefore apply our hypothetical randomised polynomial time sparse PCA algorithm to the scaled version of the matrix  $\mathbf{Y}$  to find a good estimate of  $\boldsymbol{\gamma}$ , and since  $\boldsymbol{\gamma}$  is close to  $\boldsymbol{\gamma}'$ , this accomplishes our initial goal. With high probability, the remaining vertices in the planted clique are those having high connectivity to the identified clique vertices in  $\{w_1, \dots, w_p\}$ , which contradicts the hypothesis (A1)( $\tau$ ).

4.4. *Computationally efficient optimal estimation on subparameter spaces in the high effective sample size regime.* Theorems 2, 3, 5 and 6 enable us to summarise, in Table 1 below, our knowledge of the best possible rate of estimation in different asymptotic regimes, both for arbitrary statistical procedures and for those that are computable in randomised polynomial time. (For ease of exposition, we omit here the additional, relatively mild, side constraints required for the above theorems to hold.) The fact that Theorem 6 is primarily concerned with the setting in which  $\frac{k^2 \log p}{n\theta^2} \rightarrow \infty$  raises the question of whether computationally efficient



TABLE 1  
Rate of convergence of best estimator in different asymptotic regimes

	$n \ll \frac{k \log p}{\theta^2}$	$\frac{k \log p}{\theta^2} \ll n \ll \frac{k^2 \log p}{\theta^2}$	$n \gg \frac{k^2 \log p}{\theta^2}$
All estimators	$\asymp 1$	$\asymp \sqrt{\frac{k \log p}{n \theta^2}}$	$\asymp \sqrt{\frac{k \log p}{n \theta^2}}$
Polynomial time estimators	$\asymp 1$	$\asymp 1$	$\lesssim \sqrt{\frac{k^2 \log p}{n \theta^2}}$

procedures could attain a faster rate of convergence in the high effective sample size regime where  $n \gg \frac{k^2 \log p}{\theta^2}$ .

The purpose of this section is to extend the ideas of [Amini and Wainwright \(2009\)](#) to show that, indeed, a variant of the estimator  $\hat{v}^{\text{SDP}}$  introduced in Section 3 attains the minimax optimal rate of convergence in this asymptotic regime, at least over a subclass of the distributions in  $\mathcal{P}_p(n, k, \theta)$ . [Ma \(2013\)](#) and [Yuan and Zhang \(2013\)](#) show similar results for an iterative thresholding algorithm for other subclasses of  $\mathcal{P}_p(n, k, \theta)$  under an extra upper bound condition on  $\lambda_2(P)/\lambda_1(P)$ ; see also [Wang, Lu and Liu \(2014\)](#) and [Deshpande and Montanari \(2014\)](#).

Let  $\mathcal{T}$  denote the set of nonnegative definite matrices  $\Sigma \in \mathbb{R}^{p \times p}$  of the form

$$\Sigma = \theta v_1 v_1^\top + \begin{pmatrix} I_k & 0 \\ 0 & \Gamma_{p-k} \end{pmatrix},$$

where  $v_1 \in \mathbb{R}^p$  is a unit vector such that  $S := \text{supp}(v_1)$  has cardinality  $k$  and where  $\Gamma_{p-k} \in \mathbb{R}^{(p-k) \times (p-k)}$  is nonnegative definite and satisfies  $\lambda_1(\Gamma_{p-k}) \leq 1$ . [Here, and in the proof of Theorem 7 below, the block matrix notation refers to the  $(S, S)$ ,  $(S, S^c)$ ,  $(S^c, S)$  and  $(S^c, S^c)$  blocks.] We now define a subclass of distributions

$$\tilde{\mathcal{P}}_p(n, k, \theta) := \left\{ P \in \mathcal{P}_p(n, k, \theta) : \Sigma(P) \in \mathcal{T}, \min_{j \in S} |v_{1,j}| \geq 16 \sqrt{\frac{k \log p}{n \theta^2}} \right\}.$$

We remark that  $\tilde{\mathcal{P}}_p(n, k, \theta)$  is nonempty only if  $\sqrt{\frac{k^2 \log p}{n \theta^2}} \leq \frac{1}{16}$ , since

$$1 = \|v_{1,S}\|_2 \geq k^{1/2} \min_{j \in S} |v_{1,j}| \geq 16 \sqrt{\frac{k^2 \log p}{n \theta^2}}.$$

This is one reason that the theorem below only holds in the high effective sample size regime. Our variant of  $\hat{v}^{\text{SDP}}$  is described in Algorithm 3 below. We remark that  $\hat{v}^{\text{MSDP}}$ , like  $\hat{v}^{\text{SDP}}$ , is computable in polynomial time.

**THEOREM 7.** Assume that  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} P$  for some  $P \in \tilde{\mathcal{P}}_p(n, k, \theta)$ .

---

**Algorithm 3:** Pseudo-code for computing the modified semidefinite relaxation estimator  $\hat{v}^{\text{MSDP}}$

---

**Input:**  $\mathbf{X} = (X_1, \dots, X_n)^\top \in \mathbb{R}^{n \times p}$ ,  $\lambda > 0$ ,  $\varepsilon > 0$ ,  $\tau > 0$ .

**begin**

**Step 1:** Set  $\hat{\Sigma} \leftarrow n^{-1} \mathbf{X}^\top \mathbf{X}$ .

**Step 2:** For  $f(M) := \text{tr}(\hat{\Sigma} M) - \lambda \|M\|_1$ , let  $\hat{M}^\varepsilon$  be an  $\varepsilon$ -maximiser of  $f$  in  $\mathcal{M}_1$ .

**Step 3:** Let  $\hat{S} \leftarrow \{j \in \{1, \dots, p\} : \hat{M}_{jj}^\varepsilon \geq \tau\}$  and  $\hat{v}^{\text{MSDP}} \in \mathbb{R}^p$  by

$\hat{v}_{\hat{S}^c}^{\text{MSDP}} \leftarrow 0$  and  $\hat{v}_{\hat{S}}^{\text{MSDP}} \in \text{argmax}_{u \in \mathbb{R}^{|\hat{S}|}} u^\top \hat{\Sigma}_{\hat{S}} u$ .

**end**

**Output:**  $\hat{v}^{\text{MSDP}}$

---

(a) Let  $\lambda := 4\sqrt{\frac{\log p}{n}}$ . The function  $f$  in Step 2 of Algorithm 3 has a maximiser  $\hat{M} \in \mathcal{M}_{1,1}(k^2)$  satisfying  $\text{sgn}(\hat{M}) = \text{sgn}(v_1 v_1^\top)$ .

(b) Assume that  $\log p \leq n$ ,  $\theta^2 \leq Bk^{1/2}$  for some  $B \geq 1$  and  $p \geq \theta(n/k)^{1/2}$ . We write  $\hat{v}^{\text{MSDP}}$  for the output of Algorithm 3 with input parameters  $\mathbf{X} := (X_1, \dots, X_n)^\top \in \mathbb{R}^{n \times p}$ ,  $\lambda := 4\sqrt{\frac{\log p}{n}}$ ,  $\varepsilon := (\frac{\log p}{Bn})^{5/2}$  and  $\tau := (\frac{\log p}{Bn})^2$ . Then

$$\sup_{P \in \tilde{\mathcal{P}}_p(n, k, \theta)} \mathbb{E}_P \{L(\hat{v}^{\text{MSDP}}, v_1)\} \leq 6\sqrt{\frac{k \log p}{n\theta^2}}.$$

Theorem 7 generalises Theorem 2 of [Amini and Wainwright \(2009\)](#) in two ways: first, we relax a Gaussianity assumption to an RCC condition; second, the leading eigenvector of the population covariance matrix is not required to have nonzero entries equal to  $\pm k^{-1/2}$ .

**5. Numerical experiments.** In this section, we present the results of numerical experiments to illustrate the results of Theorems 5, 6 and 7. We generate  $v_1 \in \mathbb{R}^p$  by setting  $v_{1,j} := k^{-1/2}$  for  $j = 1, \dots, k$ , and  $v_{1,j} := 0$  for  $j = k + 1, \dots, p$ . We then draw  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} N_p(0, \Sigma)$ , where  $\Sigma := I_p + \theta v_1 v_1^\top$  and  $\theta = 1$ . We apply Algorithm 1 to the data matrix  $\mathbf{X} := (X_1, \dots, X_n)^\top$  and report the average loss of the estimator  $\hat{v}^{\text{SDP}}$  over  $N_{\text{rep}} := 100$  repetitions. For  $p \in \{50, 100, 150, 200\}$  and  $k = \lfloor p^{1/2} \rfloor$ , we repeat the experiment for several choices of  $n$  to explore the three parameter regimes described in Table 1. Since the boundaries of these regimes are  $n \asymp \frac{k \log p}{\theta^2}$  and  $n \asymp \frac{k^2 \log p}{\theta^2}$ , we plot the average loss of the experiments against effective samples sizes

$$v_{\text{lin}} := \frac{n\theta^2}{k \log p} \quad \text{and} \quad v_{\text{quad}} := \frac{n\theta^2}{k^2 \log p}.$$

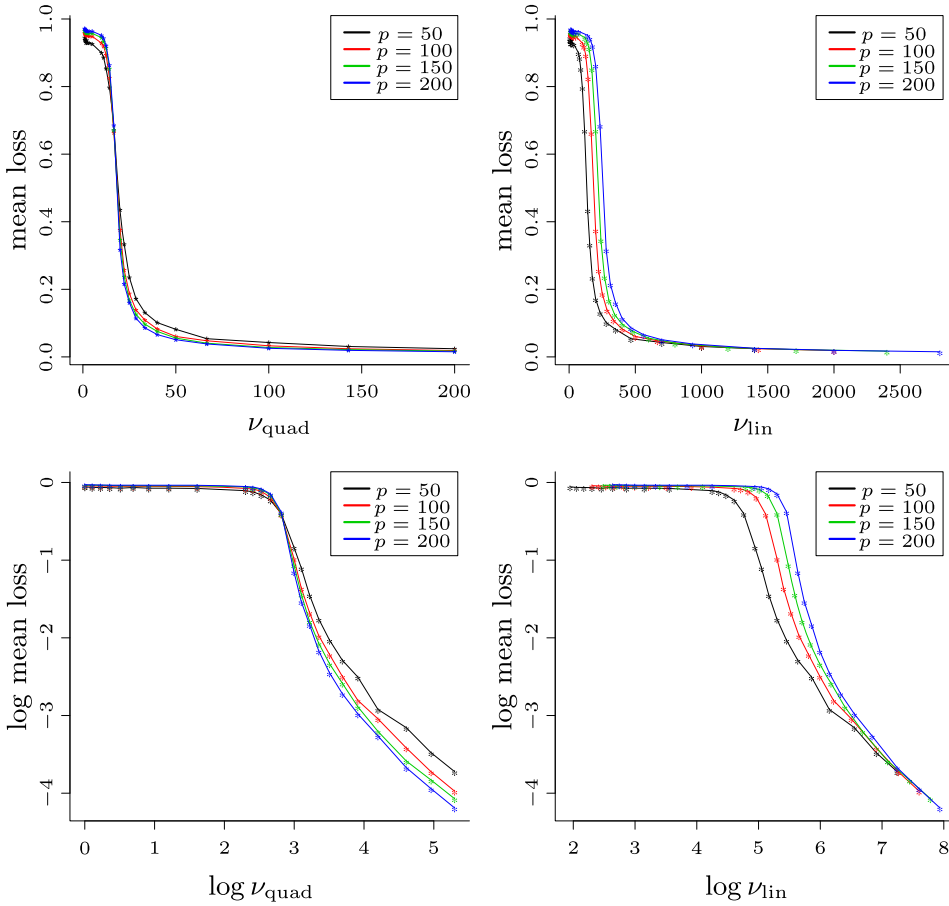


FIG. 1. Average loss of the estimator  $\hat{v}^{\text{SDP}}$  over  $N_{\text{rep}} = 100$  repetitions against effective sample sizes  $\nu_{\text{quad}}$  (top left) and  $\nu_{\text{in}}$  (top right). The tail behaviour under both scalings is examined under logarithmic scales in the bottom left and bottom right panels.

The results are shown in Figure 1. The top left panel of Figure 1 shows a sharp phase transition for the average loss, as predicted by Theorems 5 and 6. The right hand panels of Figure 1 suggest that in the high effective sample size regime,  $\hat{v}^{\text{SDP}}$  converges at rate  $\sqrt{\frac{k \log p}{n \theta^2}}$  in this setting. This is the same rate as was proved for the modified semidefinite relaxation estimator  $\hat{v}^{\text{MSDP}}$  in Theorem 7.

It is worth noting that it is relatively time-consuming to carry out the simulations for the settings in the right-hand tails of the plots in Figure 1. These extreme settings were chosen, however, to illustrate that the linear scaling is the correct one in this tail. For example, when  $\nu_{\text{quad}} = 200$  and  $p = 200$ , we require  $n = 207,694$ , and the pre-processing of the data matrix to obtain the sample covariance matrix is the time-limiting step. In general, in our experience, the semi-definite program-

ming algorithm is certainly not as fast as simpler methods such as diagonal thresholding, but is not prohibitively slow.

APPENDIX A: PROOFS FROM SECTION 2

PROOF OF PROPOSITION 1. (i) Let  $P \in \text{sub-Gaussian}_p(\sigma^2)$ , and assume that  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} P$ . Then, for any  $u \in B_0(\ell)$  and  $t \geq 0$ , we have

$$\mathbb{P}(u^\top X_1 \geq t) \leq e^{-t^2/\sigma^2} \mathbb{E}(e^{tu^\top X_1/\sigma^2}) \leq e^{-t^2/(2\sigma^2)}.$$

Similarly,  $\mathbb{P}(-u^\top X_1 \geq t) \leq e^{-t^2/(2\sigma^2)}$ . Write  $\mu_u := \mathbb{E}\{(u^\top X_1)^2\}$ ; since

$$1 + \frac{1}{2}\mu_u t^2 + o(t^2) = \mathbb{E}(e^{tu^\top X_1}) \leq e^{t^2\sigma^2/2} = 1 + \frac{1}{2}\sigma^2 t^2 + o(t^2),$$

as  $t \rightarrow 0$ , we deduce that  $\mu_u \leq \sigma^2$ . Now, for any integer  $m \geq 2$ ,

$$\begin{aligned} & \mathbb{E}(|(u^\top X_1)^2 - \mu_u|^m) \\ & \leq \int_0^\infty \mathbb{P}\{(u^\top X_1)^2 - \mu_u \geq t^{1/m}\} dt + \mu_u^m \\ & \leq 2 \int_0^\infty e^{-(t^{1/m} + \mu_u)/(2\sigma^2)} dt + \mu_u^m \\ & = m!(2\sigma^2)^m \left\{ 2e^{-\mu_u/(2\sigma^2)} + \frac{1}{m!} \left( \frac{\mu_u}{2\sigma^2} \right)^m \right\} \\ & \leq 2m!(2\sigma^2)^m, \end{aligned}$$

where the final inequality follows because the function  $x \mapsto 2e^{-x} + x^m/m!$  is decreasing on  $[0, 1/2]$ . This calculation allows us to apply Bernstein’s inequality [e.g., van de Geer (2000), Lemma 5.7, taking  $K = 2\sigma^2$ ,  $R = 4\sigma^2$  in her notation], to deduce that for any  $s \geq 0$ ,

$$\mathbb{P}(|\hat{V}(u) - V(u)| \geq s) \leq 2 \exp\left(-\frac{ns^2}{4\sigma^2 s + 32\sigma^4}\right).$$

It follows by Lemma 2 in Section 1 in the supplementary material [Wang, Berthet and Samworth (2015)], taking  $\varepsilon = 1/4$  in that result, that if  $\eta > 0$  is such that  $\ell \log(p/\eta) \leq n$ , then for  $C := 8\sigma^2$ , we have

$$\begin{aligned} & \mathbb{P}\left(\sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| \geq 2C \sqrt{\frac{\ell \log(p/\eta)}{n}}\right) \\ & \leq 2\pi \ell^{1/2} \binom{p}{\ell} \left(\frac{128}{\sqrt{255}}\right)^{\ell-1} \exp\left(-\frac{C^2 \ell \log(p/\eta)}{4C\sigma^2 \sqrt{(\ell \log(p/\eta))/n} + 32\sigma^4}\right) \\ & \leq 2\pi \ell^{1/2} \left(\frac{e}{\ell}\right)^\ell \left(\frac{128}{\sqrt{255}}\right)^{\ell-1} \eta^\ell \leq e^9 \eta. \end{aligned}$$

