

Sparse PCA Beyond Covariance Thresholding

Gleb Novikov
ETH Zurich

GLEB.NOVIKOV@INF.ETHZ.CH

Editors: Gergely Neu and Lorenzo Rosasco

Abstract

In the Wishart model for sparse PCA we are given n samples $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ drawn independently from a d -dimensional Gaussian distribution $N(0, \text{Id} + \beta vv^\top)$, where $\beta > 0$ and $v \in \mathbb{R}^d$ is a k -sparse unit vector, and we wish to recover v (up to sign).

We show that if $n \geq \Omega(d)$, then for every $t \ll k$ there exists an algorithm running in time $n \cdot d^{O(t)}$ that solves this problem as long as

$$\beta \gtrsim \frac{k}{\sqrt{nt}} \sqrt{\ln(2 + td/k^2)}.$$

Prior to this work, the best polynomial time algorithm in the regime $k \approx \sqrt{d}$, called *Covariance Thresholding* (proposed in Krauthgamer et al. (2015) and analyzed in Deshpande and Montanari (2014)), required $\beta \gtrsim \frac{k}{\sqrt{n}} \sqrt{\ln(2 + d/k^2)}$. For large enough constant t our algorithm runs in polynomial time and has better guarantees than Covariance Thresholding. Previously known algorithms with such guarantees required quasi-polynomial time $d^{O(\log d)}$.

Our idea is based on the idea of Alon et al. (1998) for reducing the clique size in the planted clique problem. Moreover, we show that it is possible to combine our techniques with recent results on sparse PCA with symmetric heavy-tailed noise d’Orsi et al. (2022). Their model generalizes both sparse PCA and the planted clique problem. In particular, in the regime $k \approx \sqrt{d}$ we get the first polynomial time algorithm that works with symmetric heavy-tailed noise, while the algorithm from d’Orsi et al. (2022) requires quasi-polynomial time in these settings. As a consequence, we get an algorithm that solves a problem that captures both sparse PCA and planted clique and achieves best known guarantees for both of them.

In addition, we show that our techniques work with sparse PCA with adversarial perturbations studied in d’Orsi et al. (2020). This model generalizes not only sparse PCA, but also the sparse planted vector problem. As a consequence, we provide polynomial time algorithms for the sparse planted vector problem that have better guarantees than the state of the art in some regimes.

Keywords: Sparse PCA, Adversarial Perturbations, Semidefinite Programming, Symmetric Noise

1. Introduction

We study sparse principal component analysis in the *Wishart* and *Wigner* models. First we describe the Wishart model (that is sometimes also called the *spiked covariance model*). In this model, we are given n samples $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ drawn¹ independently from d -dimensional Gaussian distribution $N(0, \text{Id} + \beta vv^\top)$, where $\beta > 0$ and $v \in \mathbb{R}^d$ is a k -sparse² unit vector. The goal is to compute an estimator \hat{v} such that $\|\hat{v}\| = 1$ and $|\langle \hat{v}, v \rangle|$ is close to 1 (say, is greater than 0.99) with high probability³. In this paper we mostly focus on the regime when the number of samples n is greater than the dimension d , and in this section of the paper we always assume that⁴ $n \geq \Omega(d)$ (unless stated otherwise).

1. We use boldface to denote random variables.

2. That is, this vector has at most k non-zero coordinates.

3. It is impossible to recover the sign of v from $\mathbf{Y}_1, \dots, \mathbf{Y}_n$.

4. We hide absolute constant multiplicative factors using the standard notations $O(\cdot), \Omega(\cdot), \lesssim, \gtrsim$.

Classical settings. The standard approach in covariance estimation is to consider the empirical covariance $\frac{1}{n}\mathbf{Y}^\top\mathbf{Y}$ (where \mathbf{Y} is the matrix with rows $\mathbf{Y}_1, \dots, \mathbf{Y}_n$). The top eigenvector of $\frac{1}{n}\mathbf{Y}^\top\mathbf{Y}$ is highly correlated with v or $-v$ as long as $\beta \gtrsim \sqrt{\frac{d}{n}}$, and in non-sparse settings ($k = d$) these guarantees are information theoretically optimal. For $k < d$, there exists an estimator with better guarantees. It uses exhaustive search over all $\binom{d}{k}$ candidates for the support of v and is close to v or $-v$ iff $\beta \gtrsim \sqrt{\frac{k \log(de/k)}{n}}$, and these guarantees are information theoretically optimal in sparse settings (Amini and Wainwright, 2009; Berthet and Rigollet, 2013b,c).

As was observed in Johnstone and Lu (2009), known algorithmic guarantees for sparse PCA are strictly worse than the statistical guarantees described above. In the regime $k \gg \sqrt{d}$, no polynomial time algorithm is known to work if $\beta \lesssim \sqrt{\frac{d}{n}}$ (recall that if $\beta \gtrsim \sqrt{\frac{d}{n}}$, the top eigenvector of the empirical covariance is a good estimator). Johnstone and Lu (2009) proposed a polynomial time algorithm (called *Diagonal Thresholding*) that finds an estimator that is close to v or $-v$ as long as $\beta \gtrsim k\sqrt{\frac{\log d}{n}}$, which is better than the top eigenvector of $\mathbf{Y}^\top\mathbf{Y}$ if $k \ll \sqrt{d}$, but is worse than the information-theoretically optimal estimator by a factor \sqrt{k} .

Later many computational lower bounds of different kind appeared: reductions from the planted clique problem (Berthet and Rigollet, 2013a,b; Wang et al., 2016; Gao et al., 2017; Brennan et al., 2018; Brennan and Bresler, 2019), low degree polynomial lower bounds Ding et al. (2019); d’Orsi et al. (2020), statistical query lower bounds (Brennan et al., 2021), SDP and sum-of-squares lower bounds (Krauthgamer et al., 2015; Ma and Wigderson, 2015; Potechin and Rajendran, 2022), lower bounds for Markov chain Monte Carlo methods (Arous et al., 2020). These lower bounds suggest that the algorithms described above should have optimal guarantees in the regimes $k \ll \sqrt{d}$ (Diagonal Thresholding) and $k \gg \sqrt{d}$ (the top eigenvector), so it is unlikely that there exist efficient algorithms with significantly better guarantees if $k \ll \sqrt{d}$ or $k \gg \sqrt{d}$.

The regime $k \approx \sqrt{d}$ is more interesting. For a long time no efficiently computable estimator with provable guarantees better than the top eigenvector of $\mathbf{Y}^\top\mathbf{Y}$ or than Diagonal Thresholding was known, until Deshpande and Montanari (2014) proved that a polynomial time algorithm (called *Covariance Thresholding*) computes an estimator that is close to v or $-v$ as long as $\beta \gtrsim k\sqrt{\frac{\log(2+d/k^2)}{n}}$. This estimator can exploit sparsity if $k < \sqrt{d}$ and is better than Diagonal Thresholding and the top eigenvector of the empirical covariance in the regime $d^{1/2-o(1)} < k < \sqrt{d}$.

These results show that in order to work with smaller signal strength β , one needs either to work with larger number of samples n , or to work with a sparser vector v (i.e. smaller k). Ding et al. (2019) (and independently Holtzman et al. (2020)) showed that in some regimes there is another option: one can (smoothly) increase the running time needed to compute the estimator in order to work with smaller signal strength. Concretely, they showed that for $1 \leq t \leq k/\log d$ there exists an estimator that can be computed in time $d^{O(t)}$ (via *limited brute force*) and is close to v or $-v$ as long as $\beta \gtrsim k\sqrt{\frac{\log d}{tn}}$. The following example illustrates their result: For some $n, d, k \in \mathbb{N}$, let β_{DT} be the smallest signal strength such that Diagonal Thresholding, given an instance \mathbf{Y} of sparse PCA with n samples, dimension d , sparsity k and signal strength β_{DT} , finds a unit vector \hat{v}_{DT} such that $|\langle \hat{v}_{\text{DT}}, v \rangle| \geq 0.99$ with high probability. Now suppose that for the same n, d, k , we are given an instance \mathbf{Y}' of sparse PCA with smaller signal strength $\beta_{\text{new}} = 0.01 \cdot \beta_{\text{DT}}$. Then their result implies that there exists a polynomial time algorithm that, given \mathbf{Y}' , finds a unit vector \hat{v}_{new} such that $|\langle \hat{v}_{\text{new}}, v \rangle| \geq 0.99$ with high probability.

However, the approach of [Ding et al. \(2019\)](#) and [Holtzman et al. \(2020\)](#) is not compatible with the optimal guarantees in the regime $k \approx \sqrt{d}$. More precisely, if we define β_{CT} as the smallest signal strength for Covariance Thresholding (in the same way as we defined β_{DT} for Diagonal Thresholding), then the limited brute force that works with signal strength $0.01 \cdot \beta_{\text{CT}}$ requires t to be at least $\log d$ and hence runs in quasi-polynomial time $d^{O(\log d)}$. Prior to this work it was the fastest algorithm in this regime.

Our result shows that it is possible to smoothly increase running time in order to work with smaller signal strength as long as $k \leq O(\sqrt{d})$ and $n \geq \Omega(d)$. It can be informally described as follows: Let \mathcal{A} be an arbitrary currently known polynomial time algorithm for sparse PCA. Let $\beta_{\mathcal{A}}$ be the smallest signal strength such that \mathcal{A} , given an instance \mathbf{Y} of sparse PCA with, dimension d , $n \geq \Omega(d)$ samples, sparsity $k \leq O(\sqrt{d})$, and signal strength $\beta_{\mathcal{A}}$, finds a unit vector $\hat{v}_{\mathcal{A}}$ such that $|\langle \hat{v}_{\mathcal{A}}, v \rangle| \geq 0.99$ with high probability. For arbitrary constant $C \geq 1$, let $\beta_C = \frac{1}{C} \beta_{\mathcal{A}}$. Then there exists a polynomial time⁵ algorithm, that, given an instance of sparse PCA \mathbf{Y}' with signal strength β_C and the same parameters n, d, k as for \mathbf{Y} , finds a unit vector \hat{v}_{new} such that $|\langle \hat{v}_{\text{new}}, v \rangle| \geq 0.99$ with high probability.

In particular, our result implies that there exists a polynomial time algorithm that works with signal strength $0.01 \cdot \beta_{\text{CT}}$, which is a significant improvement compared to the best previously known (quasi-polynomial time) algorithm. Moreover, our result also implies the first polynomial time algorithm that can exploit sparsity and has better guarantees than the top eigenvector even in the regime $k \geq \sqrt{d}$ (as long as $k \leq O(\sqrt{d})$).

Semidefinite programming and adversarial perturbations. [d’Aspremont et al. \(2004\)](#) introduced *basic SDP* for sparse PCA. basic SDP achieves guarantees of both the top eigenvector of the empirical covariance and Diagonal Thresholding. Later [d’Orsi et al. \(2020\)](#) proved that it also achieves the guarantees of Covariance Thresholding, and hence captures the best currently known polynomial time guarantees.