Similarly, if  $\ell \log(p/\eta) > n$ , then

$$\begin{aligned} & \mathbb{P}\left(\sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| \geq 2C \frac{\ell \log(p/\eta)}{n}\right) \\ & \leq 2\pi \ell^{1/2} \binom{p}{\ell} \left(\frac{128}{\sqrt{255}}\right)^{\ell-1} \exp\left(-\frac{C^2 \ell^2 \log^2(p/\eta)}{4C\sigma^2 \ell \log(p/\eta) + 32\sigma^4 n}\right) \leq e^9 \eta. \end{aligned}$$

Setting  $\delta := e^9 \eta$ , we find (noting that we only need to consider the case  $\delta \in (0, 1]$ ) that

$$\begin{aligned} & \mathbb{P}\left\{\sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| \geq 16\sigma^2 \left(1 + \frac{9}{\log p}\right) \max\left(\sqrt{\frac{\ell \log(p/\delta)}{n}}, \frac{\ell \log(p/\delta)}{n}\right)\right\} \\ & \leq \mathbb{P}\left\{\sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| \geq 16\sigma^2 \max\left(\sqrt{\frac{\ell \log(e^9 p/\delta)}{n}}, \frac{\ell \log(e^9 p/\delta)}{n}\right)\right\} \\ & \leq \delta. \end{aligned}$$

(ii) By Lemma 1 of [Laurent and Massart \(2000\)](#), if  $Y_1, \dots, Y_n$  are independent  $\chi_1^2$  random variables, then for all  $a > 0$ ,

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n Y_i - n\right| \geq a\right) \leq 2e^{-(n/2)(1+a-\sqrt{1+2a})} \leq 2e^{-n \min(a/4, a^2/16)}.$$

Setting  $\eta := e^{-n \min(a/4, a^2/16)}$ , we deduce that

$$\mathbb{P}\left\{\left|\frac{1}{n} \sum_{i=1}^n Y_i - n\right| \geq 4 \max\left(\sqrt{\frac{\log(1/\eta)}{n}}, \frac{\log(1/\eta)}{n}\right)\right\} \leq 2\eta.$$

Hence, using Lemma 2 again, and by a similar calculation to part (i),

$$\mathbb{P}\left\{\sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| \geq 8\lambda_1(P) \max\left(\sqrt{\frac{\log(1/\eta)}{n}}, \frac{\log(1/\eta)}{n}\right)\right\} \leq e^9 p^\ell \eta.$$

The result follows on setting  $\delta := e^9 p^\ell \eta$ .  $\square$

**PROOF OF THEOREM 2.** Fix an arbitrary  $P \in \mathcal{P}_p(n, k, \theta)$ . For notational simplicity, we write  $v := v_1(P)$  and  $\hat{v} := \hat{v}_{\max}^k(\hat{\Sigma})$  in this proof. We now exploit the curvature lemma of [Vu et al. \[\(2013\), Lemma 3.1\]](#), which is closely related to the Davis–Kahan  $\sin \theta$  theorem [[Davis and Kahan \(1970\)](#), [Yu, Wang and Samworth \(2015\)](#)]. This lemma gives that

$$\|\hat{v}\hat{v}^\top - vv^\top\|_2^2 \leq \frac{2}{\theta} \text{tr}(\Sigma(vv^\top - \hat{v}\hat{v}^\top)) \leq \frac{2}{\theta} \text{tr}((\Sigma - \hat{\Sigma})(vv^\top - \hat{v}\hat{v}^\top)).$$

When  $\hat{v}\hat{v}^\top \neq vv^\top$ , we have that  $\frac{vv^\top - \hat{v}\hat{v}^\top}{\|vv^\top - \hat{v}\hat{v}^\top\|_2}$  has rank 2, trace 0 and has nonzero entries in at most  $2k$  rows and  $2k$  columns. It follows that its nonzero eigenvalues are  $\pm 1/\sqrt{2}$ , so it can be written as  $(xx^\top - yy^\top)/\sqrt{2}$  for some  $x, y \in B_0(2k)$ . Thus,

$$\begin{aligned} \mathbb{E}L(\hat{v}, v) &= \mathbb{E}\frac{1}{\sqrt{2}}\|\hat{v}\hat{v}^\top - vv^\top\|_2 \leq \frac{1}{\theta}\mathbb{E}\text{tr}((\Sigma - \hat{\Sigma})(xx^\top - yy^\top)) \\ &\leq \frac{2}{\theta}\mathbb{E}\sup_{u \in B_0(2k)}|\hat{V}(u) - V(u)| \leq 2\sqrt{2}\left(1 + \frac{1}{\log p}\right)\sqrt{\frac{k \log p}{n\theta^2}}, \end{aligned}$$

where we have used Proposition 1 in Section 1 in the online supplementary material [Wang, Berthet and Samworth (2015)] to obtain the final inequality.  $\square$

PROOF OF THEOREM 3. Set  $\sigma^2 := \frac{1}{8(1+9/\log p)} - \theta$ . We have by Proposition 1(ii) that  $N_p(0, \sigma^2 I_p + \theta v_1 v_1^\top) \in \mathcal{P}_p(n, k, \theta)$  for any unit vector  $v_1 \in B_0(k)$ . Define  $k_0 := k - 1$  and  $p_0 := p - 1$ . Applying the variant of the Gilbert–Varshamov lemma given as Lemma 3 in Section 1 in the online supplementary material [Wang, Berthet and Samworth (2015)] with  $\alpha := 1/2$  and  $\beta := 1/4$ , we can construct a set  $\mathcal{N}_0$  of  $k_0$ -sparse vectors in  $\{0, 1\}^{p_0}$  with cardinality at least  $(p_0/k_0)^{k_0/8}$ , such that the Hamming distance between every pair of distinct points in  $\mathcal{N}_0$  is at least  $k_0$ . For  $\varepsilon \in (0, 1]$  to be chosen later, define a set of  $k$ -sparse vectors in  $\mathbb{R}^p$  by

$$\mathcal{N} := \left\{ \begin{pmatrix} \sqrt{1 - \varepsilon^2} \\ k_0^{-1/2} \varepsilon u_0 \end{pmatrix} : u_0 \in \mathcal{N}_0 \right\}.$$

Observe that if  $u, v$  are distinct elements of  $\mathcal{N}$ , then

$$L(u, v) = \{1 - (u^\top v)^2\}^{1/2} \geq \{1 - (1 - \varepsilon^2/2)^2\}^{1/2} \geq \frac{\sqrt{3}\varepsilon}{2},$$

and similarly  $L(u, v) \leq \varepsilon$ . For  $u \in \mathcal{N}$ , let  $P_u$  denote the multivariate normal distribution  $N_p(0, \sigma^2 I_p + \theta uu^\top)$ . For any estimator  $\hat{v} \in \mathcal{V}_{n,p}$ , we define  $\hat{\psi}_{\hat{v}} := \text{sargmin}_{u \in \mathcal{N}} L(\hat{v}, u)$ , where  $\text{sargmin}$  denotes the smallest element of the argmin in the lexicographic ordering. Note that  $\{\hat{\psi}_{\hat{v}} \neq u\} \subseteq \{L(\hat{v}, u) \geq \sqrt{3}\varepsilon/4\}$ . We now apply the generalised version of Fano’s lemma given as Lemma 4 in Section 1 in the online supplementary material [Wang, Berthet and Samworth (2015)]. Writing  $D(P\|Q)$  for the Kullback–Leibler divergence between two probability measures defined on the same space (a formal definition is given just prior to Lemma 4), we have

$$\begin{aligned} &\inf_{\hat{v} \in \mathcal{V}_{n,p}} \sup_{P \in \mathcal{P}_p(n,k,\theta)} \mathbb{E}_P L(\hat{v}, v_1(P)) \\ (8) \quad &\geq \inf_{\hat{v} \in \mathcal{V}_{n,p}} \max_{u \in \mathcal{N}} \mathbb{E}_{P_u} L(\hat{v}, u) \geq \frac{\sqrt{3}\varepsilon}{4} \inf_{\hat{v} \in \mathcal{V}_{n,p}} \max_{u \in \mathcal{N}} P_u^{\otimes n}(\hat{\psi}_{\hat{v}} \neq u) \\ &\geq \frac{\sqrt{3}\varepsilon}{4} \left( 1 - \frac{\max_{u,v \in \mathcal{N}, u \neq v} D(P_v^{\otimes n} \| P_u^{\otimes n}) + \log 2}{(k_0/8) \log(p_0/k_0)} \right). \end{aligned}$$

We can compute, for distinct points  $u, v \in \mathcal{N}$ ,

$$\begin{aligned}
 D(P_v^{\otimes n} \| P_u^{\otimes n}) &= nD(P_v \| P_u) = \frac{n}{2} \operatorname{tr}((\sigma^2 I_p + \theta uu^\top)^{-1}(\sigma^2 I_p + \theta vv^\top) - I_p) \\
 (9) \quad &= \frac{n}{2} \operatorname{tr}((\sigma^2 I_p + \theta uu^\top)^{-1} \theta (vv^\top - uu^\top)) \\
 &= \frac{n\theta}{2} \operatorname{tr}\left(\left(\frac{1}{\sigma^2} I_p - \frac{\theta}{\sigma^2(\sigma^2 + \theta)} uu^\top\right)(vv^\top - uu^\top)\right) \\
 &= \frac{n\theta^2}{2\sigma^2(\sigma^2 + \theta)} L^2(u, v) \leq \frac{n\theta^2 \varepsilon^2}{2\sigma^2(\sigma^2 + \theta)}.
 \end{aligned}$$

Let  $\varepsilon := \min\{\sqrt{a/(3b)}, 1\}$ , where

$$a := 1 - \frac{8 \log 2}{k_0 \log(p_0/k_0)} \quad \text{and} \quad b := \frac{4n\theta^2}{\sigma^2(\sigma^2 + \theta)k_0 \log(p_0/k_0)}.$$

Then from (8) and (9), we find that

$$\inf_{\hat{v} \in \mathcal{V}_{n,p}} \sup_{P \in \mathcal{P}_p(n,k,\theta)} \mathbb{E}_P L(\hat{v}, v_1(P)) \geq \min\left\{\frac{1}{1660} \sqrt{\frac{k \log p}{n\theta^2}}, \frac{5}{18\sqrt{3}}\right\},$$

as required.  $\square$

## APPENDIX B: PROOFS FROM SECTION 3

**PROOF OF LEMMA 4.** For convenience, we write  $v := v_1(\Sigma)$ ,  $\hat{v}$  for  $\hat{v}^{\text{SDP}}$  and  $\hat{M}$  for  $\hat{M}^\varepsilon$  in this proof. We first study  $vv^\top - \hat{M}$ , where  $\hat{M} \in \mathcal{M}_1$  is computed in Step 2 of Algorithm 1. By the curvature lemma of Vu et al. [(2013), Lemma 3.1],

$$\|vv^\top - \hat{M}\|_2^2 \leq \frac{2}{\theta} \operatorname{tr}(\Sigma(vv^\top - \hat{M})).$$

Moreover, since  $vv^\top \in \mathcal{M}_1$ , we have the basic inequality

$$\operatorname{tr}(\hat{\Sigma} \hat{M}) - \lambda \|\hat{M}\|_1 \geq \operatorname{tr}(\hat{\Sigma} vv^\top) - \lambda \|vv^\top\|_1 - \varepsilon.$$

Let  $S$  denote the set of indices corresponding to the nonzero components of  $v$ , and recall that  $|S| \leq k$ . Since by hypothesis  $\|\hat{\Sigma} - \Sigma\|_\infty \leq \lambda$ , we have

$$\begin{aligned}
 \|vv^\top - \hat{M}\|_2^2 &\leq \frac{2}{\theta} \{\operatorname{tr}(\hat{\Sigma}(vv^\top - \hat{M})) + \operatorname{tr}((\Sigma - \hat{\Sigma})(vv^\top - \hat{M}))\} \\
 &\leq \frac{2}{\theta} (\lambda \|vv^\top\|_1 - \lambda \|\hat{M}\|_1 + \varepsilon + \|\hat{\Sigma} - \Sigma\|_\infty \|vv^\top - \hat{M}\|_1) \\
 &\leq \frac{2\lambda}{\theta} (\|v_S v_S^\top\|_1 - \|\hat{M}_{S,S}\|_1 + \|v_S v_S^\top - \hat{M}_{S,S}\|_1) + \frac{2\varepsilon}{\theta} \\
 &\leq \frac{4\lambda}{\theta} \|v_S v_S^\top - \hat{M}_{S,S}\|_1 + \frac{2\varepsilon}{\theta} \leq \frac{4\lambda k}{\theta} \|vv^\top - \hat{M}\|_2 + \frac{2\varepsilon}{\theta}.
 \end{aligned}$$



We deduce that

$$\|vv^\top - \hat{M}\|_2 \leq \frac{4\lambda k}{\theta} + \sqrt{\frac{2\varepsilon}{\theta}}.$$

On the other hand,

$$\begin{aligned} \|vv^\top - \hat{M}\|_2^2 &= \text{tr}((vv^\top - \hat{M})^2) = 1 - 2v^\top \hat{M}v + \text{tr}(\hat{M}^2) \\ &\geq 1 - 2\hat{v}^\top \hat{M}\hat{v} + \text{tr}(\hat{M}^2) = \|\hat{v}\hat{v}^\top - \hat{M}\|_2^2. \end{aligned}$$

We conclude that

$$\begin{aligned} L(\hat{v}, v) &= \frac{1}{\sqrt{2}} \|\hat{v}\hat{v}^\top - vv^\top\|_2 \leq \frac{1}{\sqrt{2}} (\|\hat{v}\hat{v}^\top - \hat{M}\|_2 + \|vv^\top - \hat{M}\|_2) \\ &\leq \sqrt{2} \|vv^\top - \hat{M}\|_2 \leq \frac{4\sqrt{2}\lambda k}{\theta} + 2\sqrt{\frac{\varepsilon}{\theta}}, \end{aligned}$$

as required.  $\square$

**PROOF OF THEOREM 5.** Fix  $P \in \mathcal{P}_p(n, k, \theta)$ . By Lemma 4, and by Lemma 5 in Section 1 of the online supplementary material [Wang, Berthet and Samworth (2015)],

$$\begin{aligned} \mathbb{E}L(\hat{v}^{\text{SDP}}, v_1(P)) &= \mathbb{E}\{L(\hat{v}^{\text{SDP}}, v_1(P))\mathbb{1}_{\{\|\hat{\Sigma} - \Sigma\|_\infty \leq \lambda\}}\} \\ (10) \quad &+ \mathbb{E}\{L(\hat{v}^{\text{SDP}}, v_1(P))\mathbb{1}_{\{\|\hat{\Sigma} - \Sigma\|_\infty > \lambda\}}\} \\ &\leq \frac{4\sqrt{2}\lambda k}{\theta} + 2\sqrt{\frac{\varepsilon}{\theta}} + \mathbb{P}\left(\sup_{u \in B_0(2)} |\hat{V}(u) - V(u)| > 2\sqrt{\frac{\log p}{n}}\right). \end{aligned}$$

Since  $P \in \text{RCC}_p(n, 2, 1)$ , we have for each  $\delta > 0$  that

$$\mathbb{P}\left\{\sup_{u \in B_0(2)} |\hat{V}(u) - V(u)| > \max\left(\sqrt{\frac{2 \log(p/\delta)}{n}}, \frac{2 \log(p/\delta)}{n}\right)\right\} \leq \delta.$$

Set  $\delta := \sqrt{\frac{k^2 \log p}{n\theta^2}}$ . Since  $4 \log p \leq n$ , which in particular implies  $n \geq 3$ , we have

$$\frac{2 \log(p/\delta)}{n} \leq \frac{1}{2} + \frac{1}{n} \log\left(\frac{n\theta^2}{k^2 \log p}\right) \leq \frac{1}{2} + \frac{\log n}{n} - \frac{1}{n} \log \log 2 \leq 1.$$

Moreover, since  $n \leq k^2 p^2 \theta^{-2} \log p$ ,

$$2 \log(p/\delta) = 2 \log p + \log\left(\frac{n\theta^2}{k^2 \log p}\right) \leq 4 \log p.$$

We deduce that

$$(11) \quad \mathbb{P}\left(\sup_{u \in B_0(2)} |\hat{V}(u) - V(u)| > 2\sqrt{\frac{\log p}{n}}\right) \leq \sqrt{\frac{k^2 \log p}{n\theta^2}}.$$

The desired risk bound follows from (10), the fact that  $\theta \leq k$ , and (11).  $\square$

APPENDIX C: PROOFS FROM SECTION 4

PROOF OF THEOREM 6. Suppose, for a contradiction, that there exist an infinite subset  $\mathcal{N}$  of  $\mathbb{N}$ ,  $K_0 \in [0, \infty)$  and a sequence  $(\hat{v}^{(n)})$  of randomised polynomial time estimators of  $v_1(P)$  satisfying

$$\sup_{P \in \mathcal{P}_p(n, k, \theta)} \mathbb{E}_P L(\hat{v}^{(n)}(\mathbf{X}), v_1(P)) \leq K_0 \sqrt{\frac{k^{1+\alpha} \log p}{n\theta^2}}$$

for all  $n \in \mathcal{N}$ . Let  $L := \lceil \log p_n \rceil$ , let  $m = m_n := \lceil 10Lp_n/9 \rceil$  and let  $\kappa = \kappa_n := Lk_n$ . We claim that Algorithm 4 below is a randomised polynomial time algorithm that correctly identifies the planted clique problem on  $m_n$  vertices and a planted clique of size  $\kappa_n$  with probability tending to 1 as  $n \rightarrow \infty$ . Since  $\kappa_n = O(m_n^{1/2-\tau-\delta} \log m_n)$ , this contradicts assumption (A1)( $\tau$ ). We prove the claim below.

Let  $G \sim \mathbb{G}_{m, \kappa}$ , and let  $K \subseteq V(G)$  denote the planted clique. Note that the matrix  $\mathbf{A}$  defined in Step 1 of Algorithm 4 is the off-diagonal block of the adjacency matrix of  $G$  associated with the bipartite graph induced by the two parts  $\{u_i : i = 1, \dots, n\}$  and  $\{w_j : j = 1, \dots, p\}$ . Let  $\boldsymbol{\varepsilon}' = (\varepsilon'_1, \dots, \varepsilon'_n)^\top$  and  $\boldsymbol{\gamma}' = (\gamma'_1, \dots, \gamma'_p)^\top$ , where  $\varepsilon'_i := \mathbb{1}_{\{u_i \in K\}}$ ,  $\gamma'_j := \mathbb{1}_{\{w_j \in K\}}$ , and set  $S' := \{j : \gamma'_j = 1\}$ .

**Algorithm 4:** Pseudo-code for a planted clique algorithm based on a hypothetical randomised polynomial time sparse principal component estimation algorithm

**Input:**  $m \in \mathbb{N}$ ,  $\kappa \in \{1, \dots, m\}$ ,  $G \in \mathbb{G}_m$ ,  $L \in \mathbb{N}$

**begin**

**Step 1:** Let  $n \leftarrow \lfloor 9m/(10L) \rfloor$ ,  $p \leftarrow p_n$ ,  $k \leftarrow \lfloor \kappa/L \rfloor$ . Draw  $u_1, \dots, u_n$ ,  $w_1, \dots, w_p$  uniformly at random without replacement from  $V(G)$ . Form  $\mathbf{A} = (A_{ij}) \leftarrow (\mathbb{1}_{\{u_i \sim w_j\}}) \in \mathbb{R}^{n \times p}$  and  $\mathbf{X} \leftarrow \text{diag}(\xi_1, \dots, \xi_n)(2\mathbf{A} - \mathbf{1}_{n \times p})$ , where  $\xi_1, \dots, \xi_n$  are independent Rademacher random variables (independent of  $u_1, \dots, u_n, w_1, \dots, w_p$ ), and where every entry of  $\mathbf{1}_{n \times p} \in \mathbb{R}^{n \times p}$  is 1.

**Step 2:** Use the randomised estimator  $\hat{v}^{(n)}$  to compute  $\hat{v} = \hat{v}^{(n)}(\mathbf{X}/\sqrt{750})$ .

**Step 3:** Let  $\hat{S} = \hat{S}(\hat{v})$  be the lexicographically smallest  $k$ -subset of  $\{1, \dots, p\}$  such that  $(\hat{v}_j : j \in \hat{S})$  contains the  $k$  largest coordinates of  $\hat{v}$  in absolute value.

**Step 4:** For  $u \in V(G)$  and  $W \subseteq V(G)$ , let  $\text{nb}(u, W) := \mathbb{1}_{\{u \in W\}} + \sum_{w \in W} \mathbb{1}_{\{u \sim w\}}$ . Set  $\hat{K} := \{u \in V(G) : \text{nb}(u, \{w_j : j \in \hat{S}\}) \geq 3k/4\}$ .

**end**

**Output:**  $\hat{K}$

It is convenient at this point to introduce the notion of a *graph vector distribution*. We say  $Y$  has a  $p$ -variate graph vector distribution with parameters  $g = (g_1, \dots, g_p)^\top \in \{0, 1\}^p$  and  $\pi_0 \in [0, 1]$ , and write  $Y \sim \text{GV}_p^g(\pi_0)$ , if we can write

$$Y = \xi \{(1 - \varepsilon)R + \varepsilon(g + \tilde{R})\},$$

where  $\xi$ ,  $\varepsilon$  and  $R$  are independent, where  $\xi$  is a Rademacher random variable, where  $\varepsilon \sim \text{Bern}(\pi_0)$ , where  $R = (R_1, \dots, R_p)^\top \in \mathbb{R}^p$  has independent Rademacher components, and where  $\tilde{R} = (\tilde{R}_1, \dots, \tilde{R}_p)^\top$  with  $\tilde{R}_j := (1 - g_j)R_j$ .

Let  $(\boldsymbol{\varepsilon}, \boldsymbol{\gamma})^\top = (\varepsilon_1, \dots, \varepsilon_n, \gamma_1, \dots, \gamma_p)^\top$  be  $n + p$  independent  $\text{Bern}(\kappa/m)$  random variables. For  $i = 1, \dots, n$ , let  $Y_i := \xi_i \{(1 - \varepsilon_i)R_i + \varepsilon_i(\boldsymbol{\gamma} + \tilde{R}_i)\}$  so that, conditional on  $\boldsymbol{\gamma}$ , the random vectors  $Y_1, \dots, Y_n$  are independent, each distributed as  $\text{GV}_p^\gamma(\kappa/m)$ . As shorthand, we denote this conditional distribution as  $\mathcal{Q}_\boldsymbol{\gamma}$ , and write  $S := \{j : \gamma_j = 1\}$ . Note that by Lemma 6 in Section 1 of the online supplementary material [Wang, Berthet and Samworth (2015)],  $\mathcal{Q}_\boldsymbol{\gamma} \in \bigcap_{\ell=1}^{\lfloor 20p/(9k) \rfloor} \text{RCC}_p(\ell, 750)$ .

Let  $\mathbf{Y} := (Y_1, \dots, Y_n)^\top$ . Recall that if  $P$  and  $Q$  are probability measures on a measurable space  $(\mathcal{X}, \mathcal{B})$ , the *total variation distance* between  $P$  and  $Q$  is defined by

$$d_{\text{TV}}(P, Q) := \sup_{B \in \mathcal{B}} |P(B) - Q(B)|.$$

Writing  $\mathcal{L}(Z)$  for the distribution (or law) of a generic random element  $Z$ , and using elementary properties of the total variation distance given in Lemma 9 in Section 1 in the online supplementary material [Wang, Berthet and Samworth (2015)], we have

$$\begin{aligned} d_{\text{TV}}(\mathcal{L}(\mathbf{X}), \mathcal{L}(\mathbf{Y})) &\leq d_{\text{TV}}(\mathcal{L}(\boldsymbol{\varepsilon}', \boldsymbol{\gamma}', (R_{ij}), (\xi_i)), \mathcal{L}(\boldsymbol{\varepsilon}, \boldsymbol{\gamma}, (R_{ij}), (\xi_i))) \\ (12) \qquad \qquad \qquad &= d_{\text{TV}}(\mathcal{L}(\boldsymbol{\varepsilon}', \boldsymbol{\gamma}'), \mathcal{L}(\boldsymbol{\varepsilon}, \boldsymbol{\gamma})) \\ &\leq \frac{2(n + p)}{m} \leq \frac{9(n + p)}{5p \log p}. \end{aligned}$$

Here, the penultimate inequality follows from Diaconis and Freedman [(1980), Theorem 4]. In view of (12), we initially analyse Steps 2, 3 and 4 in Algorithm 4 with  $\mathbf{X}$  replaced by  $\mathbf{Y}$ . Observe that  $\mathbb{E}(Y_i | \boldsymbol{\gamma}) = 0$  and, writing  $\Delta := \text{diag}(\boldsymbol{\gamma}) \in \mathbb{R}^{p \times p}$ , we have

$$\begin{aligned} \Sigma_\boldsymbol{\gamma} &:= \text{Cov}(Y_i | \boldsymbol{\gamma}) = \mathbb{E}\{(1 - \varepsilon_i)R_i R_i^\top + \varepsilon_i(\boldsymbol{\gamma} + \tilde{R}_i)(\boldsymbol{\gamma} + \tilde{R}_i)^\top | \boldsymbol{\gamma}\} \\ &= I_p + \frac{\kappa}{m}(\boldsymbol{\gamma}\boldsymbol{\gamma}^\top - \Delta). \end{aligned}$$

Writing  $N_\boldsymbol{\gamma} := \sum_{j=1}^p \gamma_j$ , it follows that the largest eigenvalue of  $\Sigma_\boldsymbol{\gamma}$  is  $1 + \frac{\kappa}{m}(N_\boldsymbol{\gamma} - 1)$ , with corresponding eigenvector  $\boldsymbol{\gamma}/N_\boldsymbol{\gamma}^{1/2} \in B_0(N_\boldsymbol{\gamma})$ . The other eigenvalues

are 1, with multiplicity  $p - N_{\mathcal{Y}}$ , and  $1 - \frac{\kappa}{m}$ , with multiplicity  $N_{\mathcal{Y}} - 1$ . Hence,  $\lambda_1(\Sigma_{\mathcal{Y}}) - \lambda_2(\Sigma_{\mathcal{Y}}) = \frac{\kappa}{m}(N_{\mathcal{Y}} - 1)$ . Define

$$\Gamma_0 := \left\{ g \in \{0, 1\}^p : \left| N_g - \frac{p\kappa}{m} \right| \leq \frac{k}{20} \right\},$$

where  $N_g := \sum_{j=1}^p g_j$ . We note that by Bernstein’s inequality [e.g., [Shorack and Wellner \(1986\)](#), page 855] that

$$(13) \quad \mathbb{P}(\mathcal{Y} \in \Gamma_0) \geq 1 - 2e^{-k/800}.$$

If  $g \in \Gamma_0$ , the conditional distribution of  $Y_1/\sqrt{750}$  given  $\mathcal{Y} = g$  belongs to  $\mathcal{P}_p(n, k, \theta)$  for  $\theta \leq \frac{\kappa}{750m}(N_g - 1)$  and all large  $n \in \mathcal{N}$ . By hypothesis, it follows that for  $g \in \Gamma_0$ ,

$$\mathbb{E}\{L(\hat{v}^{(n)}(\mathbf{Y}/\sqrt{750}), v_1(Q_{\mathcal{Y}})) | \mathcal{Y} = g\} \leq K_0 \sqrt{\frac{k^{1+\alpha} \log p}{n\theta^2}}$$

for all large  $n \in \mathcal{N}$ . Then by Lemma 7 in Section 1 in the online supplementary material [[Wang, Berthet and Samworth \(2015\)](#)], for  $\hat{S}(\cdot)$  defined in Step 3 of Algorithm 4, for  $g \in \Gamma_0$ , and large  $n \in \mathcal{N}$ ,

$$\begin{aligned} \mathbb{E}\{|S \setminus \hat{S}(\hat{v}^{(n)}(\mathbf{Y}/\sqrt{750}))| | \mathcal{Y} = g\} &\leq 2N_g \mathbb{E}\{L(\hat{v}^{(n)}(\mathbf{Y}/\sqrt{750}), v_1(Q_{\mathcal{Y}}))^2 | \mathcal{Y} = g\} \\ &\leq 2N_g K_0 \sqrt{\frac{k^{1+\alpha} \log p}{n\theta^2}}. \end{aligned}$$

We deduce by Markov’s inequality that for  $g \in \Gamma_0$ , and large  $n \in \mathcal{N}$ ,

$$(14) \quad \mathbb{P}\{|S \cap \hat{S}(\hat{v}^{(n)}(\mathbf{Y}/\sqrt{750}))| \leq 16N_{\mathcal{Y}}/17 | \mathcal{Y} = g\} \leq 34K_0 \sqrt{\frac{k^{1+\alpha} \log p}{n\theta^2}}.$$

Let

$$\Omega_{0,n} := \{\mathcal{Y} \in \Gamma_0\} \cap \{|S \cap \hat{S}(\hat{v}^{(n)}(\mathbf{Y}/\sqrt{750}))| > 16N_{\mathcal{Y}}/17\},$$

$$\Omega'_{0,n} := \{\mathcal{Y}' \in \Gamma_0\} \cap \{|S \cap \hat{S}(\hat{v}^{(n)}(\mathbf{X}/\sqrt{750}))| > 16N_{\mathcal{Y}'}/17\} =: \Omega'_{1,n} \cap \Omega'_{2,n},$$

say, where  $N_{\mathcal{Y}'} := \sum_{j=1}^p \gamma'_j$ . When  $n \in \mathcal{N}$  is sufficiently large, we have on the event  $\Omega'_{0,n}$  that

$$(15) \quad |\{j \in \hat{S}(\hat{v}^{(n)}(\mathbf{X}/\sqrt{750})) : w_j \in K\}| > 3k/4.$$

Now set

$$\Omega'_{3,n} := \left\{ \text{nb}(u, \{w_j : j \in S'\}) \leq \frac{k}{2} \text{ for all } u \in V(G) \setminus K \right\}.$$

Recall the definition of  $\hat{K}$  from Step 4 of Algorithm 4. We claim that for sufficiently large  $n \in \mathcal{N}$ ,

$$\Omega'_{0,n} \cap \Omega'_{3,n} \subseteq \{\hat{K} = K\}.$$

To see this, note that for  $n \in \mathcal{N}$  sufficiently large, on  $\Omega'_{0,n}$  we have  $K \subseteq \hat{K}$  by (15). For the reverse inclusion, note that if  $u \in V(G) \setminus K$ , then on  $\Omega'_{0,n} \cap \Omega'_{3,n}$ , we have for sufficiently large  $n \in \mathcal{N}$  that

$$\begin{aligned} & \text{nb}(u, \{w_j : j \in \hat{S}(\hat{v}^{(n)}(\mathbf{X}/\sqrt{750}))\}) \\ & \leq |\{w_j : j \in \hat{S}\} \setminus K| + \text{nb}(u, \{w_j : j \in \hat{S}\} \cap K) \\ & \leq |\{w_j : j \in \hat{S}\} \setminus K| + \text{nb}(u, \{w_j : j \in S'\}) < \frac{k}{4} + \frac{k}{2} = \frac{3k}{4}. \end{aligned}$$

This establishes our claim. We conclude that for sufficiently large  $n \in \mathcal{N}$ ,

$$(16) \quad \mathbb{P}(\hat{K} \neq K) \leq \mathbb{P}((\Omega'_{0,n} \cap \Omega'_{3,n})^c) \leq \mathbb{P}((\Omega'_{0,n})^c) + \mathbb{P}(\Omega'_{1,n} \cap (\Omega'_{3,n})^c).$$

Now by Lemma 9 in Section 1 in the online supplementary material [Wang, Berthet and Samworth (2015)], we have

$$(17) \quad |\mathbb{P}(\Omega'_{0,n}) - \mathbb{P}(\Omega_{0,n})| \leq d_{\text{TV}}(\mathcal{L}(\mathbf{X}, \boldsymbol{\gamma}'), \mathcal{L}(\mathbf{Y}, \boldsymbol{\gamma})) \leq \frac{9(n+p)}{5p \log p}.$$

Moreover, by a union bound and Hoeffding’s inequality, for large  $n \in \mathcal{N}$ ,

$$(18) \quad \mathbb{P}(\Omega'_{1,n} \cap (\Omega'_{3,n})^c) \leq \sum_{g \in \Gamma_0} \mathbb{P}((\Omega'_{3,n})^c | \boldsymbol{\gamma} = g) \mathbb{P}(\boldsymbol{\gamma} = g) \leq m e^{-k/800}.$$