Moreover, [d’Orsi et al. \(2020\)](#) showed that basic SDP also works with adversarial perturbations. More precisely, if a small (adversarially chosen) value E_{ij} is added to every entry \mathbf{Y}_{ij} of an instance of sparse PCA, basic SDP still recovers v or $-v$ with high probability. Known estimators that are not based on semidefinite programming, including top eigenvector of the empirical covariance, Diagonal Thresholding, Covariance Thresholding and limited brute force, do not work with adversarial perturbations. [d’Orsi et al. \(2020\)](#) also provided a family of algorithms based on sum-of-squares relaxations that work with adversarial perturbations and achieves the guarantees of limited brute force from [Ding et al. \(2019\)](#).

Similar to non-adversarial case, basic SDP and limited brute force based on sum-of-squares are not compatible with each other: in the regime $k \approx \sqrt{d}$, the sum-of-squares approach from [d’Orsi et al. \(2020\)](#) requires degree $\log d$ in order to achieve better guarantees than basic SDP, so the corresponding estimator can be computed only in quasi-polynomial time.

We show that our technique also works with adversarial perturbations. We remark that we do not use higher degree sum-of-squares, but only basic SDP for sparse PCA (with some preprocessing and postprocessing steps).

One of the applications of our result is an improvement in the planted sparse vector problem. For this problem we focus on the regime $\Omega(d) < n < d$. In this problem, we are given an n -dimensional

5. The degree of the polynomial depends on C .

subspace of \mathbb{R}^d that contains a sparse vector, and the goal is to estimate this vector. This problem was extensively studied in literature in different settings (Hand and Demanet, 2013; Barak et al., 2014; Hopkins et al., 2016; Qu et al., 2020; Mao and Wein, 2021; Zadik et al., 2022; Diakonikolas and Kane, 2022). It is not hard to see⁶ that this problem in the *Gaussian basis* model (in the sense of Mao and Wein (2021)) is a special case of sparse PCA with small perturbations. Our result shows that as long as $k \leq \sqrt{td}$, there exists a $d^{O(t)}$ time algorithm for this problem. Previously known polynomial time algorithms in the regime $n \geq \Omega(d)$ required $k \leq C\sqrt{d}$ for some absolute constant C and did not work for $k > C\sqrt{d}$. Lower bounds against restricted computational models (d’Orsi et al., 2020; Ding et al., 2021; Ding and Hua, 2023) suggest that in the regime $n \geq \Omega(d)$ this problem is unlikely to be solvable in polynomial time if $k \gg \sqrt{d}$.

The Wigner model and symmetric noise. Our results can be naturally applied also to the Wigner model. In this model, we are given $\mathbf{Y} = \lambda vv^\top + \mathbf{W}$, where $\lambda > 0$, $v \in \mathbb{R}^d$ is a k -sparse unit vector, and $\mathbf{W} \sim N(0, 1)^{d \times d}$. For this model, Covariance Thresholding finds an estimator highly correlated with v or $-v$ as long as $\lambda \gtrsim k\sqrt{\log(2 + d/k^2)}$, while limited brute force from Ding et al. (2019) computes in time $d^{O(t)}$ an estimator close to v or $-v$ as long as $\lambda \gtrsim k\sqrt{\frac{\log d}{t}}$. As in the Wishart model, these approaches are not compatible in the regime $k \approx \sqrt{d}$. Our techniques can be naturally applied to the Wigner model, leading to the best known algorithms for this problem.

As in the Wishart model, our techniques also work with adversarial perturbations. Moreover, our approach is compatible with the recent study of Sparse PCA with symmetric noise d’Orsi et al. (2022). In this model, Gaussian noise \mathbf{W} is replaced by an arbitrary noise N with symmetric about zero independent entries that are only guaranteed to be bounded by 1 with probability⁷ $\Omega(1)$. They proposed a quasi-polynomial algorithm for sparse PCA in these settings and provided evidence that in the regime $k \ll \sqrt{d}$ this running time cannot be improved (via reduction from the planted clique problem). Combining their algorithm with our approach, we show that in the regime $k \approx \sqrt{d}$ there exists a polynomial time algorithm that solves this problem.

1.1. Results

Before stating the results, observe that one can write an instance \mathbf{Y} sparse PCA in the Wishart model as $\mathbf{Y} = \sqrt{\beta}\mathbf{u}v^\top + \mathbf{W}$, where $\mathbf{u} \sim N(0, 1)^n$ and $\mathbf{W} \sim N(0, 1)^{n \times d}$ are independent.

Classical settings. Our first result is estimating v in the Wishart model in classical settings (without perturbations).

Theorem 1 (The Wishart model) *Let $n, d, k, t \in \mathbb{N}$, $\beta > 0$. Let $\mathbf{Y} = \sqrt{\beta}\mathbf{u}v^\top + \mathbf{W}$, where $\mathbf{u} \sim N(0, 1)^n$, $v \in \mathbb{R}^d$ is a k -sparse unit vector, $\mathbf{W} \sim N(0, 1)^{n \times d}$ independent of \mathbf{u} .*

There exists an absolute constant $C > 1$, such that if $n \geq Ck$, $k \geq Ct \log^2 d$ and

$$\beta \geq C \frac{k}{\sqrt{tn}} \sqrt{\log \left(2 + \frac{td}{k^2} \left(1 + \frac{d}{n} \right) \right)},$$

then there exists an algorithm that, given \mathbf{Y} , k and t , in time $n \cdot d^{O(t)}$ outputs a unit vector \hat{v} such that with probability $1 - o(1)$ as $d \rightarrow \infty$,

$$|\langle \hat{v}, v \rangle| \geq 0.99.$$

6. See the discussion before corollary 3.

7. Note that even the first moment is not required to exist.

Let us compare our guarantees with previously known estimators. For simplicity we assume $n \geq \Omega(d)$. For this regime, best estimators known prior to this work and their guarantees are listed in Table 1.

Estimator	Signal Strength	Time Complexity
Statistically optimal estimator	$\beta \gtrsim \sqrt{\frac{k}{n} \log(ed/k)}$	$n \cdot d^{O(k)}$
Top eigenvector of the empirical covariance	$\beta \gtrsim \sqrt{\frac{d}{n}}$	$n \cdot d^{O(1)}$
Covariance Thresholding	$\beta \gtrsim \frac{k}{\sqrt{n}} \sqrt{\log(2 + d/k^2)}$	$n \cdot d^{O(1)}$
Limited brute force from Ding et al. (2019)	$\beta \gtrsim \frac{k}{\sqrt{tn}} \sqrt{\log d}$	$n \cdot d^{O(t)}$
Our estimator	$\beta \gtrsim \frac{k}{\sqrt{tn}} \sqrt{\log(2 + td/k^2)}$	$n \cdot d^{O(t)}$

Table 1: Estimators for sparse PCA in the Wishart model (assuming $n \geq \Omega(d)$ and $t \leq k/\text{polylog}(d)$).

For $k \leq d^{1/2-\Omega(1)}$ (say, $k \leq d^{0.49}$), the guarantees of the algorithm from [Ding et al. \(2019\)](#) are similar to ours (up to a constant factor). For $k \geq d^{1/2-o(1)}$ our algorithm can work with asymptotically smaller signal strength (with the same running time).

To compare with Covariance Thresholding and the top eigenvector of the empirical covariance, consider the regime $k = \Theta(\sqrt{d})$. Note in this regime both Covariance Thresholding and the top eigenvector require

$$\beta \geq c\sqrt{d/n}$$

for some specific constant c (that depends on \sqrt{d}/k), and they do not work for smaller β . Our condition on β in these settings is

$$\beta \gtrsim \frac{k}{\sqrt{tn}} \sqrt{\log t},$$

so if $\beta = \varepsilon\sqrt{d/n}$ for arbitrary constant ε , we can choose large enough constant t such that $\varepsilon\sqrt{d} \gtrsim k\sqrt{\frac{\log t}{t}}$ and get an estimator that is highly correlated with v or $-v$ in polynomial time $n \cdot d^{O(t)}$. Neither Covariance Thresholding nor the top eigenvector of the empirical covariance can work with small values of ε , and limited brute force from [Ding et al. \(2019\)](#) requires quasi-polynomial time $n \cdot d^{O(\log d)}$ in these settings.

It is also interesting to compare our upper bound with the low degree polynomial lower bound from [d’Orsi et al. \(2020\)](#). They showed that in the regime $k \leq O(\sqrt{d})$, polynomials of degree $D \leq n/\log^2 n$ cannot distinguish⁸ $\mathbf{Y}_1, \dots, \mathbf{Y}_n \sim N(0, \text{Id} + \beta vv^\top)$ from $\mathbf{Y}_1, \dots, \mathbf{Y}_n \sim N(0, \text{Id})$ if

$$\beta \lesssim \frac{k}{\sqrt{Dn}} \cdot \log\left(2 + \frac{Dd}{k^2}\right).$$

8. More precisely, they cannot *strongly distinguish* sequences of distributions in the sense of [Kunisky et al. \(2019\)](#).

and hence for such β they cannot be used to design an estimator that is close to v or $-v$ with high probability. Their lower bound does not formally imply that our upper bound is tight (that is, it does not imply that there are no better estimators than ours among low degree polynomials). However, there is an interesting similarity between the lower bound and the upper bound: They have a very similar logarithmic factor. If this similarity can be formalized, it may lead to an algorithm that works in a small sample regime $n \ll d$. The term d/n that we have in the logarithmic factor in the bound on β is necessary for our techniques. Many other algorithms, like basic SDP or Covariance Thresholding, also have similar terms. However, the low-degree lower bound does not have this term and [d’Orsi et al. \(2020\)](#) provided an algorithm based on low degree polynomials that does not have such a term and works as long as $\beta \gtrsim \frac{k}{\sqrt{n}} \sqrt{\log(2 + td/k^2)}$ even for very small n (e.g. $n = d^{0.01}$). Finding an estimator with guarantees similar to ours in the small sample regime $n \ll d$ is an interesting open question, and low degree polynomials might be useful in designing such an estimator.

Adversarial perturbations. Our approach also works in the presence of adversarial perturbations.

Theorem 2 (The Wishart model with adversarial perturbations) *Let $n, d, k, t \in \mathbb{N}$, $\beta > 0$, $\varepsilon \in (0, 1)$. Let $Y = \sqrt{\beta} \mathbf{u} v^\top + \mathbf{W} + E$, where $\mathbf{u} \sim N(0, 1)^n$, $v \in \mathbb{R}^d$ is a k -sparse unit vector, $\mathbf{W} \sim N(0, 1)^{n \times d}$ independent of \mathbf{u} and $E \in \mathbb{R}^{n \times d}$ is a matrix such that*

$$\|E\|_{1 \rightarrow 2} \leq \varepsilon \cdot \min\{\sqrt{\beta}, \beta\} \cdot \sqrt{n/k},$$

where $\|E\|_{1 \rightarrow 2}$ is the maximal norm of the columns of E and $\varepsilon < 1$.

There exists an absolute constant $C > 1$, such that if $n \geq Ck$, $k \geq Ct \log^2 d$,

$$\beta \geq C \frac{k}{\sqrt{tn}} \sqrt{\log\left(2 + \frac{td}{k^2} \left(1 + \frac{d}{n}\right)\right)}.$$

and $\varepsilon \sqrt{\log(1/\varepsilon)} \leq \frac{1}{C} \min\left\{1, \min\{\beta, \sqrt{\beta}\} \cdot \sqrt{n/d}\right\}$, then there exists an algorithm that, given Y , k and t , in time $n \cdot d^{O(t)}$ outputs a unit vector \hat{v} such that with probability $1 - o(1)$ as $d \rightarrow \infty$,

$$|\langle \hat{v}, v \rangle| \geq 0.99.$$

To illustrate how large the adversarial perturbations are allowed to be, consider the following example: Let $k = \Theta(\sqrt{d})$, $n = \Theta(d)$, $\beta = \Theta(1)$, $t \leq O(1)$. Then the columns of E can have norm as large as $\Omega(\sqrt{k})$. Note that in these settings column norms of $\sqrt{\beta} \mathbf{u} v^\top$ can be $O(\sqrt{k})$. Hence, in this regime, if we allow E to be larger by a constant factor, the adversary can choose $E = -\sqrt{\beta} \mathbf{u} v^\top$ and erase the signal. As was shown in [d’Orsi et al. \(2020\)](#), in these settings Covariance Thresholding, Diagonal Thresholding and the top eigenvector of the empirical covariance do not work with some perturbations E such that $\|E\|_{1 \rightarrow 2} \leq k^{o(1)}$.

Our assumption on E is stronger than the assumption from [d’Orsi et al. \(2020\)](#), which is $\|E\|_{1 \rightarrow 2} \lesssim \min\{\beta, \sqrt{\beta}\} \sqrt{n/k}$. Designing an estimator with guarantees similar to ours that works with larger E is an interesting problem.

Similar to the non-adversarial settings, our algorithms have the same guarantees⁹ as the sum-of-squares approach from [d’Orsi et al. \(2020\)](#) if $k \leq d^{1/2 - \Omega(1)}$ and has asymptotically better guarantees

9. Assuming our bound on the columns of E .

if $k \geq d^{1/2-o(1)}$. Similarly to Covariance Thresholding in the non-adversarial case, in the regime $n \geq \Omega(d)$ and $k = \Theta(\sqrt{d})$ basic SDP requires $\beta \geq c\sqrt{d/n}$ for some specific constant c and does not work for smaller β . Our condition on β in these settings is $\beta \gtrsim \frac{k}{\sqrt{tn}}\sqrt{\log t}$, so if $\beta = \varepsilon\sqrt{d/n}$ for arbitrary constant ε , we can choose large enough constant t such that $\varepsilon\sqrt{d} \gtrsim k\sqrt{\frac{\log t}{t}}$ and get an estimator that can be computed in polynomial time $n \cdot d^{O(t)}$. basic SDP cannot work with small values of ε , and sum-of-squares approach from d’Orsi et al. (2020) requires quasi-polynomial time in these settings.

The sparse planted vector problem. As was observed in d’Orsi et al. (2020), sparse PCA with perturbations is a generalization not only for the spiked covariance model, but also for the planted sparse vector problem. In this problem we are given an n -dimensional subspace of \mathbb{R}^d spanned by $n - 1$ random vectors and a sparse vector, and the goal is to find the sparse vector. More precisely, let $\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_n$ be standard d -dimensional Gaussian vectors and let \mathbf{B} be an $n \times d$ matrix whose first $n - 1$ rows are $\mathbf{g}_1^\top, \dots, \mathbf{g}_{n-1}^\top$ and the last row is a vector $\|\mathbf{g}_n\|v^\top$, where $v \in \mathbb{R}^d$ is k -sparse and unit. Let \mathbf{R} be a random rotation of \mathbb{R}^n independent of $\mathbf{g}_1 \dots, \mathbf{g}_n$, and let $\mathbf{Y} = \mathbf{R}\mathbf{B}$. The goal is to recover v from \mathbf{Y} .

This problem can be seen as a special case of sparse PCA with perturbation matrix $E = -\frac{1}{\|\mathbf{u}\|^2}\mathbf{u}\mathbf{u}^\top\mathbf{W}$ (see section B for the proof). Therefore, we can apply theorem 2 and get

Corollary 3 (The sparse planted vector problem) *Let $n, d, k, t \in \mathbb{N}$, $\beta > 0$. Let $\mathbf{Y} = \sqrt{\beta}\mathbf{u}v^\top + \mathbf{W} - \frac{1}{\|\mathbf{u}\|^2}\mathbf{u}\mathbf{u}^\top\mathbf{W}$, where $\mathbf{u} \sim N(0, 1)^n$, $v \in \mathbb{R}^d$ is a k -sparse unit vector, $\mathbf{W} \sim N(0, 1)^{n \times d}$ independent of \mathbf{u} , and $\sqrt{\beta} = \frac{\|\mathbf{u}^\top\mathbf{W}\|}{\|\mathbf{u}\|^2}$.*

There exists an absolute constant $C > 1$, such that if $d > n$, $n \geq Ck$, $k \geq Ct \log^2 d$ and

$$k \leq \frac{1}{C} \cdot d\sqrt{t/n},$$

then there exists an algorithm that, given \mathbf{Y} , k and t , in time $d^{O(t)}$ outputs a unit vector \hat{v} such that with probability $1 - o(1)$ as $d \rightarrow \infty$,

$$|\langle \hat{v}, v \rangle| \geq 0.99.$$

Prior to this work, in the regime $n \geq \Omega(d)$, polynomial time estimators were known only if $k \leq cd/\sqrt{n}$ for some small constant $c < 1$ (the existence of such an algorithm follows from Theorem 4.5 from d’Orsi et al. (2020)). We show that even if $k \geq 100d/\sqrt{n}$, sparsity can still be exploited and there are estimators that can be computed in polynomial time.

The Wigner model and symmetric noise. Our techniques also work with sparse PCA in the Wigner model.

Theorem 4 (The Wigner model) *Let $k, d, t \in \mathbb{N}$, $\lambda > 0$. Let $Y = \lambda vv^\top + \mathbf{W} + E$, where $v \in \mathbb{R}^d$ is a k -sparse unit vector and $\mathbf{W} \sim N(0, 1)^{d \times d}$ and $E \in \mathbb{R}^{d \times d}$.*

There exists an absolute constant $C > 1$, such that if $k \geq Ct \log d$, $\|E\|_\infty \leq \frac{1}{C}\lambda/k$, and

$$\lambda \geq Ck\sqrt{\frac{\log(2 + td/k^2)}{t}},$$

then there exists an algorithm that, given Y , k and t , in time $d^{O(t)}$ outputs a unit vector \hat{v} such that with probability $1 - o(1)$ as $d \rightarrow \infty$,

$$|\langle \hat{v}, v \rangle| \geq 0.99.$$

Note that E is allowed to be as large as possible (up to a constant factor). Similar to the Wishart model, previously known polynomial time algorithms required

$$\lambda \gtrsim \min \left\{ k \sqrt{\frac{\log d}{t}}, k \sqrt{\log(2 + d/k^2)}, \sqrt{d} \right\},$$

where $k \sqrt{\frac{\log d}{t}}$ corresponds to the limited brute force from [Ding et al. \(2019\)](#), and $\min \left\{ k \sqrt{\log(2 + d/k^2)}, \sqrt{d} \right\}$ corresponds to basic SDP. Similarly to the Wishart model, in the regime $k \geq d^{1/2 - o(1)}$ we get asymptotically better guarantees than the algorithm from [Ding et al. \(2019\)](#). In the regime $n \geq \Omega(d)$ and $k = \Theta(\sqrt{d})$ basic SDP requires $\lambda \geq c\sqrt{d}$ for some specific constant c and does not work for smaller λ . Our condition on λ in these settings is $\lambda \gtrsim k \sqrt{\frac{\log t}{t}}$, so if $\lambda = \varepsilon \sqrt{d}$ for arbitrary constant ε , we can choose large enough constant t such that $\varepsilon \sqrt{d} \gtrsim k \sqrt{\frac{\log t}{t}}$ and get an estimator that can be computed in polynomial time $d^{O(t)}$. basic SDP cannot work with small values of ε , and limited brute force requires quasi-polynomial time in these settings.

Our techniques can be also applied to more general model with symmetric noise studied in [d'Orsi et al. \(2022\)](#).

Theorem 5 (Sparse PCA with symmetric heavy-tailed noise) *Let $k, d, t \in \mathbb{N}$, $\lambda > 0$. Let $Y = \lambda v v^\top + N$, where $v \in \mathbb{R}^d$ is a k -sparse unit vector such that $\|v\|_\infty \leq 100/\sqrt{k}$ and N is a random matrix with independent (but not necessarily identically distributed) symmetric about zero entries¹⁰ such that for all $i, j \in [d]$, $\mathbb{P}[|N_{ij}| \leq 1] \geq 0.1$.*

There exists an absolute constant $C > 1$, such that if $k \geq Ct \log d$,

$$t \geq C \cdot \log(2 + d/k^2),$$

and

$$\lambda \geq k,$$

then there exists an algorithm that, given Y , k , t and λ , in time $d^{O(t)}$ outputs a unit vector \hat{v} such that with probability $1 - o(1)$ as $d \rightarrow \infty$,

$$|\langle \hat{v}, v \rangle| \geq 0.99.$$

Moreover, we get the same guarantees if we are given only the upper triangle of Y , i.e. the entries Y_{ij} such that $i < j$.

If $k = \Theta(\sqrt{d})$, this algorithm runs in polynomial time as long as $\lambda \geq \varepsilon \sqrt{d}$ for (arbitrary) constant ε . It is the first known polynomial time algorithm for this model, since the algorithm from

10. That is, N_{ij} and $-N_{ij}$ have the same distribution.

d’Orsi et al. (2022) requires quasi-polynomial time. For example, our algorithm finds an estimator close to v or $-v$ in polynomial time if $k = \sqrt{d}$ and $\lambda = k/100$ when the noise has iid Cauchy entries¹¹ (with location 0 and scale 1), while prior to this work the fastest known algorithm in this regime required quasi-polynomial time even for standard Gaussian noise.

In the special case of Gaussian noise and $\lambda = k$, this algorithm matches the best known algorithmic guarantees. Note, however, that in the regime $k \leq d^{1/2-\Omega(1)}$ this algorithm runs in quasipolynomial time. This is not surprising, since as was observed in d’Orsi et al. (2020), sparse PCA with symmetric noise actually generalizes the planted clique problem.

More precisely, let $G \sim G(d, 1/2, k)$ be a random graph with a planted clique of size k . Let A be the adjacency matrix of G . Let J be the matrix with all entries equal to 1 and let $C = 2A - J$. Note that the upper triangle of C coincides with the upper triangle of $k \cdot vv^\top + \eta$, where $\sqrt{k} \cdot v$ is the indicator vector of the vertices of the clique (so it is k -sparse), and η is the noise whose entries that correspond to the vertices of the clique are zero, and other entries are iid uniform over $\{\pm 1\}$.

The algorithm from theorem 5 solves the planted clique problem in time $n^{O(\log(2+n/k^2))}$, which matches the best known algorithmic guarantees for the planted clique¹². Moreover, for some $k = n^{\Omega(1)}$ it is conjectured to be impossible to solve it in time $n^{o(\log n)}$ (see Manurangsi et al. (2021) for more details). Note that our algorithm achieves best known algorithmic guarantees for both sparse PCA in the Wigner model and the planted clique problem¹³.