We conclude by (16), (17), (13), (14) and (18) that for large  $n \in \mathcal{N}$ ,

$$\mathbb{P}(\hat{K} \neq K) \leq \frac{9(n+p)}{5p \log p} + 2e^{-k/800} + 34K_0 \sqrt{\frac{k^{1+\alpha} \log p}{n\theta^2}} + m e^{-k/800} \rightarrow 0$$

as  $n \rightarrow \infty$ . This contradicts assumption (A1)( $\tau$ ) and, therefore, completes the proof.  $\square$

**PROOF OF THEOREM 7.** Setting  $\delta := p^{-1}$  in (3), there exist events  $\Omega_1$  and  $\Omega_2$ , each with probability at least  $1 - p^{-1}$ , such that on  $\Omega_1$  and  $\Omega_2$ , we, respectively, have

$$(19) \quad \begin{aligned} \sup_{u \in B_0(2k)} |\hat{V}(u) - V(u)| & \leq 2\sqrt{\frac{k \log p}{n}} \quad \text{and} \\ \sup_{u \in B_0(2)} |\hat{V}(u) - V(u)| & \leq 2\sqrt{\frac{\log p}{n}}. \end{aligned}$$

Let  $\Omega_0 := \Omega_1 \cap \Omega_2$ . We work on  $\Omega_0$  henceforth. The main ingredient for proving both parts of the theorem is the following weak-duality inequality:

$$(20) \quad \begin{aligned} \max_{M \in \mathcal{M}_1} \text{tr}(\hat{\Sigma}M) - \lambda \|M\|_1 & = \max_{M \in \mathcal{M}_1} \min_{U \in \mathcal{U}} \text{tr}((\hat{\Sigma} - U)M) \\ & \leq \min_{U \in \mathcal{U}} \max_{M \in \mathcal{M}_1} \text{tr}((\hat{\Sigma} - U)M) \\ & = \min_{U \in \mathcal{U}} \lambda_1(\hat{\Sigma} - U). \end{aligned}$$

It is convenient to denote  $\gamma := \sqrt{\frac{k^2 \log p}{n\theta^2}}$ , and note that

$$\gamma \leq \frac{\sqrt{k}}{16} \min_{j \in S} |v_{1,j}| \leq \frac{1}{16} \|v_{1,S}\|_2 = \frac{1}{16}.$$

PROOF OF (a). From (20), it suffices to exhibit a primal-dual pair  $(\hat{M}, \hat{U}) \in \mathcal{M}_1 \times \mathcal{U}$ , such that:

$$\begin{aligned} \text{(C1)} \quad & \hat{M} = \hat{v} \hat{v}^\top \text{ with } \text{sgn}(\hat{v}) = \text{sgn}(v_1). \\ \text{(C2)} \quad & \text{tr}(\hat{\Sigma} \hat{M}) - \lambda \|\hat{M}\|_1 = \lambda_1(\hat{\Sigma} - \hat{U}). \end{aligned}$$

We construct the primal-dual pair as follows. Define

$$\hat{U} := \begin{pmatrix} \lambda \text{sgn}(v_{1,S}) \text{sgn}(v_{1,S})^\top & \hat{\Sigma}_{SS^c} - \Sigma_{SS^c} \\ \hat{\Sigma}_{S^cS} - \Sigma_{S^cS} & \hat{\Sigma}_{S^cS^c} - \Sigma_{S^cS^c} \end{pmatrix}.$$

By (19) and Lemma 5, we have that  $\|\hat{\Sigma} - \Sigma\|_\infty \leq 4\sqrt{\frac{\log p}{n}} \leq \lambda$ , so  $U \in \mathcal{U}$ . Let  $w = (w_1, \dots, w_k)$  be a unit-length leading eigenvector of  $\Sigma_{SS} - \hat{U}_{SS}$  such that  $w^\top v_{1,S} \geq 0$ . Then define  $\hat{v}$  componentwise by

$$\hat{v}_S \in \underset{\substack{u \in \mathbb{R}^k, \|u\|_2=1 \\ u^\top w \geq 0}}{\text{argmax}} u^\top (\hat{\Sigma}_{SS} - \hat{U}_{SS})u, \quad \hat{v}_{S^c} = 0,$$

and set  $\hat{M} := \hat{v} \hat{v}^\top$ . Note that our choices above ensure that  $\hat{M} \in \mathcal{M}_1$ . To verify (C1), we now show that  $\text{sgn}(\hat{v}_S) = \text{sgn}(w) = \text{sgn}(v_{1,S})$ . By a variant of the Davis–Kahan theorem [Yu, Wang and Samworth (2015), Theorem 2],

$$\begin{aligned} \|w - \hat{v}_S\|_\infty &\leq \|w - \hat{v}_S\|_2 \leq \sqrt{2}L(\hat{v}_S, w) \leq \frac{2\sqrt{2}\|\hat{\Sigma}_{SS} - \Sigma_{SS}\|_{\text{op}}}{\theta} \\ (21) \quad &\leq \frac{2\sqrt{2}}{\theta} \sup_{u \in B_0(2k)} |\hat{V}(u) - V(u)| \leq 4\sqrt{2}\gamma k^{-1/2}, \end{aligned}$$

where the final inequality uses (19). But  $w$  is also a leading eigenvector of

$$\frac{1}{\theta}(\Sigma_{SS} - \hat{U}_{SS} - I_k) = v_{1,S}v_{1,S}^\top - 4\gamma s s^\top,$$

where  $s := \frac{\text{sgn}(v_{1,S})}{\|\text{sgn}(v_{1,S})\|}$ . Write  $s = \alpha v_{1,S} + \beta v_\perp$  for some  $\alpha, \beta \in \mathbb{R}$  with  $\alpha^2 + \beta^2 = 1$ , and a unit vector  $v_\perp \in \mathbb{R}^k$  orthogonal to  $v_{1,S}$ . Then

$$\begin{aligned} v_{1,S}v_{1,S}^\top - 4\gamma s s^\top &= (v_{1,S} \quad v_\perp) \begin{pmatrix} 1 - 4\gamma\alpha^2 & -4\gamma\alpha\beta \\ -4\gamma\alpha\beta & -4\gamma\beta^2 \end{pmatrix} \begin{pmatrix} v_{1,S}^\top \\ v_\perp^\top \end{pmatrix} \\ &= (v_{1,S} \quad v_\perp) \begin{pmatrix} a_1 & b_1 \\ a_2 & b_2 \end{pmatrix} \begin{pmatrix} d_1 & 0 \\ 0 & d_2 \end{pmatrix} \begin{pmatrix} a_1 & a_2 \\ b_1 & b_2 \end{pmatrix} \begin{pmatrix} v_{1,S}^\top \\ v_\perp^\top \end{pmatrix}, \end{aligned}$$

where  $d_1 \geq d_2$  and  $(a_1 \ a_2)^\top, (b_1 \ b_2)^\top$  are eigenvalues and corresponding unit-length eigenvectors of the middle matrix on the right-hand side of the first line. Direct computation yields that  $d_1 \geq 1/2 > 0 \geq d_2$  and

$$\begin{pmatrix} a_1 \\ a_2 \end{pmatrix} \propto \begin{pmatrix} 1 - 4\gamma\alpha^2 + 4\gamma\beta^2 + \sqrt{16\gamma\beta^2 + (1 - 4\gamma)^2} \\ -8\gamma\alpha\beta \end{pmatrix}.$$

Consequently,  $w$  is a scalar multiple of

$$(22) \quad a_1 v_{1,S} + a_2 v_\perp = \{1 + 4\gamma + \sqrt{16\gamma\beta^2 + (1 - 4\gamma)^2}\} v_{1,S} - 8\gamma\alpha s.$$

Since

$$\begin{aligned} \{1 + 4\gamma + \sqrt{16\gamma\beta^2 + (1 - 4\gamma)^2}\} \min_{j \in S} |v_{1,j}| &\geq 2 \min_{j \in S} |v_{1,j}| \geq 32\gamma k^{-1/2} \\ &> 8\gamma\alpha \|s\|_\infty, \end{aligned}$$

we have  $\text{sgn}(w) = \text{sgn}(v_{1,S})$ . Hence, by (22),

$$\begin{aligned} \min_{j=1,\dots,k} |w_j| &\geq \frac{\{1 + 4\gamma + \sqrt{16\gamma\beta^2 + (1 - 4\gamma)^2}\} \min_{j \in S} |v_{1,j}| - 8\gamma\alpha \|s\|_\infty}{\|a_1 v_{1,S} + a_2 v_\perp\|_2} \\ (23) \quad &\geq \frac{(32 - 8\alpha)\gamma k^{-1/2}}{1 + 4\gamma + \sqrt{16\gamma\beta^2 + (1 - 4\gamma)^2}} \\ &\geq \frac{12\gamma k^{-1/2}}{1 + 4\gamma} \geq \frac{48}{5}\gamma k^{-1/2}. \end{aligned}$$

By (21) and (23), we have  $\min_j |w_j| > \|w - \hat{v}_S\|_\infty$ . So  $\text{sgn}(\hat{v}_S) = \text{sgn}(w) = \text{sgn}(v_{1,S})$  as desired.

It remains to check condition (C2). Since  $\text{sgn}(\hat{v}_S) = \text{sgn}(v_{1,S})$ , we have

$$\begin{aligned} \text{tr}(\hat{\Sigma} \hat{M}) - \lambda \|\hat{M}\|_1 &= \text{tr}(\hat{\Sigma}_{SS} \hat{v}_S \hat{v}_S^\top) - \text{tr}(\hat{U}_{SS} \hat{v}_S \hat{v}_S^\top) \\ &= \hat{v}_S^\top (\hat{\Sigma}_{SS} - \hat{U}_{SS}) \hat{v}_S = \lambda_1(\hat{\Sigma}_{SS} - \hat{U}_{SS}). \end{aligned}$$

Moreover,

$$\hat{\Sigma} - \hat{U} = \begin{pmatrix} \hat{\Sigma}_{SS} - \hat{U}_{SS} & 0 \\ 0 & \Gamma_{p-k} \end{pmatrix}.$$

As  $\lambda_1(\Gamma_{p-k}) \leq 1$  by assumption, it suffices to show that  $\lambda_1(\hat{\Sigma}_{SS} - \hat{U}_{SS}) \geq 1$ . By Weyl's inequality [see, e.g., [Horn and Johnson \(2012\)](#), Theorem 4.3.1]

$$\begin{aligned} \lambda_1(\hat{\Sigma}_{SS} - \hat{U}_{SS}) &\geq \lambda_1(\Sigma_{SS} - \hat{U}_{SS}) - \|\hat{\Sigma}_{SS} - \Sigma_{SS}\|_{\text{op}} \\ (24) \quad &\geq 1 + \theta \lambda_1(v_{1,S} v_{1,S}^\top - 4\gamma s s^\top) - 2\sqrt{\frac{k \log p}{n}} \\ &\geq 1 + \frac{3\theta}{8} > 1, \end{aligned}$$

as required.



PROOF OF (b). We claim first that  $\hat{S} = S$ . Let  $\phi^* := f(\hat{M})$  be the optimal value of the semidefinite programme (5). From (24), we have  $\phi^* \geq 1 + 3\theta/8$ . The proof strategy here is to use dual matrices  $\hat{U}$  defined in part (a) and  $\hat{U}'$  to be defined below to respectively bound  $\text{tr}(\hat{M}_{S^c S^c}^\varepsilon)$  from above and bound  $\hat{M}_{rr}^\varepsilon$  from below for each  $r \in S$ . We then check that for the choice of  $\varepsilon$  we have in the theorem, the diagonal entries of  $\hat{M}^\varepsilon$  are above the threshold  $\log p/(6n)$  precisely when they belong to the  $(S, S)$ -block of the matrix.

From (20), and using the fact that  $\text{tr}(AB) \leq \text{tr}(A)\lambda_1(B)$  for all symmetric matrices  $A$  and  $B$ , we have

$$\begin{aligned} \text{tr}(\hat{\Sigma} \hat{M}^\varepsilon) - \lambda \|\hat{M}^\varepsilon\|_1 &\leq \text{tr}((\hat{\Sigma} - \hat{U}) \hat{M}^\varepsilon) \\ &= \text{tr}((\hat{\Sigma}_{SS} - \hat{U}_{SS}) \hat{M}_{SS}^\varepsilon) + \text{tr}(\Sigma_{S^c S^c} \hat{M}_{S^c S^c}^\varepsilon) \\ &\leq \text{tr}(\hat{M}_{SS}^\varepsilon) \phi^* + \text{tr}(\hat{M}_{S^c S^c}^\varepsilon) \lambda_1(\Gamma_{p-k}) \\ &= \phi^* - \text{tr}(\hat{M}_{S^c S^c}^\varepsilon) (\phi^* - 1) \leq \phi^* - 3\theta \text{tr}(\hat{M}_{S^c S^c}^\varepsilon)/8. \end{aligned}$$

On the other hand,  $\text{tr}(\hat{\Sigma} \hat{M}^\varepsilon) - \lambda \|\hat{M}^\varepsilon\|_1 \geq \phi^* - \varepsilon$ . It follows that

$$(25) \quad \text{tr}(\hat{M}_{S^c S^c}^\varepsilon) \leq \frac{8\varepsilon}{3\theta} \leq \frac{1}{6} \left( \frac{\log p}{Bn} \right)^2 < \tau.$$

Next, fix an arbitrary  $r \in S$  and define  $S_0 := S \setminus \{r\}$ . Define  $\hat{U}'$  by

$$\hat{U}'_{ij} := \begin{cases} \lambda \text{sgn}(\hat{M}_{ij}), & \text{if } i, j \in S_0, \\ \hat{\Sigma}_{ij} - \Sigma_{ij}, & \text{otherwise.} \end{cases}$$

We note that on  $\Omega_0$ , we have  $\hat{U}' \in \mathcal{U}$ . Again by (20),

$$\begin{aligned} \text{tr}(\hat{\Sigma} \hat{M}^\varepsilon) - \lambda \|\hat{M}^\varepsilon\|_1 &\leq \text{tr}((\hat{\Sigma} - \hat{U}') \hat{M}^\varepsilon) \\ &= \text{tr}((\hat{\Sigma}_{S_0 S_0} - \hat{U}_{S_0 S_0}) \hat{M}_{S_0 S_0}^\varepsilon) + \sum_{\substack{(i,j) \in S \times S \\ i=r \text{ or } j=r}} \Sigma_{ij} \hat{M}_{ji}^\varepsilon \\ (26) \quad &+ \text{tr}(\Sigma_{S^c S^c} \hat{M}_{S^c S^c}^\varepsilon) \\ &\leq \text{tr}(\hat{M}_{S_0 S_0}^\varepsilon) \lambda_1(\hat{\Sigma}_{S_0 S_0} - \hat{U}_{S_0 S_0}) + \sum_{\substack{(i,j) \in S \times S \\ i=r \text{ or } j=r}} \Sigma_{ij} \hat{M}_{ji}^\varepsilon \\ &+ \text{tr}(\hat{M}_{S^c S^c}^\varepsilon) \lambda_1(\Gamma_{p-k}). \end{aligned}$$

We bound the three terms of (26) separately. By Lemma 8 in Section 1 in the online supplementary material [Wang, Berthet and Samworth (2015)],

$$\begin{aligned} &\lambda_1(\hat{\Sigma}_{S_0 S_0} - \hat{U}_{S_0 S_0}) \\ &\leq \lambda_1(\hat{\Sigma}_{SS} - \hat{U}_{SS}) - \{\lambda_1(\hat{\Sigma}_{SS} - \hat{U}_{SS}) - \lambda_2(\hat{\Sigma}_{SS} - \hat{U}_{SS})\} \min_{j \in S} v_j^2. \end{aligned}$$

From (21) and (23),

$$\min_j |\hat{v}_j| \geq \min_j |w_j| - \|w - \hat{v}_S\|_\infty \geq 3.9\gamma k^{-1/2}.$$

Also, by Weyl's inequality,

$$\begin{aligned} & \lambda_1(\hat{\Sigma}_{SS} - \hat{U}_{SS}) - \lambda_2(\hat{\Sigma}_{SS} - \hat{U}_{SS}) \\ & \geq \lambda_1(\Sigma_{SS} - \hat{U}_{SS}) - \lambda_2(\Sigma_{SS} - \hat{U}_{SS}) - 2\|\hat{\Sigma}_{SS} - \Sigma_{SS}\|_{\text{op}} \\ & \geq \theta\{\lambda_1(v_{1,S}v_{1,S}^\top - 4\gamma_{SS}^\top) - \lambda_2(v_{1,S}v_{1,S}^\top - 4\gamma_{SS}^\top)\} - 4\sqrt{\frac{k \log p}{n}} \\ & \geq \theta(1/2 - 4\gamma k^{-1/2}) \geq \theta/4. \end{aligned}$$