2. Techniques

The idea of our approach is similar to the well-known technique of reducing the constant in the planted clique problem. Recall that an instance of the planted clique problem is a random graph G sampled according to the following distribution: First, a graph is sampled from Erdős-Rényi distribution $\mathcal{G}(m, 1/2)$ (i.e. each pair of vertices is chosen independently to be an edge with probability $1/2$), and then a random subset of vertices of size k is chosen (uniformly from the sets of size k and independently from the graph) and the clique corresponding to these vertices is added to the graph. The goal is to find the clique. The problem can be solved in quasi-polynomial time, however, no polynomial time algorithm is known in the regime $k \leq o(\sqrt{m})$.

Alon et al. (1998) proposed a spectral algorithm that can be used to find the clique in polynomial time if $k \gtrsim \sqrt{m}$. They also introduced a technique that allows to find the clique in polynomial time if $k \geq \varepsilon\sqrt{m}$ for arbitrary constant $\varepsilon > 0$. The idea is to look at every subset T of vertices of size $t \gtrsim \log(1/\varepsilon)$ and consider the subgraph $H(T)$ induced by the vertices of G that are adjacent to T (i.e. adjacent to every vertex of T). This subgraph has approximately $m' = 2^{-t}m \lesssim \varepsilon^2 m$ vertices, and if T was a part of the clique, then the clique is preserved in $H(T)$, and since $k \gtrsim \sqrt{m'}$, we can find a clique applying the spectral algorithm to $H(T)$. The running time of the algorithm is $m^{O(t)}$, so it is polynomial for constant ε .

A similar (but technically more challenging) idea can be also used for sparse PCA. Recall that the instance of sparse PCA (in the Wishart model) is $Y = \sqrt{\beta}uv^\top + W$, where $u \sim N(0, 1)^n$, $v \in \mathbb{R}^d$ is a k -sparse unit vector, $W \sim N(0, 1)^{n \times d}$ independent of u . To illustrate the idea, we assume that v is flat, i.e. its nonzero entries are $\pm 1/\sqrt{k}$. Instead of the adjacency matrix of the graph, we have the empirical covariance $\frac{1}{n}Y^\top Y$. For simplicity, let us ignore cross terms and

11. Cauchy noise is very heavy-tailed, the entries do not even have a finite first moment.

12. Up to a constant factor in the degree.

13. Assuming $\lambda = k$ for sparse PCA.

assume that

$$\frac{1}{n} \mathbf{Y}^\top \mathbf{Y} \approx \frac{\|\mathbf{u}\|^2}{n} \beta v v^\top + \frac{1}{n} \mathbf{W}^\top \mathbf{W}.$$

Since $\mathbf{u} \sim N(0, \text{Id})$, $\|\mathbf{u}\|^2 \approx n$. So we assume that we are given $\beta v v^\top + \frac{1}{n} \mathbf{W}^\top \mathbf{W}$, and the goal is to recover v . Similar to the planted clique, if $\beta \gtrsim \sqrt{d/n}$, there is a spectral algorithm for this problem (that computes the top eigenvector of the empirical covariance).

Suppose that $n \geq d$, $k \approx \sqrt{d}$ and $\beta = \varepsilon \sqrt{d/n} \approx \varepsilon k / \sqrt{n}$ for some constant $\varepsilon > 0$. We can look at each subset T of entries of size $t < k$ and try to find a (principal) submatrix of $\frac{1}{n} \mathbf{Y}^\top \mathbf{Y}$ that would be an analogue of the graph $\mathbf{H}(T)$ from the algorithm for the planted clique. One option is to say that an entry i is “adjacent” to T if the sum of the elements in the i -th row of $\frac{1}{n} \mathbf{Y}^\top \mathbf{Y}$ is large. However, since v has both positive and negative entries, the sum can be small even if here were no noise and $T \subset \text{supp}(v)$. Hence we also need to take the signs of the entries of v into account.

Let \mathcal{S}_t be the set of all t -sparse vectors with entries from $\{0, \pm 1\}$. Let us call $s \in \mathcal{S}_t$ *correct* if $\text{supp}(s) \subset \text{supp}(v)$ and for all nonzero s_i , $\text{sign}(s_i) = \text{sign}(v_i)$. If $s \in \mathcal{S}_t$ is correct, then

$$\left| \sum_{j \in \text{supp}(s)} \beta v_i v_j s_j \right| = \beta |v_i v^\top s| = \beta t / k.$$

For $s \in \mathcal{S}$ let us call an entry $i \in [d]$ *adjacent* to s if either $|(\frac{1}{n} \mathbf{Y}^\top \mathbf{Y} s)_i| \geq \beta t / (2k)$ or $i \in \text{supp}(s)$, and let $\mathbf{H}(s)$ be a principal submatrix induced by indices adjacent to s . The size of $\mathbf{H}(s)$ is close to pd , where p is the probability that $|(\frac{1}{n} \mathbf{W}^\top \mathbf{W} s)_i|$ is greater than $\beta t / (2k)$. We need to count $i \notin \text{supp}(s)$ adjacent to s . The vector $\mathbf{W} s$ has distribution $N(0, t \cdot \text{Id})$, hence $\|\mathbf{W} s\| \approx \sqrt{tn}$. Since $i \notin \text{supp}(s)$, the i -th row of \mathbf{W}^\top is independent of $\mathbf{W} s$, and the distribution of $(\frac{1}{n} \mathbf{W}^\top \mathbf{W} s)_i$ is close to $N(0, t/n)$. By the tail bound for the Gaussian distribution,

$$\mathbb{P} \left[\left| \left(\frac{1}{n} \mathbf{W}^\top \mathbf{W} s \right)_i \right| \geq x \sqrt{t/n} \right] \leq \exp(-x^2/2).$$

In our case, $x = \frac{\beta \sqrt{tn}}{2k}$. Hence for

$$\beta \gtrsim \frac{k}{\sqrt{tn}} \sqrt{\log(td/k^2)},$$

we get $x \gtrsim \sqrt{\log(td/k^2)}$ and $p = \exp(-x^2/2) \lesssim k^2/(td)$. Therefore, $\mathbf{H}(s)$ has $d' \lesssim k^2/t$ entries. Moreover, by the same argument, $k' \approx (1-p)k \geq 0.999k$ entries of v are adjacent to correct s , so the signal part of $\mathbf{H}(s)$ is close to $\beta v v^\top$. Since for correct s we get $\beta \gtrsim \sqrt{d'/n}$, we can try to use the spectral algorithm to recover the sparse vector from $\mathbf{H}(s)$.

Here we see the difference between planted clique and sparse PCA. In the planted clique problem, if we take a subset of the clique, we can easily recover the whole clique from the output of the spectral algorithm and we do not need to consider other sets after that. In sparse PCA, since we do not know β exactly, it might not be easy to understand if the observed s was correct or not from the output of the spectral algorithm.

We use the following observation: if we have computed the list $\mathbf{L} = \{\tilde{v}(s)\}$ of the top eigenvectors of $\mathbf{H}(s)$ for all $s \in \mathcal{S}_t$, we can compute a vector close to v (or to $-v$) from this list. Indeed, if we erase all but the largest $O(k)$ entries (in absolute value) of the vectors from \mathbf{L} , we get a new

list \mathbf{L}' of $O(k)$ -sparse vectors. It turns out that for correct s not only $\tilde{v}(s)$, but also the corresponding $O(k)$ -sparse vector $v'(s) \in \mathbf{L}'$ is close to v . Moreover, for all $O(k)$ -sparse unit vectors x (in particular, for all vectors in \mathbf{L}'),

$$x\left(\frac{1}{n}\mathbf{W}^\top\mathbf{W}\right)x = 1 \pm \tilde{O}\left(\sqrt{k/n}\right).$$

Hence we can just compute $\hat{v} \in \operatorname{argmax}_{x \in \mathbf{L}'} x\left(\frac{1}{n}\mathbf{Y}^\top\mathbf{Y}\right)x$, and it is close to v , since

$$x\left(\frac{1}{n}\mathbf{Y}^\top\mathbf{Y}\right)x \approx 1 + \beta\langle v, x \rangle^2 \pm \tilde{O}\left(\sqrt{k/n}\right),$$

and for $\beta > \frac{k}{\sqrt{tn}}$, the term $\tilde{O}\left(\sqrt{k/n}\right)$ is smaller than β (as long as $t \ll k$).

Note that in the definition of adjacent entries we used β , which might be unknown. But that is not a problem: if we use some value $\beta/2 < \beta' \leq \beta$ instead of β , the algorithm still works. Hence we can use all possible candidates from $n^{-O(1)}$ to $n^{O(1)}$ such that one of them differs from β by at most factor of 2, and in the end work not with the list \mathbf{L} , but with a list of size $O(\log n) \cdot |\mathbf{L}|$.

Remark [*Comparison with the Covariance Thresholding analysis from [Deshpande and Montanari \(2014\)](#)*] Our algorithm for $t = 1$ has running time $O(nd^2) + \tilde{O}(d^3)$. For $n \geq d$, the running time is comparable to the running time of Covariance Thresholding $O(nd^2)$. The guarantees of both algorithms are the same (up to a constant factor). One advantage of our algorithm is that it is much easier to analyze. The crucial difference is that in Covariance Thresholding one has to bound the spectral norm of thresholded Wishart matrix, which requires a sophisticated probabilistic argument. In our algorithm, we only need to bound principal submatrices of the Wishart matrix, and such bounds easily follow from concentration of the spectral norm of $(\frac{1}{n}\mathbf{W}^\top\mathbf{W} - \text{Id})$ and a union bound argument. Another advantage of our algorithms is that we can get better guarantees than Covariance Thresholding (by increasing t and hence also the running time) and use the same analysis for all t .

Remark [*Comparison with the algorithms from [Ding et al. \(2019\)](#)*] The algorithms from [Ding et al. \(2019\)](#) also use vectors $s \in \mathcal{S}_t$. However, the crucial difference between our approaches is that they work with $s' \in \mathcal{S}_t$ that maximizes $s\left(\frac{1}{n}\mathbf{Y}^\top\mathbf{Y}\right)s$. This approach works only if $\beta \gtrsim \frac{k}{\sqrt{t}}\sqrt{\log d}$, since for smaller β the maximizer of $s\left(\frac{1}{n}\mathbf{Y}^\top\mathbf{Y}\right)s$ might be completely unrelated to v . For our analysis it is not a problem, since the correct s is only determined in the end from the list \mathbf{L}' .

Adversarial Perturbations. Similar approach also works in the presence of adversarial perturbations, that is, if the input is $Y = \sqrt{\beta}\mathbf{u}v^\top + \mathbf{W} + E$ such that the columns of E have norm bounded by $b \ll \beta\sqrt{n/k}$. This is interesting, since known algorithms for sparse PCA that use thresholding techniques and are not based on semidefinite programming, like Diagonal Thresholding, Covariance Thresholding, or the algorithms from [Ding et al. \(2019\)](#), do not work in these settings (see [d’Orsi et al. \(2020\)](#) for more details).

As in the non-adversarial case, we can compute the submatrices $H(s)$ that are induced by indices adjacent to s , that is, indices i such that either $\left|(\frac{1}{n}\mathbf{Y}^\top\mathbf{Y}s)_i\right| \geq \beta t/(2k)$ or $i \in \operatorname{supp}(s)$. Then, instead of computing the top eigenvector of $H(s)$, we compute $\tilde{X}(s) \in \operatorname{argmax}_{X \in \mathcal{P}_k} \langle X, H(s) \rangle$, where

$$\mathcal{P}_k = \left\{ X \in \mathbb{R}^{d \times d} \mid X \succeq 0, \operatorname{Tr} X = 1, \|X\|_1 \leq k \right\}$$

is the feasible region of the basic SDP for sparse PCA. Then, we can compute the list \mathbf{L} of top eigenvectors $\tilde{v}(s)$ of $\tilde{X}(s)$, and perform the same procedure as in the non-adversarial case to recover \hat{v} from \mathbf{L} .

In order to show the correctness, we need to bound all terms of $\frac{1}{n}Y^\top Y s$. In adversarial settings cross terms can be large, and the most problematic term is $\frac{1}{n}E^\top \mathbf{W} s$. Rows of E^\top do not have large norm, but $\mathbf{W} s$ has large norm $\|\mathbf{W} s\| \approx \sqrt{tn}$, and since E can depend on \mathbf{W} , the entries of $E^\top \mathbf{W} s$ can be large. In particular, for each correct s , the adversary can always choose E such that the term $\frac{1}{n}E^\top \mathbf{W} s$ is large enough to make $H(s)$ useless for recovering v .

To resolve this issue, we work with some probability distribution over the set of correct s and show that $\frac{1}{n}E^\top \mathbf{W} s$ has small expectation with respect to this distribution. In particular, it implies that for each E there exists some s' such that the term $\frac{1}{n}E^\top \mathbf{W} s'$ is small¹⁴. For flat v , we just divide the support of v into $m = k/t$ blocks of size t , and then each block corresponds to some correct $s \in \mathcal{S}_t$. Then it is enough to consider uniform distribution \mathcal{U} over the set $\{s_1, \dots, s_m\}$ of such s . By the concentration of spectral norm of \mathbf{W} , with high probability

$$\frac{1}{n^2} \mathbb{E}_{s \sim \mathcal{U}} \langle E_i, \mathbf{W} s \rangle^2 = \frac{1}{n^2 m} \sum_{j=1}^m \langle E_i, \mathbf{W} s_j \rangle^2 \leq O\left(\frac{b^2 t}{n^2 m} (m+n)\right) \leq O\left(\frac{b^2 t^2}{nk}\right) \ll \frac{\beta^2 t^2}{k^2}.$$

Hence there exists some $s' \in \{s_1, \dots, s_m\}$ such that

$$\left\| \left(\frac{1}{n} E^\top \mathbf{W} s' \right)_{\text{supp}(v)} \right\|^2 \ll \|\beta v v^\top s'\|^2.$$

The other terms of $\frac{1}{n}Y^\top Y s$ can be bounded only assuming that s' is correct (so it is not needed to use properties of \mathcal{U} anymore), hence the signal part of $H(s')$ is close to $\beta v v^\top$.

In addition to pd entries $i \in [d] \setminus \text{supp}(v)$ adjacent to s' that appear due to the term $\frac{1}{n} \mathbf{W}^\top \mathbf{W} s$, there could be some entries adjacent to s' that appear from $\frac{1}{n}(E^\top Y + Y^\top E) s'$. We show that the number of such entries is at most $\varepsilon^2 \log(1/\varepsilon)d$, where ε is the same as in theorem 2. Assuming our bound¹⁵ on the maximal norm of columns of E , we get $\varepsilon^2 \log(1/\varepsilon)d \lesssim \beta n$, and by standard properties of basic SDP for sparse PCA, the top eigenvector of $\tilde{X}(s')$ is close to v .

To finish the argument, we need to show that we can still compute \hat{v} close to v or $-v$ from \mathbf{L} even in the presence of perturbations. It is not hard since our argument depends only on the upper bound on $x(\frac{1}{n}Y^\top Y - \text{Id} - \beta v v^\top)x$ for all $O(k)$ -sparse x . As was shown in d'Orsi et al. (2020), our assumption on the maximal norm of columns of E is enough to obtain the desired upper bound.

The Wigner model and symmetric noise. In the Wigner model the input is $n \mathbf{Y} = \lambda v v^\top + \mathbf{W}$. The same argument as for the Wishart model works in these settings, and the proof is technically simpler since there are no cross terms and \mathbf{W} is easier to analyze than $\frac{1}{n} \mathbf{W}^\top \mathbf{W}$ that appears in the Wishart model. Our approach for sparse PCA with perturbations also works for Wigner model, and the proof is much easier in this case since the adversary cannot exploit magnitude of the columns of \mathbf{W} .

The Wigner model with symmetric noise is more challenging. In these settings we assume that λ is known. We cannot work with $\mathbf{Y} s$, since the noise is unbounded. So we first threshold the

14. More precisely, this term has small norm, and it is enough for our analysis.

15. This is the reason why our bound on E is worse than the bound from d'Orsi et al. (2020).

entries of \mathbf{Y} . Concretely, for $h > 0$ and $x \in \mathbb{R}$, let $\tau_h(x)$ be x if $x \in [-h, h]$ and $\text{sign}(x) \cdot h$ otherwise. We apply this transformation for some $h = \Theta(\lambda/k)$ to the entries of \mathbf{Y} and get a new matrix \mathbf{T} . Then we use our approach for matrix \mathbf{T} and for all $s \in \mathcal{S}_t$ we compute the submatrices $\mathbf{H}(s)$. As long as $k \gtrsim d$, their algorithm applied to \mathbf{Y} outputs a matrix that is close to λvv^\top in polynomial time. As in the Gaussian case, the submatrices $\mathbf{H}(s)$ have small size, and we can apply their result to every $\mathbf{H}(s)$. However, since $\mathbf{H}(s)$ depends on \mathbf{Y} , the noise part of $\mathbf{H}(s)$ might not have the same distribution as \mathbf{N} . Fortunately, the error probability in d’Orsi et al. (2022) is very small, which allows us to use union bound and conclude that for correct s , the output of the algorithm from d’Orsi et al. (2022) on $\mathbf{H}(s)$ is close to λvv^\top . Moreover, since we know λ , we do not even need to work with the list of candidates in these settings: It is enough to check the norm of the output, and if it is close to λ , the output is close to λvv^\top , and we can recover v from it (up to sign).

Acknowledgments

This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No 815464).

I would like to thank David Steurer for helpful conversations. Thanks also Lucas Slot for his comments and suggestions, which greatly enhanced the clarity of the paper. Lastly, I would like to thank the anonymous reviewers whose suggestions substantially improved this paper.

References

- Noga Alon, Michael Krivelevich, and Benny Sudakov. Finding a large hidden clique in a random graph. In Howard J. Karloff, editor, *Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, 25-27 January 1998, San Francisco, California, USA*, pages 594–598. ACM/SIAM, 1998. URL <http://dl.acm.org/citation.cfm?id=314613.315014>.
- Arash A. Amini and Martin J. Wainwright. High-dimensional analysis of semidefinite relaxations for sparse principal components. *Ann. Statist.*, 37(5B):2877–2921, 10 2009. doi: 10.1214/08-AOS664. URL <https://doi.org/10.1214/08-AOS664>.
- G erard Ben Arous, Alexander S. Wein, and Ilias Zadik. Free energy wells and overlap gap property in sparse PCA. In Jacob D. Abernethy and Shivani Agarwal, editors, *Conference on Learning Theory, COLT 2020, 9-12 July 2020, Virtual Event [Graz, Austria]*, volume 125 of *Proceedings of Machine Learning Research*, pages 479–482. PMLR, 2020. URL <http://proceedings.mlr.press/v125/ben-arous20a.html>.
- Boaz Barak, Jonathan A. Kelner, and David Steurer. Rounding sum-of-squares relaxations. In David B. Shmoys, editor, *Symposium on Theory of Computing, STOC 2014, New York, NY, USA, May 31 - June 03, 2014*, pages 31–40. ACM, 2014. doi: 10.1145/2591796.2591886. URL <https://doi.org/10.1145/2591796.2591886>.
- Quentin Berthet and Philippe Rigollet. Complexity theoretic lower bounds for sparse principal component detection. In *COLT*, volume 30 of *JMLR Workshop and Conference Proceedings*, pages 1046–1066. JMLR.org, 2013a.

- Quentin Berthet and Philippe Rigollet. Computational lower bounds for sparse PCA. *CoRR*, abs/1304.0828, 2013b.
- Quentin Berthet and Philippe Rigollet. Optimal detection of sparse principal components in high dimension. *Ann. Statist.*, 41(4):1780–1815, 08 2013c. doi: 10.1214/13-AOS1127. URL <https://doi.org/10.1214/13-AOS1127>.
- Matthew Brennan, Guy Bresler, and Wasim Huleihel. Reducibility and computational lower bounds for problems with planted sparse structure. In Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet, editors, *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pages 48–166. PMLR, 06–09 Jul 2018. URL <https://proceedings.mlr.press/v75/brennan18a.html>.
- Matthew S. Brennan and Guy Bresler. Optimal average-case reductions to sparse PCA: from weak assumptions to strong hardness. In Alina Beygelzimer and Daniel Hsu, editors, *Conference on Learning Theory, COLT 2019, 25-28 June 2019, Phoenix, AZ, USA*, volume 99 of *Proceedings of Machine Learning Research*, pages 469–470. PMLR, 2019. URL <http://proceedings.mlr.press/v99/brennan19b.html>.
- Matthew S. Brennan, Guy Bresler, Samuel B. Hopkins, Jerry Li, and Tselil Schramm. Statistical query algorithms and low degree tests are almost equivalent. In Mikhail Belkin and Samory Kpotufe, editors, *Conference on Learning Theory, COLT 2021, 15-19 August 2021, Boulder, Colorado, USA*, volume 134 of *Proceedings of Machine Learning Research*, page 774. PMLR, 2021. URL <http://proceedings.mlr.press/v134/brennan21a.html>.
- Alexandre d’Aspremont, Laurent El Ghaoui, Michael I. Jordan, and Gert R. G. Lanckriet. A direct formulation for sparse PCA using semidefinite programming. In *Advances in Neural Information Processing Systems 17 [Neural Information Processing Systems, NIPS 2004, December 13-18, 2004, Vancouver, British Columbia, Canada]*, pages 41–48, 2004. URL <https://proceedings.neurips.cc/paper/2004/hash/8e065119c74efe3a47aec8796964cf8b-Abstract.html>.
- Yash Deshpande and Andrea Montanari. Sparse PCA via covariance thresholding. In *NIPS*, pages 334–342, 2014.
- Ilias Diakonikolas and Daniel Kane. Non-gaussian component analysis via lattice basis reduction. In Po-Ling Loh and Maxim Raginsky, editors, *Conference on Learning Theory, 2-5 July 2022, London, UK*, volume 178 of *Proceedings of Machine Learning Research*, pages 4535–4547. PMLR, 2022. URL <https://proceedings.mlr.press/v178/diakonikolas22d.html>.
- Jingqiu Ding and Yiding Hua. SQ lower bounds for random sparse planted vector problem. *CoRR*, abs/2301.11124, 2023. doi: 10.48550/arXiv.2301.11124. URL <https://doi.org/10.48550/arXiv.2301.11124>.
- Yunzi Ding, Dmitriy Kunisky, Alexander S. Wein, and Afonso S. Bandeira. Subexponential-time algorithms for sparse PCA. *CoRR*, abs/1907.11635, 2019. URL <http://arxiv.org/abs/1907.11635>.