It follows that

$$(27) \quad \lambda_1(\hat{\Sigma}_{S_0 S_0} - \hat{U}_{S_0 S_0}) \leq \phi^* - 3.8\gamma^2 k^{-1}\theta.$$

For the second term in (26), observe that

$$\begin{aligned} & \sum_{\substack{(i,j) \in S \times S \\ i=r \text{ or } j=r}} \Sigma_{ij} \hat{M}_{ij}^\varepsilon \leq (1 + \theta v_{1,r}^2) \hat{M}_{rr}^\varepsilon + 2 \sum_{i \in S, i \neq r} \theta v_{1,i} v_{1,r} \hat{M}_{i,r}^\varepsilon \\ (28) \quad & \leq \hat{M}_{rr}^\varepsilon + 2\theta |v_{1,r}| \cdot \|v_1\|_1 \sqrt{\hat{M}_{rr}^\varepsilon} \\ & \leq \hat{M}_{rr}^\varepsilon + 2\theta \sqrt{k} \sqrt{\hat{M}_{rr}^\varepsilon}, \end{aligned}$$

where the penultimate inequality uses the fact that  $\hat{M}_{ir}^\varepsilon \leq \sqrt{\hat{M}_{ii}^\varepsilon \hat{M}_{rr}^\varepsilon} \leq \sqrt{\hat{M}_{rr}^\varepsilon}$  for a nonnegative definite matrix  $\hat{M}^\varepsilon$ . Substituting (27) and (28) into (26),

$$\begin{aligned} & \text{tr}(\hat{\Sigma} \hat{M}^\varepsilon) - \lambda \|\hat{M}^\varepsilon\|_1 \\ & \leq \text{tr}(\hat{M}_{S_0 S_0}^\varepsilon) \left( \phi^* - \frac{3.8\gamma^2 \theta}{k} \right) + \hat{M}_{rr}^\varepsilon + 2\theta \sqrt{k \hat{M}_{rr}^\varepsilon} + \text{tr}(\hat{M}_{S^c S^c}^\varepsilon) \\ & \leq \phi^* - 3.8\gamma^2 k^{-1} \theta \text{tr}(\hat{M}_{S_0 S_0}^\varepsilon) + 2\theta \sqrt{k \hat{M}_{rr}^\varepsilon} \\ & \leq \phi^* - 3.8\gamma^2 k^{-1} \theta \{1 - \text{tr}(\hat{M}_{S^c S^c}^\varepsilon)\} + 2\theta (\sqrt{k} + 1.9\gamma^2) \sqrt{\hat{M}_{rr}^\varepsilon}. \end{aligned}$$

By definition,  $\text{tr}(\hat{\Sigma} \hat{M}^\varepsilon) - \lambda \|\hat{M}^\varepsilon\|_1 \geq \phi^* - \varepsilon$ , so together with (25), we have

$$\begin{aligned} & \sqrt{\hat{M}_{rr}^\varepsilon} \geq \frac{3.8\gamma^2 k^{-1} \theta (1 - (8\varepsilon)/(3\theta)) - \varepsilon}{2\theta (\sqrt{k} + 1.9\gamma^2)} \\ (29) \quad & \geq \frac{1.9\gamma^2 k^{-1} (1 - (8\varepsilon)/(3\theta))}{(\sqrt{k} + 1.9/256)} - \frac{\varepsilon}{2\theta} \end{aligned}$$

$$\begin{aligned} &\geq 1.8\gamma^2 k^{-3/2} \left(1 - \frac{8\varepsilon}{3\theta}\right) - \frac{\varepsilon}{2\theta} \\ &\geq \frac{1.8k^{1/2} \log p}{n\theta^2} \left\{1 - \frac{1}{6} \left(\frac{\log p}{Bn}\right)^2\right\} - \frac{1}{32} \left(\frac{\log p}{Bn}\right)^2 \\ &\geq \frac{1.4 \log p}{Bn} > \tau^{1/2}. \end{aligned}$$

From (25) and (29), we conclude that  $\hat{S} = S$ , as claimed.

To conclude, by [Yu, Wang and Samworth \[\(2015\), Theorem 2\]](#), on  $\Omega_0$ ,

$$L(\hat{v}^{\text{MSDP}}, v_1) = L(\hat{v}_S^{\text{MSDP}}, v_{1,S}) \leq \frac{2\|\hat{\Sigma}_{SS} - \Sigma_{SS}\|_{\text{op}}}{\lambda_1(\Sigma_{SS}) - \lambda_2(\Sigma_{SS})} \leq 4\sqrt{\frac{k \log p}{n\theta^2}},$$

where we used (19) and Lemma 5 in the online supplementary material [[Wang, Berthet and Samworth \(2015\)](#)] in the final bound.

For the final part of the theorem, when  $p \geq \theta\sqrt{n/k}$ ,

$$\begin{aligned} \sup_{P \in \tilde{\mathcal{P}}_p(n,k,\theta)} \mathbb{E}_P\{L(\hat{v}^{\text{MSDP}}, v_1)\} &\leq 4\sqrt{\frac{k \log p}{n\theta^2}} + \mathbb{P}(\Omega_0^c) \\ &\leq 4\sqrt{\frac{k \log p}{n\theta^2}} + 2p^{-1} \leq 6\sqrt{\frac{k \log p}{n\theta^2}}, \end{aligned}$$

as desired.  $\square$

**Acknowledgements.** We thank the anonymous reviewers for helpful and constructive comments on an earlier draft.

SUPPLEMENTARY MATERIAL

**Supplementary material to “Statistical and computational trade-offs in estimation of sparse principal components”** (DOI: [10.1214/15-AOS1369SUPP.pdf](https://doi.org/10.1214/15-AOS1369SUPP.pdf)). Ancillary results and a brief introduction to computational complexity theory.

REFERENCES

ALLEN, G. I. and MALETIĆ-SAVATIĆ, M. (2011). Sparse non-negative generalized PCA with applications to metabolomics. *Bioinformatics* **27** 3029–3035.

ALON, N., KRIVELEVICH, M. and SUDAKOV, B. (1998). Finding a large hidden clique in a random graph. In *Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms (San Francisco, CA, 1998)* 594–598. ACM, New York. [MR1642973](#)

ALON, N., ANDONI, A., KAUFMAN, T., MATULEF, K., RUBINFELD, R. and XIE, N. (2007). Testing  $k$ -wise and almost  $k$ -wise independence. In *STOC’07—Proceedings of the 39th Annual ACM Symposium on Theory of Computing* 496–505. ACM, New York. [MR2402475](#)

- AMES, B. P. W. and VAVASIS, S. A. (2011). Nuclear norm minimization for the planted clique and biclique problems. *Math. Program.* **129** 69–89. [MR2831403](#)
- AMINI, A. A. and WAINWRIGHT, M. J. (2009). High-dimensional analysis of semidefinite relaxations for sparse principal components. *Ann. Statist.* **37** 2877–2921. [MR2541450](#)
- APPLEBAUM, B., BARAK, B. and WIGDERSON, A. (2010). Public-key cryptography from different assumptions. In *STOC'10—Proceedings of the 2010 ACM International Symposium on Theory of Computing* 171–180. ACM, New York. [MR2743266](#)
- BACH, F., AHİPAŞAOĞLU, S. D. and d'ASPREMONT, A. (2010). Convex relaxations for subset selection. Available at [arXiv:1006.3601](#).
- BAIK, J., BEN AROUS, G. and PÉCHÉ, S. (2005). Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. *Ann. Probab.* **33** 1643–1697. [MR2165575](#)
- BERTHET, Q. (2015). Optimal testing for planted satisfiability problems. *Electron. J. Stat.* **9** 298–317. [MR3319521](#)
- BERTHET, Q. and RIGOLLET, P. (2013a). Optimal detection of sparse principal components in high dimension. *Ann. Statist.* **41** 1780–1815. [MR3127849](#)
- BERTHET, Q. and RIGOLLET, P. (2013b). Complexity theoretic lower bounds for sparse principal component detection. *J. Mach. Learn. Res. W&CP* **30** 1046–1066.
- BIRNBAUM, A., JOHNSTONE, I. M., NADLER, B. and PAUL, D. (2013). Minimax bounds for sparse PCA with noisy high-dimensional data. *Ann. Statist.* **41** 1055–1084. [MR3113803](#)
- CAI, T. T., MA, Z. and WU, Y. (2013). Sparse PCA: Optimal rates and adaptive estimation. *Ann. Statist.* **41** 3074–3110. [MR3161458](#)
- CHAN, Y.-B. and HALL, P. (2010). Using evidence of mixed populations to select variables for clustering very high-dimensional data. *J. Amer. Statist. Assoc.* **105** 798–809. [MR2724862](#)
- CHANDRASEKARAN, V. and JORDAN, M. I. (2013). Computational and statistical tradeoffs via convex relaxation. *Proc. Natl. Acad. Sci. USA* **110** E1181–E1190. [MR3047651](#)
- CHEN, Y. and XU, J. (2014). Statistical-computational tradeoffs in planted problems and submatrix localization with a growing number of clusters and submatrices. Available at [arXiv:1402.1267](#).
- CHUN, H. and SÜNDÜZ, K. (2009). Expression quantitative trait loci mapping with multivariate sparse partial least squares regression. *Genetics* **182** 79–90.
- d'ASPREMONT, A., EL GHAOU, L., JORDAN, M. I. and LANCKRIET, G. R. G. (2007). A direct formulation for sparse PCA using semidefinite programming. *SIAM Rev.* **49** 434–448 (electronic). [MR2353806](#)
- DAVIS, C. and KAHAN, W. M. (1970). The rotation of eigenvectors by a perturbation. III. *SIAM J. Numer. Anal.* **7** 1–46. [MR0264450](#)
- DESHPANDE, Y. and MONTANARI, A. (2014). Sparse PCA via covariance thresholding. Preprint. Available at [arXiv:1311.5179](#).
- DIACONIS, P. and FREEDMAN, D. (1980). Finite exchangeable sequences. *Ann. Probab.* **8** 745–764. [MR0577313](#)
- FEIGE, U. and KRAUTHGAMER, R. (2000). Finding and certifying a large hidden clique in a semi-random graph. *Random Structures Algorithms* **16** 195–208. [MR1742351](#)
- FEIGE, U. and KRAUTHGAMER, R. (2003). The probable value of the Lovász–Schrijver relaxations for maximum independent set. *SIAM J. Comput.* **32** 345–370 (electronic). [MR1969394](#)
- FEIGE, U. and RON, D. (2010). Finding hidden cliques in linear time. In *21st International Meeting on Probabilistic, Combinatorial, and Asymptotic Methods in the Analysis of Algorithms (AofA'10)* 189–203. Assoc. Discrete Math. Theor. Comput. Sci., Nancy. [MR2735341](#)
- FELDMAN, V., PERKINS, W. and VEMPALA, S. (2015). On the complexity of random satisfiability problems with planted solutions. In *STOC'15—Proceedings of the 2015 ACM Symposium on Theory of Computing* 77–86 ACM, New York.
- FELDMAN, V., GRIGORESCU, E., REYZIN, L., VEMPALA, S. S. and XIAO, Y. (2013). Statistical algorithms and a lower bound for detecting planted cliques. In *STOC'13—Proceedings of the 2013 ACM Symposium on Theory of Computing* 655–664. ACM, New York. [MR3210827](#)

- GAO, C., MA, Z. and ZHOU, H. H. (2014). Sparse CCA: Adaptive estimation and computational barriers. Available at [arXiv:1409.8565](https://arxiv.org/abs/1409.8565).
- GOLUB, G. H. and VAN LOAN, C. F. (1996). *Matrix Computations*, 3rd ed. Johns Hopkins Univ. Press, Baltimore, MD. [MR1417720](https://doi.org/10.1137/1417720)
- GRIMMETT, G. R. and MCDIARMID, C. J. H. (1975). On colouring random graphs. *Math. Proc. Cambridge Philos. Soc.* **77** 313–324. [MR0369129](https://doi.org/10.1017/S03050041000069129)
- HAJEK, B., WU, Y. and XU, J. (2014). Computational lower bounds for community detection on random graphs. Preprint. Available at [arXiv:1406.6625](https://arxiv.org/abs/1406.6625).
- HAZAN, E. and KRAUTHGAMER, R. (2011). How hard is it to approximate the best Nash equilibrium? *SIAM J. Comput.* **40** 79–91. [MR2765712](https://doi.org/10.1137/092765712)
- HORN, R. A. and JOHNSON, C. R. (2012). *Matrix Analysis*. Cambridge Univ. Press, Cambridge.
- JERRUM, M. (1992). Large cliques elude the Metropolis process. *Random Structures Algorithms* **3** 347–359. [MR1179827](https://doi.org/10.1002/rsa.327)
- 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. [MR2751448](https://doi.org/10.1198/016214508000000000)
- JOLLIFFE, I. T., TRENDAFILOV, N. T. and UDDIN, M. (2003). A modified principal component technique based on the LASSO. *J. Comput. Graph. Statist.* **12** 531–547. [MR2002634](https://doi.org/10.1198/106186003000000000)
- JOURNÉE, M., NESTEROV, Y., RICHTÁRIK, P. and SEPULCHRE, R. (2010). Generalized power method for sparse principal component analysis. *J. Mach. Learn. Res.* **11** 517–553. [MR2600619](https://doi.org/10.1162/JMLR.2010.11.1.1779310)
- JUELS, A. and PEINADO, M. (2000). Hiding cliques for cryptographic security. *Des. Codes Cryptogr.* **20** 269–280. [MR1779310](https://doi.org/10.1007/s001470100000)
- KARP, R. M. (1972). Reducibility among combinatorial problems. In *Complexity of Computer Computations* (R. M. Miller et al., eds.) 85–103. Plenum, New York. [MR0378476](https://doi.org/10.1007/978-1-4613-4666-6_2)
- KUČERA, L. (1995). Expected complexity of graph partitioning problems. *Discrete Appl. Math.* **57** 193–212. [MR1327775](https://doi.org/10.1016/0167-5060(95)00000-0)
- LANCZOS, C. (1950). An iteration method for the solution of the eigenvalue problem of linear differential and integral operators. *J. Res. Natl. Bur. Stand.* **45** 255–282. [MR0042791](https://doi.org/10.1126/science.1126271)
- LAURENT, B. and MASSART, P. (2000). Adaptive estimation of a quadratic functional by model selection. *Ann. Statist.* **28** 1302–1338. [MR1805785](https://doi.org/10.1214/aos/1013203485)
- MA, Z. (2013). Sparse principal component analysis and iterative thresholding. *Ann. Statist.* **41** 772–801. [MR3099121](https://doi.org/10.1214/12-AOS1121)
- MA, Z. and WU, Y. (2015). Computational barriers in minimax submatrix detection. *Ann. Statist.* **43** 1089–1116. [MR3346698](https://doi.org/10.1214/14-AOS1268)
- MAJUMDAR, A. (2009). Image compression by sparse PCA coding in curvelet domain. *Signal Image Video Process.* **3** 27–34.
- NAIKAL, N., YANG, A. Y. and SASTRY, S. S. (2011). Informative feature selection for object recognition via sparse PCA. In *Computer Vision (ICCV), 2011 IEEE International Conference* 818–825. IEEE, Barcelona, Spain.
- NEMIROVSKI, A. (2004). Prox-method with rate of convergence  $O(1/t)$  for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM J. Optim.* **15** 229–251 (electronic). [MR2112984](https://doi.org/10.1137/S1089395403000000)
- NESTEROV, Y. (2005). Smooth minimization of non-smooth functions. *Math. Program.* **103** 127–152. [MR2166537](https://doi.org/10.1007/s101070100000)
- PARKHOMENKO, E., TRITCHLER, D. and BEYENE, J. (2009). Sparse canonical correlation analysis with application to genomic data integration. *Stat. Appl. Genet. Mol. Biol.* **8** Art. 1, 36. [MR2471148](https://doi.org/10.1186/1471-2108-8-36)
- PAUL, D. (2007). Asymptotics of sample eigenstructure for a large dimensional spiked covariance model. *Statist. Sinica* **17** 1617–1642. [MR2399865](https://doi.org/10.1007/s11464-007-0000-0)
- SAMWORTH, R. J. (2016). Peter Hall’s work on high-dimensional data and classification. *Ann. Statist.* To appear.