- Yunzi Ding, Dmitriy Kunisky, Alexander S. Wein, and Afonso S. Bandeira. The average-case time complexity of certifying the restricted isometry property. *IEEE Trans. Inf. Theory*, 67(11):7355–7361, 2021. doi: 10.1109/TIT.2021.3112823. URL <https://doi.org/10.1109/TIT.2021.3112823>.
- Tommaso d’Orsi, Pravesh K. Kothari, Gleb Novikov, and David Steurer. Sparse PCA: algorithms, adversarial perturbations and certificates. In Sandy Irani, editor, *61st IEEE Annual Symposium on Foundations of Computer Science, FOCS 2020, Durham, NC, USA, November 16-19, 2020*, pages 553–564. IEEE, 2020. doi: 10.1109/FOCS46700.2020.00058. URL <https://doi.org/10.1109/FOCS46700.2020.00058>.
- Tommaso d’Orsi, Rajai Nasser, Gleb Novikov, and David Steurer. Higher degree sum-of-squares relaxations robust against oblivious outliers. *CoRR*, abs/2211.07327, 2022. doi: 10.48550/arXiv.2211.07327. URL <https://doi.org/10.48550/arXiv.2211.07327>.
- Chao Gao, Zongming Ma, and Harrison H. Zhou. Sparse cca: Adaptive estimation and computational barriers. *The Annals of Statistics*, 45(5):2074–2101, 2017. ISSN 00905364. URL <http://www.jstor.org/stable/26362895>.
- Paul Hand and Laurent Demanet. Recovering the sparsest element in a subspace. *Information and Inference*, 3, 10 2013. doi: 10.1093/imaiai/iau007.
- Guy Holtzman, Adam Soffer, and Dan Vilenchik. A greedy anytime algorithm for sparse pca. pages 1939–1956, 2020. URL <http://learningtheory.org/colt2020/cfp.html>.
- Samuel B. Hopkins, Tselil Schramm, Jonathan Shi, and David Steurer. Fast spectral algorithms from sum-of-squares proofs: tensor decomposition and planted sparse vectors. In Daniel Wichs and Yishay Mansour, editors, *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2016, Cambridge, MA, USA, June 18-21, 2016*, pages 178–191. ACM, 2016. doi: 10.1145/2897518.2897529. URL <https://doi.org/10.1145/2897518.2897529>.
- Iain M. Johnstone and Arthur Yu Lu. On consistency and sparsity for principal components analysis in high dimensions. *Journal of the American Statistical Association*, 104(486):682–693, 2009. doi: 10.1198/jasa.2009.0121. URL <https://doi.org/10.1198/jasa.2009.0121>. PMID: 20617121.
- Robert Krauthgamer, Boaz Nadler, and Dan Vilenchik. Do semidefinite relaxations solve sparse PCA up to the information limit? *The Annals of Statistics*, 43(3):1300 – 1322, 2015. doi: 10.1214/15-AOS1310. URL <https://doi.org/10.1214/15-AOS1310>.
- Dmitriy Kunisky, Alexander S. Wein, and Afonso S. Bandeira. Notes on computational hardness of hypothesis testing: Predictions using the low-degree likelihood ratio. *CoRR*, abs/1907.11636, 2019. URL <http://arxiv.org/abs/1907.11636>.
- B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *Ann. Statist.*, 28(5):1302–1338, 10 2000. doi: 10.1214/aos/1015957395. URL <https://doi.org/10.1214/aos/1015957395>.

- Tengyu Ma and Avi Wigderson. Sum-of-squares lower bounds for sparse PCA. In Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 1612–1620, 2015. URL <https://proceedings.neurips.cc/paper/2015/hash/ec5aa0b7846082a2415f0902f0da88f2-Abstract.html>.
- Pasin Manurangsi, Aviad Rubinfeld, and Tselil Schramm. The strongish planted clique hypothesis and its consequences. In James R. Lee, editor, *12th Innovations in Theoretical Computer Science Conference, ITCS 2021, January 6-8, 2021, Virtual Conference*, volume 185 of *LIPICs*, pages 10:1–10:21. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021. doi: 10.4230/LIPICs.ITCS.2021.10. URL <https://doi.org/10.4230/LIPICs.ITCS.2021.10>.
- Cheng Mao and Alexander S. Wein. Optimal spectral recovery of a planted vector in a subspace. *CoRR*, abs/2105.15081, 2021. URL <https://arxiv.org/abs/2105.15081>.
- Aaron Potechin and Goutham Rajendran. Sub-exponential time sum-of-squares lower bounds for principal components analysis. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=D45iCWZYcff>.
- Qing Qu, Zihui Zhu, Xiao Li, Manolis C. Tsakiris, John Wright, and René Vidal. Finding the sparsest vectors in a subspace: Theory, algorithms, and applications. *CoRR*, abs/2001.06970, 2020. URL <https://arxiv.org/abs/2001.06970>.
- Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018. doi: 10.1017/9781108231596.
- Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019. doi: 10.1017/9781108627771.
- Tengyao Wang, Quentin Berthet, and Richard J. Samworth. Statistical and computational trade-offs in estimation of sparse principal components. *The Annals of Statistics*, 44(5):1896 – 1930, 2016. doi: 10.1214/15-AOS1369. URL <https://doi.org/10.1214/15-AOS1369>.
- Ilias Zadik, Min Jae Song, Alexander S. Wein, and Joan Bruna. Lattice-based methods surpass sum-of-squares in clustering. In Po-Ling Loh and Maxim Raginsky, editors, *Conference on Learning Theory, 2-5 July 2022, London, UK*, volume 178 of *Proceedings of Machine Learning Research*, pages 1247–1248. PMLR, 2022. URL <https://proceedings.mlr.press/v178/zadik22a.html>.

Appendix A. The Wishart Model

Notation. For $m_1, m_2 \in \mathbb{N}$, we use the notation $\mathbb{R}^{m_1 \times m_2}$ for the set of $m_1 \times m_2$ matrices with entries from \mathbb{R} . We denote by $N(0, 1)^m$ an m -dimensional random vector with iid standard Gaussian entries. Similarly, we denote by $N(0, 1)^{m_1 \times m_2}$ an $m_1 \times m_2$ random matrix with iid standard

Gaussian entries. For $m \in \mathbb{N}$, we denote by $[m]$ the set $\{1, 2, \dots, m-1, m\}$. For a vector $v \in \mathbb{R}^m$, we denote by $\|v\|$ its ℓ_2 norm, and for $p \in [1, \infty]$ we denote by $\|v\|_p$ its ℓ_p norm. For a matrix $M \in \mathbb{R}^{m_1 \times m_2}$, we denote by $\|M\|$ its spectral norm, and by $\|M\|_F$ its Frobenius norm. We write \log for the logarithm to the base e .

Recall that the input is an $n \times d$ matrix $\mathbf{Y} = \sqrt{\beta} \mathbf{u} \mathbf{v}^\top + \mathbf{W}$, where $\mathbf{u} \sim N(0, 1)^n$, $v \in \mathbb{R}^d$ is a k -sparse unit vector, $\mathbf{W} \sim N(0, 1)^{n \times d}$ independent of \mathbf{u} . The goal is to compute a unit vector \hat{v} such that $|\langle \hat{v}, v \rangle| \geq 0.99$ with high probability (with respect to the randomness of \mathbf{u} and \mathbf{W}).

First we define vectors $\mathbf{z}_s(r)$ that we will use in the algorithm. Let $t \in \mathbb{N}$ be such that $1 \leq t \leq k$ and let \mathcal{S}_t be the set of all d -dimensional vectors whose entries are in $-\infty, t, \infty$ that have exactly t nonzero coordinates. For $s \in \mathcal{S}_t$ and $r > 0$ let $\mathbf{z}_s(r)$ be the d -dimensional (random) vector defined as

$$\mathbf{z}_{si}(r) = \begin{cases} \mathbf{1}[(\mathbf{Y}^\top \mathbf{Y} s)_i \geq r \cdot t \cdot n] & \text{if } s_i = 0 \\ 1 & \text{otherwise} \end{cases}$$

The following theorem is a restatement of theorem 1.

Theorem 6 *Let $n, d, k, t \in \mathbb{N}$, $\beta > 0$, $0 < \delta < 0.1$. Let $\mathbf{Y} = \sqrt{\beta} \mathbf{u} \mathbf{v}^\top + \mathbf{W}$, where $\mathbf{u} \sim N(0, 1)^n$, $v \in \mathbb{R}^d$ is a k -sparse unit vector, $\mathbf{W} \sim N(0, 1)^{n \times d}$ independent of \mathbf{u} .*

Suppose that $n \gtrsim k + \frac{t \ln^2 d}{\delta^4}$, $k \gtrsim \frac{t \ln d}{\delta^2}$ and

$$\beta \gtrsim \frac{k}{\delta^2 \sqrt{tn}} \sqrt{\ln \left(2 + \frac{td}{k^2} \cdot \left(1 + \frac{d}{n} \right) \right)}.$$

Then there exists an algorithm that, given \mathbf{Y} , k , t and δ , in time $n \cdot d^{O(t)}$ outputs a unit vector \hat{v} such that with probability $1 - o(1)$ as $d \rightarrow \infty$,

$$|\langle \hat{v}, v \rangle| \geq 1 - \delta.$$

Lemma 7 *For all $s \in \mathcal{S}_t$ such that $s_i = 0$ and for all $\tau > 0$,*

$$\mathbb{P} \left[\left(\mathbf{W}^\top \mathbf{Y} s \right)_i \geq \tau \cdot \left(\|\mathbf{W} s\| + |\langle v, s \rangle| \cdot \sqrt{\beta} \cdot \|\mathbf{u}\| \right) \mid \|\mathbf{u}\|, \|\mathbf{W} s\| \right] \leq \exp(-\tau^2/2).$$

Proof *Denote $\mathbf{g} = \mathbf{W} s \sim N(0, t \cdot \text{Id})$.*

$$\left(\mathbf{W}^\top \mathbf{Y} s \right)_i = \langle \mathbf{W}_i, \mathbf{g} \rangle + \sqrt{\beta} \cdot \langle v, s \rangle \cdot \langle \mathbf{W}_i, \mathbf{u} \rangle.$$

Since \mathbf{W}_i , \mathbf{g} and \mathbf{u} are independent, conditional distribution of $(\mathbf{W}^\top \mathbf{Y} s)_i$ given $\|\mathbf{u}\|$ and $\|\mathbf{g}\|$ is $N\left(0, \|\mathbf{g} + \sqrt{\beta} \cdot |\langle v, s \rangle| \cdot \mathbf{u}\|^2\right)$. The lemma follows from the triangle inequality and the tail bound for Gaussian distribution (fact 38). \blacksquare

Lemma 8 *Let $0 < \delta < 0.1$, $r \leq \frac{\beta \delta}{100k}$ and suppose that $n \gtrsim t + \ln^2 d$, $\beta \gtrsim \frac{k}{\delta^2 \sqrt{tn}}$ and $k \gtrsim \frac{\ln d}{\delta^2}$.*

Let $s^* \in \operatorname{argmax}_{s \in \mathcal{S}_t} \langle s, v \rangle$. Then,¹⁶

$$\|v \circ \mathbf{z}_{s^*}(r)\|^2 = \langle v \circ \mathbf{z}_{s^*}(r), v \rangle \geq 1 - \delta,$$

with probability $1 - d^{-10}$.

Proof To simplify the notation we write \mathbf{z}_{si} instead of $\mathbf{z}_{si}(r)$. Let $\mathbf{g} = \mathbf{W} s^*$. By fact 39 with probability at least $1 - \exp(-n/10)$, $n/2 \leq \|\mathbf{u}\|^2 \leq 2n$ and $tn/2 \leq \|\mathbf{g}\|^2 \leq 2tn$, and in this proof we assume that these bounds on $\|\mathbf{u}\|$ and $\|\mathbf{g}\|$ are satisfied.

Let $T = \operatorname{supp}(s^*)$. Note that it is the set of t largest entries of v (by absolute value). Since v is unit, there are two possibilities: either $\|v_T\|^2 > 1 - \delta$, or $\|v_T\|^2 \leq 1 - \delta$. If $\|v_T\|^2 > 1 - \delta$, then the statement is true since $\|v \circ \mathbf{z}_{s^*}\|^2 \geq \|v_T\|^2$. If $\|v_T\|^2 \leq 1 - \delta$, then for every $j \in T$, $|v_j| \geq \sqrt{\delta}/\sqrt{k}$. Hence

$$\langle s^*, v \rangle = \|v_T\|_1 \geq \sqrt{\delta} \cdot t/\sqrt{k}.$$

Let $\mathcal{L} = \left\{ i \in [n] \mid |v_i| \geq \frac{\sqrt{\delta}}{10\sqrt{k}} \right\}$. For all $i \in \mathcal{L}$, $\beta \|\mathbf{u}\|^2 \langle v, s^* \rangle |v_i| \geq 10rtn$. Note that $\|v_{\mathcal{L}}\|^2 \geq 1 - \delta/10$.

Let us write $\mathbf{Y}^\top \mathbf{Y} s^*$ as follows

$$\mathbf{Y}^\top \mathbf{Y} s^* = \beta \|\mathbf{u}\|^2 \langle v, s^* \rangle v + \mathbf{W}^\top \mathbf{W} s^* + \sqrt{\beta} v \mathbf{u}^\top \mathbf{W} s^* + \sqrt{\beta} \langle v, s^* \rangle \mathbf{W}^\top \mathbf{u}.$$

We will bound each term separately

Consider the term $\mathbf{W}^\top \mathbf{W} s^* = \mathbf{W}^\top \mathbf{g}$. By the Chi-squared tail bound (fact 39), with probability at least $1 - \exp(-\tau/2)$,

$$\left\| \left(\mathbf{W}^\top \mathbf{g} \right)_{\mathcal{L} \setminus \operatorname{supp}(s^*)} \right\| \leq 2\sqrt{tn \cdot (|\mathcal{L}| + \tau)}.$$

Since $|\mathcal{L}| \leq k$, with probability at least $1 - \exp(-k)$,

$$\left\| \left(\mathbf{W}^\top \mathbf{g} \right)_{\mathcal{L} \setminus \operatorname{supp}(s^*)} \right\| \leq 10\sqrt{kt} \leq \frac{\delta}{100} \|\beta \|\mathbf{u}\|^2 \langle v, s^* \rangle \cdot v\|,$$

where we used $\beta \gtrsim \frac{k}{\delta^2 \sqrt{nt}}$.

Consider the term $\sqrt{\beta} \langle v, s^* \rangle \mathbf{W}^\top \mathbf{u}$. Since \mathbf{W} and \mathbf{u} are independent, by fact 39 with probability at least $1 - \exp(-k/10)$,

$$\left\| \left(\sqrt{\beta} \langle v, s^* \rangle \mathbf{W}^\top \mathbf{u} \right)_{\mathcal{L}} \right\| \leq 2\sqrt{\beta} \|\mathbf{u}\| \sqrt{k} \cdot \langle v, s^* \rangle \leq \frac{\delta}{100} \|\beta \|\mathbf{u}\|^2 \langle v, s^* \rangle \cdot v\|,$$

where we used $\beta n \gtrsim \frac{k\sqrt{n}}{\delta^2 \sqrt{t}}$ and the fact that $n \geq t$.

Consider the term $\sqrt{\beta} v \mathbf{u}^\top \mathbf{W} s^*$. With probability at least $1 - \exp(-0.1\sqrt{nt})$,

$$\left\| \sqrt{\beta} v \mathbf{u}^\top \mathbf{W} s^* \right\| \leq 2\sqrt{\beta} \cdot \sqrt{t} \cdot \|\mathbf{u}\| \cdot (nt)^{1/4} \leq \frac{\delta}{100} \|\beta \|\mathbf{u}\|^2 \langle v, s^* \rangle \cdot v\|,$$

where we used the fact that the distribution of $\mathbf{u}^\top \mathbf{W} s^*$ given $\|\mathbf{u}\|$ is $N(0, t \cdot \|\mathbf{u}\|^2)$.

By lemma 26,

$$\|v_{\mathcal{L}} \circ \mathbf{z}_{s^*}\|^2 \geq (1 - \delta/10) \|v_{\mathcal{L}}\|^2 \geq 1 - \delta.$$

■

16. Here and further we denote by $a \circ b$ the entrywise product of vectors $a, b \in \mathbb{R}^d$.

Lemma 9 Let $r > 0$ and $s \in \mathcal{S}_t$. With probability at least $1 - \delta$, number of entries i such that

$$|(\mathbf{W}^T Y s)_i| \geq t \cdot r \cdot n$$

is bounded by $pd + 2\sqrt{pd \ln(1/\delta)}$, where $p = \exp(-tnr^2/2)$.

Proof By lemma 7 and Chernoff bound fact 35, number of such entries is at most $pd + \tau d$ with probability at least $1 - \exp(-\frac{\tau^2 d}{2p})$. With $\tau = 2\sqrt{\ln(1/\delta)p/d}$ we get the desired bound. ■

Lemma 10 For $s \in \mathcal{S}_t$, let $\mathbf{N}(s) = (\mathbf{Y}^\top \mathbf{Y} - n \cdot \text{Id} - \beta \|\mathbf{u}\|^2 v^\top v) \circ (\mathbf{z}_s(r) \mathbf{z}_s^\top(r))$ and let $p = \exp(-tnr^2/2)$. Suppose that $k \geq \ln d$.

Then for each $s \in \mathcal{S}_t$, with probability $1 - 2d^{-10}$,

$$\|\mathbf{N}(s)\| \leq 10 \sqrt{(n + \beta n) \cdot (pd + k) \ln\left(\frac{ed}{pd + k}\right)} + 10(pd + k) \ln\left(\frac{ed}{pd + k}\right)$$

Proof To simplify the notation we write z_{si} instead of $z_{si}(r)$. By lemma 9, number of nonzero z_{si} is at most $2pd + 2\sqrt{pd \ln(d)} + 2k \leq 4pd + 4k$ with probability at least $1 - d^{-10}$. Applying lemma 42, we get the desired bound. ■

Lemma 11 Let $0 < \delta < 0.1$ and $\frac{\delta\beta}{200k} \leq r \leq \frac{\delta\beta}{100k}$. Let $p = \exp(-tnr^2/2)$. Suppose that $n \geq t \ln^2 d$, $k \geq t \ln d$ and

$$\beta \gtrsim \frac{k}{\delta^2 \sqrt{tn}} \sqrt{\ln\left(2 + \frac{td}{k^2} \cdot \left(1 + \frac{d}{n}\right)\right)}.$$

Then

$$\sqrt{(n + \beta n) \cdot (pd + k) \ln\left(\frac{ed}{pd + k}\right)} + (pd + k) \ln\left(\frac{ed}{pd + k}\right) \leq \delta\beta n / 10.$$

Proof First note that since

$$\ln\left(2 + \frac{td}{k^2} + \frac{td^2}{k^2 n}\right) \geq 0.5 \cdot \ln\left(2 + \frac{td}{k^2} + \sqrt{\frac{td^2}{k^2 n}}\right),$$

we get

$$\beta \gtrsim \frac{k}{\delta^2 \sqrt{tn}} \sqrt{\ln\left(2 + \frac{td}{k^2} \cdot \left(1 + \frac{k}{\sqrt{nt}}\right)\right)}.$$

Since

$$\sqrt{(n + \beta n) \cdot (pd + k) \ln\left(\frac{ed}{pd + k}\right)} \leq \sqrt{(n + \beta n) \cdot pd \ln\left(\frac{e}{p}\right)} + \sqrt{(n + \beta n) \cdot k \ln\left(\frac{ed}{pd + k}\right)},$$

we can bound the term with pd and the term with k separately. Let us bound the first term. Note that

$$\exp(-tnr^2/4) \leq \exp\left[-\frac{1}{\delta^2} \ln\left(2 + \frac{td}{k^2} \cdot \left(1 + \frac{k}{\sqrt{nt}}\right)\right)\right] \leq \frac{\delta^2 k}{\sqrt{td + kd\sqrt{t/n}}} \lesssim \min \beta, \sqrt{\beta} \cdot \delta^2 \sqrt{n/d}.$$

Hence

$$\left(\sqrt{n} + \sqrt{\beta n}\right) \cdot \left(\exp\left(-\frac{tnr^2}{2}\right) tnr^2 \sqrt{dn}\right) \leq \left(\sqrt{n} + \sqrt{\beta n}\right) \cdot \left(\exp\left(-\frac{tnr^2}{4}\right) \sqrt{dn}\right) \leq 0.01\delta^2 \beta n.$$

Since $t \leq \frac{k}{\ln(ed/k)}$ and $n \geq k \ln(ed/k)$, the second term can be bounded as follows

$$\sqrt{(n + \beta n) \cdot k \ln\left(\frac{ed}{pd + k}\right)} \leq \sqrt{nk \ln\left(\frac{ed}{k}\right)} + \sqrt{\beta nk \ln\left(\frac{ed}{pd + k}\right)} \leq 0.01\delta \beta n.$$

Now, let us bound

$$(pd + k) \ln\left(\frac{ed}{pd + k}\right) \leq pd \ln(e/p) + k \ln(ed/k).$$

The first term can be bounded as follows:

$$pd \ln(e/p) \leq d \cdot \frac{k^2}{td + kd\sqrt{t/n}} \leq k\sqrt{n/t} \leq 0.01\delta \beta n.$$

For the second term,

$$k \ln(ed/k) \leq k\sqrt{n/t} \leq 0.01\delta \beta n. \quad \blacksquare$$

Lemma 12 *Let $0 < \delta < 0.1$ and $\frac{\delta\beta}{200k} \leq r \leq \frac{\delta\beta}{100k}$. For $s \in \mathcal{S}_t$ let $\hat{v}(s)$ be the top¹⁷ eigenvector of $(\mathbf{Y}^\top \mathbf{Y}) \circ (\mathbf{z}_s(r) \mathbf{z}_s^\top(r))$.*