- SHEN, D., SHEN, H. and MARRON, J. S. (2013). Consistency of sparse PCA in high dimension, low sample size contexts. *J. Multivariate Anal.* **115** 317–333. [MR3004561](#)
- SHORACK, G. R. and WELLNER, J. A. (1986). *Empirical Processes with Applications to Statistics*. Wiley, New York. [MR0838963](#)
- TAN, K. M., PETERSEN, A. and WITTEN, D. (2014). Classification of RNA-seq data. In *Statistical Analysis of Next Generation Sequencing Data* (S. Datta and D. Witten, eds.) 219–246. Springer, Cham. [MR3307374](#)
- VAN DE GEER, S. (2000). *Empirical Processes in M-Estimation*. Cambridge Univ. Press, Cambridge.
- VU, V. Q. and LEI, J. (2013). Minimax sparse principal subspace estimation in high dimensions. *Ann. Statist.* **41** 2905–2947. [MR3161452](#)
- VU, V. Q., CHO, J., LEI, J. and ROHE, K. (2013). Fantope projection and selection: A near-optimal convex relaxation of sparse PCA. *Advances in Neural Information Processing Systems (NIPS)* **26** 2670–2678.
- WANG, T., BERTHET, Q. and SAMWORTH, R. J. (2015). Supplement to “Statistical and computational trade-offs in estimation of sparse principal components”. DOI:[10.1214/15-AOS1369SUPP](#).
- WANG, Z., LU, H. and LIU, H. (2014). Tighten after relax: Minimax-optimal sparse PCA in polynomial time. *Advances in Neural Information Processing Systems (NIPS)* **27** 3383–3391.
- WANG, D., LU, H. and YANG, M.-H. (2013). Online object tracking with sparse prototypes. *IEEE Trans. Image Process.* **22** 314–325. [MR3017466](#)
- WITTEN, D. M., TIBSHIRANI, R. and HASTIE, T. (2009). A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics* **10** 515–534.
- YU, Y., WANG, T. and SAMWORTH, R. J. (2015). A useful variant of the Davis–Kahan theorem for statisticians. *Biometrika* **102** 315–323. [MR3371006](#)
- YUAN, X.-T. and ZHANG, T. (2013). Truncated power method for sparse eigenvalue problems. *J. Mach. Learn. Res.* **14** 899–925. [MR3063614](#)
- ZHANG, Y., WAINWRIGHT, M. J. and JORDAN, M. I. (2014). Lower bounds on the performance of polynomial-time algorithms for sparse linear regression. *J. Mach. Learn. Res. W&CP* **35** 921–948.
- ZOU, H., HASTIE, T. and TIBSHIRANI, R. (2006). Sparse principal component analysis. *J. Comput. Graph. Statist.* **15** 265–286. [MR2252527](#)

STATISTICAL LABORATORY  
WILBERFORCE ROAD  
CAMBRIDGE, CB3 0WB  
UNITED KINGDOM  
E-MAIL: [t.wang@statslab.cam.ac.uk](mailto:t.wang@statslab.cam.ac.uk)  
[q.berthet@statslab.cam.ac.uk](mailto:q.berthet@statslab.cam.ac.uk)  
[r.samworth@statslab.cam.ac.uk](mailto:r.samworth@statslab.cam.ac.uk)  
URL: <http://www.statslab.cam.ac.uk/~rjs57>  
<http://www.statslab.cam.ac.uk/~tw389>  
<http://www.statslab.cam.ac.uk/~qb204>

**SUPPLEMENTARY MATERIAL TO ‘STATISTICAL AND  
COMPUTATIONAL TRADE-OFFS IN ESTIMATION OF  
SPARSE PRINCIPAL COMPONENTS’**

BY TENG YAO WANG\* QUENTIN BERTHET\*,† RICHARD J. SAMWORTH\*

*University of Cambridge\**  
*California Institute of Technology†*

**1. Ancillary results.** We collect here various results used in the proofs in Appendices A, B and C in the main document Wang, Berthet and Samworth (2016).

PROPOSITION 1. *Let  $P \in \text{RCC}_p(n, \ell, C)$  and suppose that  $\ell \log p \leq n$ . Then*

$$\mathbb{E} \sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| \leq \left(1 + \frac{1}{\log p}\right) C \sqrt{\frac{\ell \log p}{n}}.$$

PROOF. By setting  $\delta = p^{1-t}$  in the RCC condition, we find that

$$\mathbb{P}\left(\sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| \geq C \max\left\{\sqrt{\frac{t\ell \log p}{n}}, \frac{t\ell \log p}{n}\right\}\right) \leq \min(1, p^{1-t})$$

for all  $t \geq 0$ . It follows that

$$\begin{aligned} \mathbb{E} \sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| &= \int_0^\infty \mathbb{P}\left(\sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| \geq s\right) ds \\ &\leq C \sqrt{\frac{\ell \log p}{n}} + C \sqrt{\frac{\ell \log p}{n}} \int_1^{\frac{n}{\ell \log p}} \frac{1}{2} p^{1-t} t^{-1/2} dt + C \frac{\ell \log p}{n} \int_{\frac{n}{\ell \log p}}^\infty p^{1-t} dt \\ &\leq C \sqrt{\frac{\ell \log p}{n}} \left\{1 + \int_1^\infty p^{1-t} dt\right\} = \left(1 + \frac{1}{\log p}\right) C \sqrt{\frac{\ell \log p}{n}}, \end{aligned}$$

as required. □

LEMMA 2. *Let  $\epsilon \in (0, 1/2)$ , let  $\ell \in \{1, \dots, p\}$  and let  $A \in \mathbb{R}^{p \times p}$  be a symmetric matrix. Then there exists  $\mathcal{N}_\epsilon \subseteq B_0(\ell)$  with cardinality at most  $\binom{p}{\ell} \pi \ell^{1/2} (1 - \epsilon^2/16)^{-(\ell-1)/2} (2/\epsilon)^{\ell-1}$  such that*

$$\sup_{u \in B_0(\ell)} |u^\top A u| \leq (1 - 2\epsilon)^{-1} \max_{u \in \mathcal{N}_\epsilon} |u^\top A u|.$$

PROOF. Let  $\mathcal{I}_\ell := \{I \subseteq \{1, \dots, p\} : |I| = \ell\}$ , and for  $I \in \mathcal{I}_\ell$ , let  $B_I := \{u \in B_0(\ell) : u_{I^c} = 0\}$ . Thus

$$B_0(\ell) = \bigcup_{I \in \mathcal{I}_\ell} B_I.$$

For each  $I \in \mathcal{I}_\ell$ , by Lemma 10 of [Kim and Samworth \(2014\)](#), there exists  $\mathcal{N}_{I,\epsilon} \subseteq B_I$  such that  $|\mathcal{N}_{I,\epsilon}| \leq \pi \ell^{1/2} (1 - \epsilon^2/16)^{-(\ell-1)/2} (2/\epsilon)^{\ell-1}$  and such that for any  $x \in B_I$ , there exists  $x' \in \mathcal{N}_{I,\epsilon}$  with  $\|x - x'\| \leq \epsilon$ . Let  $u_I \in \operatorname{argmax}_{u \in B_I} |u^\top A u|$  and find  $v_I \in \mathcal{N}_{I,\epsilon}$  such that  $\|u_I - v_I\| \leq \epsilon$ . Then

$$\begin{aligned} |u_I^\top A u_I| &\leq |v_I^\top A v_I| + |(u_I - v_I)^\top A v_I| + |u_I^\top A (u_I - v_I)| \\ &\leq \max_{u \in \mathcal{N}_{I,\epsilon}} |u^\top A u| + 2\epsilon |u_I^\top A u_I|. \end{aligned}$$

Writing  $\mathcal{N}_\epsilon := \bigcup_{I \in \mathcal{I}_\ell} \mathcal{N}_{I,\epsilon}$ , we note that  $\mathcal{N}_\epsilon$  has cardinality no larger than  $\binom{p}{\ell} \pi \ell^{1/2} (1 - \epsilon^2/16)^{-(\ell-1)/2} (2/\epsilon)^{\ell-1}$  and that

$$\begin{aligned} \sup_{u \in B_0(\ell)} |u^\top A u| &= \max_{I \in \mathcal{I}_\ell} \sup_{u \in B_I} |u^\top A u| \leq (1 - 2\epsilon)^{-1} \max_{I \in \mathcal{I}_\ell} \max_{u \in \mathcal{N}_{I,\epsilon}} |u^\top A u| \\ &= (1 - 2\epsilon)^{-1} \max_{u \in \mathcal{N}_\epsilon} |u^\top A u|, \end{aligned}$$

as required.  $\square$

LEMMA 3 (Variant of the Gilbert–Varshamov Lemma). *Let  $\alpha, \beta \in (0, 1)$  and  $k, p \in \mathbb{N}$  be such that  $k \leq \alpha\beta p$ . Writing  $\mathcal{S} := \{x = (x_1, \dots, x_p)^\top \in \{0, 1\}^p : \sum_{j=1}^p x_j = k\}$ , there exists a subset  $\mathcal{S}_0$  of  $\mathcal{S}$  such that for all distinct  $x = (x_1, \dots, x_p)^\top, y = (y_1, \dots, y_p)^\top \in \mathcal{S}_0$ , we have  $\sum_{j=1}^p \mathbb{1}_{\{x_j \neq y_j\}} \geq 2(1 - \alpha)k$  and such that*

$$\log |\mathcal{S}_0| \geq \rho k \log(p/k),$$

where  $\rho := \frac{\alpha}{-\log(\alpha\beta)} (-\log \beta + \beta - 1)$ .

PROOF. See [Massart \(2007, Lemma 4.10\)](#).  $\square$

Let  $P$  and  $Q$  be two probability measures on a measurable space  $(\mathcal{X}, \mathcal{B})$ . Recall that if  $P$  is absolutely continuous with respect to  $Q$ , then the Kullback–Leibler divergence between  $P$  and  $Q$  is  $D(P\|Q) := \int_{\mathcal{X}} \log(dP/dQ) dP$ , where  $dP/dQ$  denotes the Radon–Nikodym derivative of  $P$  with respect to  $Q$ . If  $P$  is not absolutely continuous with respect to  $Q$ , we set  $D(P\|Q) := \infty$ .



LEMMA 4 (Generalised Fano's Lemma). *Let  $P_1, \dots, P_M$  be probability distributions on a measurable space  $(\mathcal{X}, \mathcal{B})$ , and assume that  $D(P_i \| P_j) \leq \beta$  for all  $i \neq j$ . Then any measurable function  $\hat{\psi} : \mathcal{X} \rightarrow \{1, \dots, M\}$  satisfies*

$$\max_{1 \leq i \leq M} P_i(\hat{\psi} \neq i) \geq 1 - \frac{\beta + \log 2}{\log M}.$$

PROOF. See Yu (1997, Lemma 3).  $\square$

LEMMA 5. *Suppose that  $P \in \mathcal{P}$  and that  $X_1, \dots, X_n \stackrel{iid}{\sim} P$ . Let  $\Sigma := \int_{\mathbb{R}^p} xx^\top dP(x)$  and  $\hat{\Sigma} := n^{-1} \sum_{i=1}^n X_i X_i^\top$ . If  $V(u) := \mathbb{E}\{(u^\top X_1)^2\}$  and  $\hat{V}(u) := n^{-1} \sum_{i=1}^n (u^\top X_i)^2$  for  $u \in B_0(2)$ , then*

$$\|\hat{\Sigma} - \Sigma\|_\infty \leq 2 \sup_{u \in B_0(2)} |\hat{V}(u) - V(u)|.$$

PROOF. Let  $e_r$  denote the  $r$ th standard basis vector in  $\mathbb{R}^p$  and write  $X_i = (X_{i,1}, \dots, X_{i,p})^\top$ . Then

$$\begin{aligned} \|\hat{\Sigma} - \Sigma\|_\infty &= \max_{r,s \in \{1, \dots, p\}} \left| \frac{1}{n} \sum_{i=1}^n (X_{i,r} X_{i,s}) - \mathbb{E}(X_{1,r} X_{1,s}) \right| \\ &\leq \max_{r,s \in \{1, \dots, p\}} \left| \frac{1}{n} \sum_{i=1}^n \left\{ \left( \frac{1}{2} e_r + \frac{1}{2} e_s \right)^\top X_i \right\}^2 - \mathbb{E} \left[ \left\{ \left( \frac{1}{2} e_r + \frac{1}{2} e_s \right)^\top X_1 \right\}^2 \right] \right| \\ &\quad + \max_{r,s \in \{1, \dots, p\}} \left| \frac{1}{n} \sum_{i=1}^n \left\{ \left( \frac{1}{2} e_r - \frac{1}{2} e_s \right)^\top X_i \right\}^2 - \mathbb{E} \left[ \left\{ \left( \frac{1}{2} e_r - \frac{1}{2} e_s \right)^\top X_1 \right\}^2 \right] \right| \\ &\leq 2 \sup_{u \in B_0(2)} |\hat{V}(u) - V(u)|, \end{aligned}$$

as required.  $\square$

Recall the definition of the Graph Vector distribution  $\text{GV}_p^g(\pi_0)$  from the proof of Theorem 6 in the main document Wang, Berthet and Samworth (2016).

LEMMA 6. *Let  $g = (g_1, \dots, g_p)^\top \in \{0, 1\}^p$ , and let  $Y_1, \dots, Y_n$  be independent random vectors, each distributed as  $\text{GV}_p^g(\pi_0)$  for some  $\pi_0 \in (0, 1/2]$ . For any  $u \in B_0(\ell)$ , let  $V(u) := \mathbb{E}\{(u^\top Y_1)^2\}$  and  $\hat{V}(u) := n^{-1} \sum_{i=1}^n (u^\top Y_i)^2$ . Then for every  $1 \leq \ell \leq 2/\pi_0$ , every  $n \in \mathbb{N}$  and every  $\delta > 0$ ,*

$$\mathbb{P} \left[ \sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| \geq 750 \max \left\{ \sqrt{\frac{\ell \log(p/\delta)}{n}}, \frac{\ell \log(p/\delta)}{n} \right\} \right] \leq \delta.$$

*In other words,  $\text{GV}_p^g(\pi_0) \in \text{RCC}_p(\ell, 750)$  for all  $\pi_0 \in (0, 1/2]$  and  $\ell \leq 2/\pi_0$ .*

PROOF. We can write

$$Y_i = \xi_i \{(1 - \epsilon_i)R_i + \epsilon_i(g + \tilde{R}_i)\},$$

where  $\xi_i$ ,  $\epsilon_i$  and  $R_i$  are independent, where  $\xi_i$  is a Rademacher random variable, where  $\epsilon_i \sim \text{Bern}(\pi_0)$ , where  $R_i = (r_{i1}, \dots, r_{ip})^\top$  has independent Rademacher coordinates, and where  $\tilde{R}_i = (\tilde{r}_{i1}, \dots, \tilde{r}_{ip})^\top$  with  $\tilde{r}_{ij} := (1 - g_j)r_{ij}$ . Thus, for any  $u \in B_0(\ell)$ , we have

$$(u^\top Y_i)^2 = (1 - \epsilon_i)(u^\top R_i)^2 + \epsilon_i(u^\top g)^2 + \epsilon_i(u^\top \tilde{R}_i)^2 + 2\epsilon_i(u^\top \tilde{R}_i)(u^\top g).$$

Hence, writing  $S := \{j : g_j = 1\}$ ,

$$\begin{aligned} |\hat{V}(u) - V(u)| &\leq \left| \frac{1}{n} \sum_{i=1}^n (1 - \epsilon_i)(u^\top R_i)^2 - (1 - \pi_0) \right| + \frac{(u^\top g)^2}{n} \left| \sum_{i=1}^n (\epsilon_i - \pi_0) \right| \\ &\quad + \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (u^\top \tilde{R}_i)^2 - \pi_0 \|u_{S^c}\|_2^2 \right| + \left| \frac{2u^\top g}{n} \sum_{i=1}^n \epsilon_i (u^\top \tilde{R}_i) \right| \\ &\leq \left| \frac{1}{n} \sum_{i=1}^n (1 - \epsilon_i) \{(u^\top R_i)^2 - 1\} \right| + \frac{1 + (u^\top g)^2 + \|u_{S^c}\|_2^2}{n} \left| \sum_{i=1}^n (\epsilon_i - \pi_0) \right| \\ (1) \quad &\quad + \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \{(u^\top \tilde{R}_i)^2 - \|u_{S^c}\|_2^2\} \right| + \left| \frac{2u^\top g}{n} \sum_{i=1}^n \epsilon_i (u^\top \tilde{R}_i) \right|. \end{aligned}$$

We now control the four terms on the right-hand side of (1) separately. For the first term, note that the distribution of  $R_i$  is subgaussian with parameter 1. Writing  $N_\epsilon := \sum_{i=1}^n \epsilon_i$ , it follows by the same argument as in the proof of Proposition 1(i) in Wang, Berthet and Samworth (2016) that for any  $s > 0$ ,

$$\begin{aligned} &\mathbb{P} \left( \sup_{u \in B_0(\ell)} \left| \frac{1}{n} \sum_{i=1}^n (1 - \epsilon_i) \{(u^\top R_i)^2 - 1\} \right| \geq 2s \right) \\ &= \mathbb{E} \left\{ \mathbb{P} \left( \sup_{u \in B_0(\ell)} \left| \frac{1}{n - N_\epsilon} \sum_{i: \epsilon_i=0} \{(u^\top R_i)^2 - 1\} \right| \geq \frac{2ns}{n - N_\epsilon} \mid N_\epsilon \right) \right\} \\ &\leq e^9 p^\ell \mathbb{E} \left[ \exp \left\{ -\frac{n \left( \frac{ns}{n - N_\epsilon} \right)^2}{4 \left( \frac{ns}{n - N_\epsilon} \right) + 32} \right\} \right] \leq e^9 p^\ell \exp \left( -\frac{ns^2}{4s + 32} \right). \end{aligned}$$

We deduce that for any  $\delta > 0$ ,

$$\begin{aligned} &\mathbb{P} \left( \sup_{u \in B_0(\ell)} \left| \frac{1}{n} \sum_{i=1}^n (1 - \epsilon_i) \{(u^\top R_i)^2 - 1\} \right| \geq 16 \max \left\{ \sqrt{\frac{\ell \log(p/\delta)}{n}}, \frac{\ell \log(p/\delta)}{n} \right\} \right) \\ (2) \quad &\leq e^9 \delta. \end{aligned}$$

For the second term on the right-hand side of (1), note first that for any  $u \in B_0(\ell)$ , we have by Cauchy–Schwarz that

$$(u^\top g)^2 \leq \|u_S\|_0 \|u_S\|_2^2 \leq \|u_S\|_0 \leq \ell.$$

We deduce using Bernstein’s inequality for Binomial random variables (e.g. [Shorack and Wellner, 1986](#), p. 855) that for any  $s > 0$ ,

$$\begin{aligned} & \mathbb{P} \left\{ \sup_{u \in B_0(\ell)} \frac{1 + (u^\top g)^2 + \|u_{S^c}\|_2^2}{n} \left| \sum_{i=1}^n (\epsilon_i - \pi_0) \right| \geq s \right\} \\ & \leq \mathbb{P} \left\{ \frac{1}{n} \left| \sum_{i=1}^n (\epsilon_i - \pi_0) \right| \geq \frac{s}{3\ell} \right\} \leq 2 \exp \left( -\frac{ns^2}{18\ell^2\pi_0 + 2s\ell} \right) \\ & \leq 2 \max \left\{ \exp \left( -\frac{ns^2}{(19 + \sqrt{37})\ell^2\pi_0} \right), \exp \left( -\frac{ns}{(1 + \sqrt{37})\ell} \right) \right\}. \end{aligned}$$

By assumption,  $\ell\pi_0 \leq 2$ . Hence, for any  $\delta > 0$ ,

$$(3) \quad \mathbb{P} \left\{ \sup_{u \in B_0(\ell)} \frac{1 + (u^\top g)^2 + \|u_{S^c}\|_2^2}{n} \left| \sum_{i=1}^n (\epsilon_i - \pi_0) \right| \geq (1 + \sqrt{37}) \max \left( \sqrt{\frac{\ell \log(1/\delta)}{n}}, \frac{\ell \log(1/\delta)}{n} \right) \right\} \leq 2\delta.$$

The third term on the right-hand side of (1) can be handled in a very similar way to the first. We find that for every  $\delta > 0$ ,

$$(4) \quad \mathbb{P} \left( \sup_{u \in B_0(\ell)} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \{ (u^\top \tilde{R}_i)^2 - \|u_{S^c}\|_2^2 \} \right| \geq 16 \max \left\{ \sqrt{\frac{\ell \log(p/\delta)}{n}}, \frac{\ell \log(p/\delta)}{n} \right\} \right) \leq e^9 \delta.$$

For the final term, by definition of  $\tilde{R}_i$ , we have for any  $u \in B_0(\ell)$  that

$$\left| \frac{2u^\top g}{n} \sum_{i=1}^n \epsilon_i (u^\top \tilde{R}_i) \right| \leq \frac{2\ell^{1/2}}{n} \left| \sum_{j:g_j=0} u_j \sum_{i:\epsilon_i=1} r_{ij} \right| \leq \frac{2\ell}{n} \max_{j:g_j=0} \left| \sum_{i:\epsilon_i=1} r_{ij} \right|.$$

Hence by Hoeffding's inequality, for any  $s > 0$ ,

$$\begin{aligned} \mathbb{P}\left\{\sup_{u \in B_0(\ell)} \left| \frac{2u^\top g}{n} \sum_{i=1}^n \epsilon_i(u^\top \tilde{R}_i) \right| \geq s\right\} &\leq \mathbb{E}\left\{\mathbb{P}\left(\max_{1 \leq j \leq p} \left| \sum_{i: \epsilon_i=1} r_{ij} \right| \geq \frac{ns}{2\ell} \mid N_\epsilon\right)\right\} \\ &\leq 2p\mathbb{E}\left\{\exp\left(-\frac{n^2 s^2}{8\ell^2 N_\epsilon}\right)\right\} \leq 2p \inf_{t>0} \left\{\exp\left(-\frac{n^2 s^2}{8\ell^2 t}\right) + \mathbb{P}(N_\epsilon > t)\right\} \\ &\leq 2p \inf_{t>0} \left\{\exp\left(-\frac{n^2 s^2}{8\ell^2 t}\right) + \exp\left(-t \log \frac{t}{n\pi_0} + t - n\pi_0\right)\right\}, \end{aligned}$$

where the final line follows by Bennett's inequality (e.g. [Shorack and Wellner, 1986](#), p. 440). Choosing  $t = \max(e^2 n\pi_0, \frac{ns}{2^{3/2}\ell})$ , we find

$$\begin{aligned} \mathbb{P}\left\{\sup_{u \in B_0(\ell)} \left| \frac{2u^\top g}{n} \sum_{i=1}^n \epsilon_i(u^\top \tilde{R}_i) \right| \geq s\right\} \\ \leq 2p \max\left\{\exp\left(-\frac{ns^2}{8e^2\ell^2\pi_0}\right) + \exp\left(-\frac{ns}{2^{3/2}\ell}\right), 2\exp\left(-\frac{ns}{2^{3/2}\ell}\right)\right\} \\ \leq 4p \max\left\{\exp\left(-\frac{ns^2}{16e^2\ell}\right), \exp\left(-\frac{ns}{2^{3/2}\ell}\right)\right\}. \end{aligned}$$

We deduce that for any  $\delta > 0$ ,

$$(5) \quad \mathbb{P}\left[\sup_{u \in B_0(\ell)} \left| \frac{2u^\top g}{n} \sum_{i=1}^n \epsilon_i(u^\top \tilde{R}_i) \right| \geq 4e \max\left\{\sqrt{\frac{\ell \log(p/\delta)}{n}}, \frac{\ell \log(p/\delta)}{n}\right\}\right] \leq 4\delta.$$

We conclude from (1), (2), (3), (4) and (5) that for any  $\delta > 0$ ,

$$\mathbb{P}\left[\sup_{u \in B_0(\ell)} |\hat{V}(u) - V(u)| \geq 750 \max\left\{\sqrt{\frac{\ell \log(p/\delta)}{n}}, \frac{\ell \log(p/\delta)}{n}\right\}\right] \leq \delta,$$

as required.  $\square$

**LEMMA 7.** *Let  $v = (v_1, \dots, v_p)^\top \in B_0(k)$  and let  $\hat{v} = (\hat{v}_1, \dots, \hat{v}_p)^\top \in \mathbb{R}^p$  be such that  $\|\hat{v}\|_2 = 1$ . Let  $S := \{j \in \{1, \dots, p\} : v_j \neq 0\}$ . Then for any  $\hat{S} \in \operatorname{argmax}_{1 \leq j_1 < \dots < j_k \leq p} \sum_{r=1}^k |\hat{v}_{j_r}|$ , we have*

$$L(\hat{v}, v)^2 \geq \frac{1}{2} \sum_{j \in S \setminus \hat{S}} v_j^2.$$

PROOF. By the Cauchy–Schwarz inequality, and then by definition of  $\hat{S}$ ,

$$\begin{aligned} 1 - L(\hat{v}, v)^2 &= \left( \sum_{j \in S \setminus \hat{S}} \hat{v}_j v_j + \sum_{j \in S \cap \hat{S}} \hat{v}_j v_j \right)^2 \\ &\leq \left( 2 \sum_{j \in S \setminus \hat{S}} \hat{v}_j^2 + \sum_{j \in S \cap \hat{S}} \hat{v}_j^2 \right) \left( \frac{1}{2} \sum_{j \in S \setminus \hat{S}} v_j^2 + \sum_{j \in S \cap \hat{S}} v_j^2 \right) \\ &\leq \left( \sum_{j \in \hat{S} \setminus S} \hat{v}_j^2 + \sum_{j \in S \setminus \hat{S}} \hat{v}_j^2 + \sum_{j \in S \cap \hat{S}} \hat{v}_j^2 \right) \left( 1 - \frac{1}{2} \sum_{j \in S \setminus \hat{S}} v_j^2 \right) \leq 1 - \frac{1}{2} \sum_{j \in S \setminus \hat{S}} v_j^2, \end{aligned}$$

as required.  $\square$

LEMMA 8. *Let  $A \in \mathbb{R}^{d \times d}$  be a symmetric matrix. Let  $A^{(r)}$  be the principal submatrix of  $A$  obtained by deleting the  $r$ th row and  $r$ th column of  $A$ . If  $A$  has a unique (up to sign) leading eigenvector  $v$ , then*

$$\lambda_2(A) \leq \lambda_1(A^{(r)}) \leq \lambda_1(A) - v_{1,r}^2(\lambda_1(A) - \lambda_2(A))$$

PROOF. The first inequality in the lemma is implied by Cauchy’s Interlacing Theorem (see, e.g. [Horn and Johnson \(2012, Theorem 4.3.17\)](#)). It remains to show the second inequality. Let  $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_d$  be eigenvalues of  $A$  (counting multiplicities), and  $v_1, \dots, v_d$  be unit-length eigenvectors of  $A$  such that  $Av_i = \lambda_i v_i$  and  $v_i^\top v_j = 0$  for all  $i \neq j$ . We have

$$\begin{aligned} \lambda_1(A^{(r)}) &= \max_{\substack{\|u\|_2=1 \\ u_r=0}} u^\top A u = \max_{\substack{\|u\|_2=1 \\ u_r=0}} u^\top \left( \sum_{i=1}^d \lambda_i v_i v_i^\top \right) u \\ &\leq \max_{\substack{\|u\|_2=1 \\ u_r=0}} \left\{ (\lambda_1 - \lambda_2) u^\top v_1 v_1^\top u + \lambda_2 u^\top \left( \sum_{i=1}^d v_i v_i^\top \right) u \right\} \\ &\leq \max_{\substack{\|u\|_2=1 \\ u_r=0}} (\lambda_1 - \lambda_2) |u^\top v_1|^2 + \lambda_2 \\ &\leq (\lambda_1 - \lambda_2)(1 - v_{1,r}^2) + \lambda_2 \\ &= \lambda_1 - v_{1,r}^2(\lambda_1 - \lambda_2), \end{aligned}$$

where we used Cauchy–Schwarz inequality in the penultimate line.  $\square$

Recall the definition of the total variation distance  $d_{\text{TV}}$  given in the proof of Theorem 6 in the main document [Wang, Berthet and Samworth \(2016\)](#).

LEMMA 9. *Let  $X$  and  $Y$  be random elements taking values in a measurable space  $(F, \mathcal{F})$ , and let  $(G, \mathcal{G})$  be another measurable space.*

(a) *If  $\phi : F \rightarrow G$  is measurable, then*

$$d_{\text{TV}}(\mathcal{L}(\phi(X)), \mathcal{L}(\phi(Y))) \leq d_{\text{TV}}(\mathcal{L}(X), \mathcal{L}(Y)).$$

(b) *Let  $Z$  be a random element taking values in  $(G, \mathcal{G})$ , and suppose that  $Z$  is independent of  $(X, Y)$ . Then*

$$d_{\text{TV}}(\mathcal{L}(X, Z), \mathcal{L}(Y, Z)) = d_{\text{TV}}(\mathcal{L}(X), \mathcal{L}(Y)).$$

PROOF. (a) For any  $A \in \mathcal{G}$ , we have

$$\begin{aligned} |\mathbb{P}\{\phi(X) \in A\} - \mathbb{P}\{\phi(Y) \in A\}| &= |\mathbb{P}\{X \in \phi^{-1}(A)\} - \mathbb{P}\{Y \in \phi^{-1}(A)\}| \\ &\leq d_{\text{TV}}(\mathcal{L}(X), \mathcal{L}(Y)). \end{aligned}$$

Since  $A \in \mathcal{G}$  was arbitrary, the result follows.

(b) Define  $\phi : F \times G \rightarrow F$  by  $\phi(w, z) := w$ . Then  $\phi$  is measurable, and using the result of part (a),

$$\begin{aligned} d_{\text{TV}}(\mathcal{L}(X), \mathcal{L}(Y)) &= d_{\text{TV}}(\mathcal{L}(\phi(X, Z)), \mathcal{L}(\phi(Y, Z))) \\ &\leq d_{\text{TV}}(\mathcal{L}(X, Z), \mathcal{L}(Y, Z)). \end{aligned}$$

For the other inequality, let  $\mathcal{A}$  denote the set of subsets  $A$  of  $\mathcal{F} \otimes \mathcal{G}$  with the property that given  $\epsilon > 0$ , there exist sets  $B_{1,F}, \dots, B_{n,F} \in \mathcal{F}$  and disjoint sets  $B_{1,G}, \dots, B_{n,G} \in \mathcal{G}$  such that, writing  $B := \cup_{i=1}^n (B_{i,F} \times B_{i,G})$ , we have  $\mathbb{P}((X, Z) \in A \Delta B) < \epsilon$  and  $\mathbb{P}((Y, Z) \in A \Delta B) < \epsilon$ . Here, the binary operator  $\Delta$  denotes the symmetric difference of two sets, so that  $A \Delta B := (A \cap B^c) \cup (A^c \cap B)$ . Note that  $\mathcal{F} \times \mathcal{G} \subseteq \mathcal{A}$ . Now suppose  $A \in \mathcal{A}$  so that, given  $\epsilon > 0$ , we can find sets  $B_{1,F}, \dots, B_{n,F} \in \mathcal{F}$  and disjoint sets  $B_{1,G}, \dots, B_{n,G} \in \mathcal{G}$  with the properties above. Observe that we can write

$$B^c = \bigcup_{I \subseteq \{1, \dots, n\}} \left( \bigcap_{i \in I} B_{i,F}^c \times \bigcap_{i \in I} B_{i,G} \cap \bigcap_{i \in I^c} B_{i,G}^c \right).$$

For each  $I \subseteq \{1, \dots, n\}$ , the sets  $\cap_{i \in I} B_{i,F}^c$  belong to  $\mathcal{F}$ , and  $\{\cap_{i \in I} B_{i,G} \cap \cap_{i \in I^c} B_{i,G}^c : I \subseteq \{1, \dots, n\}\}$  is a family of disjoint sets in  $\mathcal{G}$ . Moreover,

$$\mathbb{P}((X, Z) \in A^c \Delta B^c) = \mathbb{P}((X, Z) \in A \Delta B) < \epsilon,$$

and similarly  $\mathbb{P}((Y, Z) \in A^c \Delta B^c) < \epsilon$ . We deduce that  $A^c \in \mathcal{A}$ . Finally, if  $(A_n)$  is a disjoint sequence in  $\mathcal{A}$ , then let  $A := \cup_{n=1}^{\infty} A_n$ , and given  $\epsilon > 0$ , find

$m \in \mathbb{N}$  such that  $\mathbb{P}((X, Z) \in A \setminus \cup_{i=1}^m A_i) < \epsilon/2$  and  $\mathbb{P}((Y, Z) \in A \setminus \cup_{i=1}^m A_i) < \epsilon/2$ . Now, for each  $i = 1, \dots, m$ , find sets  $B_{i1,F}, \dots, B_{in_i,F} \in \mathcal{F}$  and disjoint sets  $B_{i1,G}, \dots, B_{in_i,G} \in \mathcal{G}$  such that, writing  $B_i := \cup_{j=1}^{n_i} (B_{ij,F} \times B_{ij,G})$ , we have  $\mathbb{P}((X, Z) \in A_i \Delta B_i) < \epsilon/(2m)$  and  $\mathbb{P}((Y, Z) \in A_i \Delta B_i) < \epsilon/(2m)$ . It is convenient to relabel the sets  $\{(B_{ij,F}, B_{ij,G}) : i = 1, \dots, m, j = 1, \dots, n_i\}$  as  $\{(C_{1,F}, C_{1,G}), \dots, (C_{N,F}, C_{N,G})\}$ , where  $N := \sum_{i=1}^m n_i$ . This means that we can write

$$\bigcup_{i=1}^m B_i = \bigcup_{k=1}^N (C_{k,F} \times C_{k,G}) = \bigcup_{K \subseteq \{1, \dots, N\}, K \neq \emptyset} \left( \bigcup_{k \in K} C_{k,F} \times \bigcap_{k \in K} C_{k,G} \cap \bigcap_{k \in K^c} C_{k,G}^c \right).$$

Now, for each non-empty subset  $K$  of  $\{1, \dots, N\}$ , the set  $\cup_{k \in K} C_{k,F}$  belongs to  $\mathcal{F}$ , and  $\{\bigcap_{k \in K} C_{k,G} \cap \bigcap_{k \in K^c} C_{k,G}^c : K \subseteq \{1, \dots, N\}, K \neq \emptyset\}$  is a family of disjoint sets in  $\mathcal{G}$ . Moreover,

$$\mathbb{P}((X, Z) \in A \Delta \cup_{i=1}^m B_i) \leq \sum_{i=1}^m \mathbb{P}((X, Z) \in A_i \Delta B_i) + \frac{\epsilon}{2} < \epsilon,$$

and similarly,  $\mathbb{P}((Y, Z) \in A \Delta \cup_{i=1}^m B_i) < \epsilon$ . We deduce that  $A \in \mathcal{A}$ , so  $\mathcal{A}$  is a  $\sigma$ -algebra containing  $\mathcal{F} \times \mathcal{G}$ , so  $\mathcal{A}$  contains  $\mathcal{F} \otimes \mathcal{G}$ .

Now suppose that  $A \in \mathcal{F} \otimes \mathcal{G}$ . By the argument above, given  $\epsilon > 0$ , there exist sets  $B_{1,F}, \dots, B_{n,F} \in \mathcal{F}$  and disjoint sets  $B_{1,G}, \dots, B_{n,G} \in \mathcal{G}$  such that  $\mathbb{P}((X, Z) \in A \Delta \cup_{i=1}^n (B_{i,F} \times B_{i,G})) < \epsilon/2$  and  $\mathbb{P}((Y, Z) \in A \Delta \cup_{i=1}^n (B_{i,F} \times B_{i,G})) < \epsilon/2$ . It follows that

$$\begin{aligned} & |\mathbb{P}((X, Z) \in A) - \mathbb{P}((Y, Z) \in A)| \\ & \leq \sum_{i=1}^n |\mathbb{P}(X \in B_{i,F}, Z \in B_{i,G}) - \mathbb{P}(Y \in B_{i,F}, Z \in B_{i,G})| + \epsilon \\ & = \sum_{i=1}^n \mathbb{P}(Z \in B_{i,G}) |\mathbb{P}(X \in B_{i,F}) - \mathbb{P}(Y \in B_{i,F})| + \epsilon \leq d_{\text{TV}}(\mathcal{L}(X), \mathcal{L}(Y)) + \epsilon. \end{aligned}$$

Since  $A \in \mathcal{A}$  and  $\epsilon > 0$  were arbitrary, we conclude that

$$d_{\text{TV}}(\mathcal{L}(X, Z), \mathcal{L}(Y, Z)) \leq d_{\text{TV}}(\mathcal{L}(X), \mathcal{L}(Y)),$$

as required.  $\square$

## 2. A brief introduction to computational complexity theory.

The following is intended to give a short introduction to notions in computational complexity theory referred to in [Wang, Berthet and Samworth](#)

(2016). A good reference for further information is [Arora and Barak \(2009\)](#), from which much of the following is inspired.

A *computational problem* is the task of generating a desired output based on a given input. Formally, defining  $\{0, 1\}^* := \cup_{k=1}^{\infty} \{0, 1\}^k$  to be the set of all finite strings of zeros and ones, we can view a computational problem as a function  $F : \{0, 1\}^* \rightarrow \mathcal{P}(\{0, 1\}^*)$ , where  $\mathcal{P}(A)$  denotes the power set of a set  $A$ . The interpretation is that  $F(s)$  describes the set of acceptable output strings (solutions) for a particular input string  $s$ .

Loosely speaking, an *algorithm* is a collection of instructions for performing a task. Despite the widespread use of algorithms in mathematics throughout history, it was not until 1936 that Alonzo Church and Alan Turing formalised the notion by defining notational systems called the  $\lambda$ -calculus and Turing machines respectively ([Church, 1936](#); [Turing, 1936](#)). Here we define an algorithm to be a *Turing machine*:

DEFINITION 1. A Turing machine  $M$  is a pair  $(Q, \delta)$ , where

- $Q$  is a finite set of states, among which are two distinguished states  $q_{\text{start}}$  and  $q_{\text{halt}}$ .
- $\delta$  is a ‘transition’ function from  $Q \times \{0, 1, \sqcup\}$  to  $Q \times \{0, 1, \sqcup\} \times \{L, R\}$ .

A Turing Machine can be thought of as having a reading head that can access a tape consisting of a countably infinite number of squares, labelled  $0, 1, 2, \dots$ . When the Turing machine is given an input  $s \in \{0, 1\}^*$ , the tape is initialised with the components of  $s$  in its first  $|s|$  tape squares (where  $|\cdot|$  denotes the length of a string in  $\{0, 1\}^*$ ) and with ‘blank symbols’  $\sqcup$  in its remaining squares. The Turing machine starts in the state  $q_{\text{start}} \in Q$  with its head on the 0th square and operates according to its transition function  $\delta$ . When the machine is in state  $q \in Q$  with its head over the  $i$ th tape square that contains the symbol  $a \in \{0, 1, \sqcup\}$ , and if  $\delta(q, a) = (q', a', L)$ , the machine overwrites  $a$  with  $a'$ , updates its state to  $q'$ , and moves to square  $i - 1$  (or to square  $i + 1$  if the third component of the transition function is R instead of L). The Turing machine stops if it reaches state  $q_{\text{halt}} \in Q$  and outputs the vector of symbols on the tape before the first blank symbol. If the Turing machine  $M$  terminates (in finitely many steps) with input  $s$ , we write  $M(s)$  for its output.

We say an algorithm (Turing machine)  $M$  *solves a computational problem*  $F$  if  $M$  terminates for every input  $s \in \{0, 1\}^*$ , and  $M(s) \in F(s)$ . A computational problem is *solvable* if there exists a Turing machine that solves it. It turns out that other notions of an algorithm (including Church’s  $\lambda$ -calculus and modern computer programming languages) are equivalent in the sense



that the set of solvable problems is the same.

A *polynomial time algorithm* is a Turing machine  $M$  for which there exist  $a, b > 0$  such that for all input strings  $s \in \{0, 1\}^*$ ,  $M$  terminates after at most  $a|s|^b$  transitions. We say a problem  $F$  is *polynomial time solvable*, written  $F \in \mathbf{P}$ , if there exists a polynomial time algorithm that solves it<sup>1</sup>.

A *nondeterministic Turing machine* has the same definition as that for a Turing machine except that the transition function  $\delta$  becomes a set-valued function  $\delta : Q \times \{0, 1, \sqcup\} \rightarrow \mathcal{P}(Q \times \{0, 1, \sqcup\} \times \{L, R\})$ . The idea is that, while in state  $q$  with its head over symbol  $a$ , a nondeterministic Turing machine replicates  $|\delta(q, a)|$  copies of itself (and its tape) in the current configuration, each exploring a different possible future configuration in the set  $\delta(q, a)$ . Each replicate branches to further replicates in the next step. The process continues until one of its replicates reaches the state  $q_{\text{halt}}$ . At that point, the Turing machine replicate that has halted outputs its tape content and all replicates stop computation. A *nondeterministic polynomial time algorithm* is a nondeterministic Turing machine  $M_{\text{nd}}$  for which there exist  $a, b > 0$  such that for all input strings  $s \in \{0, 1\}^*$ ,  $M_{\text{nd}}$  terminates after at most  $a|s|^b$  steps. (We count all replicates of  $M_{\text{nd}}$  making one parallel transition as one step.) We say a computational problem  $F$  is *nondeterministically polynomial time solvable*, written  $F \in \mathbf{NP}$ , if there exists a nondeterministic polynomial time algorithm that solves it<sup>2</sup>.

Clearly  $\mathbf{P} \subseteq \mathbf{NP}$ , but it is not currently known if these classes are equal. It is widely believed that  $\mathbf{P} \neq \mathbf{NP}$ , and many computational lower bounds for particular computational problems have been proved conditional under this assumption. Working under this hypothesis, a common strategy is to relate the algorithmic complexity of one computational problem to another. We say a computational problem  $F$  is *polynomial time reducible* to another problem  $G$ , written as  $F \leq_{\mathbf{P}} G$ , if there exist polynomial time algorithms  $M_{\text{in}}$  and  $M_{\text{out}}$  such that  $M_{\text{out}} \circ G \circ M_{\text{in}}(s) \subseteq F(s)$ . In other words,  $F \leq_{\mathbf{P}} G$  if we can convert an input of  $F$  to an input of  $G$  through  $M_{\text{in}}$ , and translate every solution of  $G$  back to a solution for  $F$  through  $M_{\text{out}}$ .

**DEFINITION 2.** *A computational problem  $G$  is NP-hard if  $F \leq_{\mathbf{P}} G$  for all  $F \in \mathbf{NP}$ . It is NP-complete if it is in NP and is NP-hard.*

[Karp \(1972\)](#) showed that a large number of natural computational prob-

<sup>1</sup>In fact, some authors write FP (short for ‘Functional Polynomial Time’) for the class we have denoted as  $\mathbf{P}$  here. The notation  $\mathbf{P}$  is then reserved for the subset of computational problems consisting of so-called *decision problems*  $F$ , where  $F(s) \in \{\{0\}, \{1\}\}$  for all  $s \in \{0, 1\}^*$ .

<sup>2</sup>Again, some authors write FNP for the class we have denoted as  $\mathbf{NP}$  here.

lems are NP-complete, including the Clique problem mentioned in Section 4. The Turing machines and nondeterministic Turing machines introduced above are both non-random. In some situations (e.g. statistical problems), it is useful to consider random procedures:

DEFINITION 3. A probabilistic Turing machine  $M_{\text{pr}}$  is a triple  $(Q, \delta, X)$ , where

- $Q$  is a finite set of states, among which are two distinguished states  $q_{\text{start}}$  and  $q_{\text{halt}}$ .
- $\delta$  is a transition function from  $Q \times \{0, 1, \sqcup\} \times \{0, 1\}$  to  $Q \times \{0, 1, \sqcup\} \times \{L, R\}$ .
- $X = (X_1, X_2, \dots)$  is an infinite sequence of independent Bern(1/2) random variables.

In its  $t$ th step, if a probabilistic Turing machine  $M_{\text{pr}}$  is in state  $q$  with its reading head over symbol  $a$ , and  $\delta(q, a, X_t) = (q', a', L)$ , then  $M_{\text{pr}}$  overwrites  $a$  with  $a'$ , updates its state to  $q'$  and moves its reading head to the left (or to the right if  $\delta(q, a, X_t) = (q', a', R)$ ). A *randomised polynomial time algorithm* is a probabilistic Turing machine  $M_{\text{pr}}$  for which there exist  $a, b > 0$  such that for any  $s \in \{0, 1\}^*$ ,  $M_{\text{pr}}$  terminates in at most  $a|s|^b$  steps. We say a computational problem  $F$  is *solvable in randomised polynomial time*, written as  $F \in \text{BPP}$ , if, given  $\epsilon > 0$ , there exists a randomised polynomial time algorithm  $M_{\text{pr}, \epsilon}$  such that  $\mathbb{P}(M_{\text{pr}, \epsilon}(s) \in F(s)) \geq 1 - \epsilon$ .

In the above discussion, the classes P, NP, BPP are all defined through worst-case performance of an algorithm, since we require the time bound to hold for every input string  $s$ . However, in many statistical applications, the input string  $s$  is drawn from some distribution  $\mathcal{D}$  on  $\{0, 1\}^*$ , and it is the average performance of the algorithm, rather than the worst case scenario, that is of more interest. We say such a random problem is solvable in randomised polynomial time if, given  $\epsilon > 0$ , there exists a randomised polynomial time algorithm  $M_{\text{pr}, \epsilon}$  such that, when  $s \sim \mathcal{D}$ , independent of  $X$ , we have  $\mathbb{P}(M_{\text{pr}, \epsilon}(s) \in F(s)) \geq 1 - \epsilon$ . Note that the probability here is taken over both the randomness in  $s$  and the randomness in  $X$ . Similar to the non-random cases, we can talk about randomised polynomial time reduction. If  $M_F$  is a randomised polynomial time algorithm for a computational problem  $F$ , then  $M_{\text{out}} \circ M_F \circ M_{\text{in}}$  is a potential randomised polynomial time algorithm for another problem  $G$  for suitably constructed randomised polynomial time algorithms  $M_{\text{in}}$  and  $M_{\text{out}}$ . One such construction is the key to the proof of Theorem 6 in the main document Wang, Berthet and Samworth (2016).

## References.

- Arora, S. and Barak, B. (2009) *Computational Complexity: A Modern Approach*. Cambridge University Press, Cambridge.
- Church, A. (1936) An unsolvable problem of elementary number theory. *Amer. J. Math.*, **58**, 345–363.
- Horn, R. A. and Johnson, C. R. (2012) *Matrix Analysis*. Cambridge University Press.
- Karp, R. M. (1972) Reducibility among combinatorial problems. In R. E. Miller et al. (Eds.), *Complexity of Computer Computations*, 85–103. Springer, New York.
- Kim, A. K.-H. and Samworth R. J. (2014) Global rates of convergence in log-concave density estimation. Available at <http://arxiv.org/abs/1404.2298>.
- Massart, P. (2007) *Concentration Inequalities and Model Selection: Ecole d'Été de Probabilités de Saint-Flour XXXIII - 2003*. Springer, Berlin/Heidelberg.
- Shorack, G. R. and Wellner, J. A. (1986) *Empirical Processes with Applications to Statistics*. Wiley, New York.
- Turing, A. (1936) On computable numbers, with an application to the Entscheidungsproblem. *Proc. London Math. Soc.*, **2**, 230–265.
- Wang, T., Berthet, Q. and Samworth, R. J. (2016) Statistical and computational trade-offs in estimation of sparse principal components. *Ann. Statist.*, **44**, this issue.
- Yu, B. (1997) Assouad, Fano and Le Cam. In Pollard, D., Torgersen, E. and Yang G. L. (Eds.) *Festschrift for Lucien Le Cam: Research Papers in Probability and Statistics*, 423–435. Springer, New York.

STATISTICAL LABORATORY  
WILBERFORCE ROAD  
CAMBRIDGE, CB3 0WB  
UNITED KINGDOM  
E-MAIL: [r.samworth@statslab.cam.ac.uk](mailto:r.samworth@statslab.cam.ac.uk)  
E-MAIL: [t.wang@statslab.cam.ac.uk](mailto:t.wang@statslab.cam.ac.uk)  
E-MAIL: [q.berthet@statslab.cam.ac.uk](mailto:q.berthet@statslab.cam.ac.uk)  
URL: <http://www.statslab.cam.ac.uk/~rjs57>  
URL: <http://www.statslab.cam.ac.uk/~tw389>  
URL: <http://www.statslab.cam.ac.uk/~qb204>