Suppose that $n \gtrsim \frac{t \ln^2 d}{\delta^4}$, $k \gtrsim \frac{t \ln d}{\delta^2}$ and

$$\beta \gtrsim \frac{k}{\delta^2 \sqrt{tn}} \sqrt{\ln\left(2 + \frac{td}{k^2} \cdot \left(1 + \frac{d}{n}\right)\right)}.$$

Then there exists $s' \in \mathcal{S}_t$ such that with probability $1 - 10d^{-10}$,

$$|\langle \hat{v}(s'), v \rangle| \geq 1 - 20\delta.$$

Proof Let $s' \in \operatorname{argmax}_{s \in \mathcal{S}_t} \langle s, v \rangle$ and let $\tilde{v} = v \circ \mathbf{z}_{s'}(r)$. Then

$$\left(\mathbf{Y}^\top \mathbf{Y}\right) \circ \left(\mathbf{z}_{s'}(r) \mathbf{z}_{s'}^\top(r)\right) = \beta \|\mathbf{u}\|^2 \tilde{v} \tilde{v}^\top + \mathbf{N}(s'),$$

17. A unit eigenvector that corresponds to the largest eigenvalue

where $\mathbf{N}(s')$ is the same as in lemma 10. By lemma 11 and lemma 10, with probability $1 - 2d^{-10}$,

$$\|\mathbf{N}(s')\| \leq \delta\beta n.$$

Since with probability at least $1 - \exp(-n/10)$, $\|\mathbf{u}\|^2 \geq n/2$, we can use standard (non-sparse) PCA. Concretely, by lemma 32, $|\langle \hat{\mathbf{v}}(s'), \tilde{\mathbf{v}} \rangle| \geq \|\tilde{\mathbf{v}}\|^2 - 2\delta$. By lemma 8, with probability $1 - 10d^{-10}$, $\|\tilde{\mathbf{v}}\|^2 = \langle \tilde{\mathbf{v}}, v \rangle \geq 1 - \delta$. Using fact 33 we get the desired bound. \blacksquare

Proof [Proof of theorem 6] Since $\sqrt{n/2} \leq \|\mathbf{u}\| \leq \sqrt{2n}$ with probability at least $1 - \exp(-n/10)$, in this proof we assume that this bound on $\|\mathbf{u}\|$ holds.

Let $\tilde{\delta} = \delta/100$. Note that we can apply lemma 12 if $\frac{\tilde{\delta}\beta}{200k} \leq r \leq \frac{\tilde{\delta}\beta}{100k}$. Since we do not know β , we can create a list of candidates for r of size at most $2 \ln n$ (from $r = 1/n$ to $r = 1$).

For all $s \in \mathcal{S}_t$ and for all candidates for r we compute $\hat{\mathbf{v}}(s)$ as in lemma 12. By lemma 12, we get a list of vectors \mathbf{L} of size $2 \ln(n) \cdot |\mathcal{S}_t|$ such that with probability $1 - 20 \ln(n)d^{-10} \geq 1 - d^{-9}$ there exists $\mathbf{v}^* \in \mathbf{L}$ such that $|\langle \mathbf{v}^*, v \rangle| \geq 1 - 20\tilde{\delta} \geq 1 - \delta/5$.

By fact 41 and lemma 43, k' -sparse norm of $\mathbf{Y}^\top \mathbf{Y} - n\text{Id} - \beta\|\mathbf{u}\|^2 vv^\top$ is bounded by

$$10\sqrt{nk' \ln(ed/k')} + 10\sqrt{\beta nk' \ln(ed/k')}$$

with probability at least $1 - 2d^{-9}$. Let us show that if $k' \lesssim \frac{\delta^2 \beta^2 n}{(1+\beta) \ln d}$, then this k' -sparse norm is bounded by $\delta\beta n/10$. Indeed, if $\beta \geq 1$, then

$$\frac{\delta^2 \beta^2 n}{(1+\beta) \ln d} \geq \frac{\delta^2 \beta n}{2 \ln d} \gtrsim k \sqrt{\frac{n}{t \ln^2 d}} \gtrsim k/\delta^2,$$

and if $\beta < 1$,

$$\frac{\delta^2 \beta^2 n}{(1+\beta) \ln d} \geq \frac{\delta^2 \beta^2 n}{2 \ln^2 d} \gtrsim \frac{k^2}{\delta^2 t \ln d} \gtrsim k/\delta^2.$$

Applying lemma 28 with $k' = \lceil 100k/\delta^2 \rceil$, we can compute a unit vector $\hat{\mathbf{v}}$ such that

$$|\langle \hat{\mathbf{v}}, v \rangle| \geq 1 - \delta. \quad \blacksquare$$

Appendix B. Adversarial Perturbations

The sparse planted vector problem. Let us show that the sparse planted vector problem is a special case of sparse PCA with perturbation matrix $E = -\frac{1}{\|\mathbf{u}\|^2} \mathbf{u}\mathbf{u}^\top \mathbf{W}$. Let $\mathbf{Y} = \mathbf{R}\mathbf{B}$. It follows that

$$\mathbf{Y} = \|\mathbf{g}_n\| \mathbf{R}_n v^\top + \sum_{i=1}^{n-1} \mathbf{R}_i \mathbf{g}_i^\top.$$

Note that $\sum_{i=1}^{n-1} \mathbf{R}_i \mathbf{g}_i^\top$ has the same distribution as $(\text{Id} - \mathbf{R}_n \mathbf{R}_n^\top) \mathbf{W}$, where $\mathbf{W} \sim N(0, 1)^{n \times d}$ is independent of \mathbf{R} . Hence for Gaussian vector \mathbf{u} such that $\frac{1}{\|\mathbf{u}\|} \mathbf{u} = \mathbf{R}_n$, we get

$$\mathbf{Y} = \frac{\|\mathbf{g}_n\|}{\|\mathbf{u}\|} \mathbf{u} v^\top + \mathbf{W} - \frac{1}{\|\mathbf{u}\|^2} \mathbf{u}\mathbf{u}^\top \mathbf{W}.$$

With high probability $\beta := \frac{\|g_n\|^2}{\|\mathbf{u}\|^2} = (1 + o(1))d/n$. Note that columns of $-\frac{1}{\|\mathbf{u}\|^2}\mathbf{u}\mathbf{u}^\top\mathbf{W}$ have norm at most $10\sqrt{\log d}$ with high probability. In these settings, ε from theorem 1 is allowed to be as large as $\Omega(1)$, and we get a bound $\|E\|_{1 \rightarrow 2} \lesssim \sqrt{d/k}$. Hence for $d \gtrsim k \log d$, E is allowed to have columns of norm $10\sqrt{\log d}$.

The Wishart model with perturbations. The following theorem is a restatement of theorem 2.

Theorem 13 *Let $n, d, k, t \in \mathbb{N}$, $\beta > 0$, $0 < \delta < 0.1$. Let*

$$\tilde{Y} = \sqrt{\beta}\mathbf{u}v^\top + \mathbf{W} + E,$$

where $\mathbf{u} \sim N(0, 1)^n$, $v \in \mathbb{R}^d$ is a k -sparse unit vector, $\mathbf{W} \sim N(0, 1)^{n \times d}$ independent of \mathbf{u} and $E \in \mathbb{R}^{n \times d}$ is a matrix such that

$$\|E\|_{1 \rightarrow 2} = b \leq \varepsilon \cdot \min \sqrt{\beta}, \beta \cdot \sqrt{n/k},$$

where $\|E\|_{1 \rightarrow 2}$ is the maximal norm of the columns of E and $\varepsilon\sqrt{\ln(1/\varepsilon)} \lesssim \delta^6 \min\{1, \min\{\beta, \sqrt{\beta}\} \cdot \sqrt{n/d}\}$.

Suppose that $n \gtrsim k + \frac{t \ln^2 d}{\delta^4}$, $k \gtrsim \frac{t \ln d}{\delta^2}$ and

$$\beta \gtrsim \frac{k}{\delta^6 \sqrt{tn}} \sqrt{\ln\left(2 + \frac{td}{k^2} \left(1 + \frac{d}{n}\right)\right)}.$$

Then there exists an algorithm that, given \mathbf{Y} , k and t , in time $n \cdot d^{O(t)}$ outputs a unit vector \hat{v} such that with probability $1 - o(1)$ as $d \rightarrow \infty$,

$$|\langle \hat{v}, v \rangle| \geq 1 - \delta.$$

Lemma 14 *Let $\mathcal{L} \subset [d]$ be the set of indices of the largest (in absolute value) ℓ entries of v , where $\ell = \min\{\lceil \|v\|_1^2 / \delta^2 \rceil, k\}$. Suppose that $\beta \gtrsim \frac{k}{\delta^6 \sqrt{nt}}$ and $n \gtrsim k + t \ln^2 d$, $k \gtrsim \ln d$.*

Then with probability $1 - d^{-10}$ there exists $s' \in \mathcal{S}_t$ such that either

$$\left\| \left(\left(\tilde{Y}^\top \tilde{Y} - \beta \|\mathbf{u}\|^2 v v^\top \right) s' \right)_{\mathcal{L} \setminus \text{supp}(s')} \right\| \leq 10\delta \cdot \left\| \beta \|\mathbf{u}\|^2 v v^\top s' \right\|,$$

or

$$\|v_{\text{supp}(s')}\| \geq 1 - \delta.$$

Moreover,

$$\left\| \left(E^\top \mathbf{Y} + \mathbf{Y}^\top E \right) s' \right\| \leq 10^4 \cdot \frac{\varepsilon}{\delta^2} \cdot \sqrt{d} \cdot \frac{\beta n t}{\sqrt{k} \cdot \|v\|_1}.$$

Proof First note that

$$\|v_{\mathcal{L}}\|_1 \geq \|v\|_1 - \delta \geq (1 - \delta)\|v\|_1.$$

and

$$\|v_{\mathcal{L}}\| \geq \|v\| - \|v_{[d] \setminus \mathcal{L}}\|_1 \geq \|v\| - \delta \geq (1 - \delta)\|v\|.$$

Also note that for all $i \in \mathcal{L}$, $|v_i| \geq \frac{\delta}{\|v\|_1}$.

We will define a distribution over $s \in \mathcal{S}_t$, and then we show via probabilistic method that s' with desired properties exists by bounding terms of $\tilde{Y}^\top \tilde{Y} s$ one by one. Also, with probability at least $1 - \exp(-n/2)$, $2\sqrt{n} \geq \|u\| \geq \sqrt{n}/2$ and in the proof we will always assume that $2\sqrt{n} \geq \|u\| \geq \sqrt{n}/2$.

Consider the partition of \mathcal{L} into $m = \lceil |\mathcal{L}|/t \rceil$ disjoint blocks b_1, \dots, b_m of size¹⁸ t . Each block b_j corresponds to $s(j) \in \mathcal{S}_t$ such that $s(j)_i = 0$ for $i \notin b_j$ and $s(j)_i = \text{sign}(v_i)$ for $i \in b_j \cap \text{supp}(v)$.

If $|\mathcal{L}| \leq t$, then $\|v_{\text{supp}(s(1))}\| \geq \|v_{\mathcal{L}}\| \geq 1 - \delta$.

If $|\mathcal{L}| > t$, consider the uniform distribution \mathcal{U} over $s(j)$. We get

$$\mathbb{E}_{s \sim \mathcal{U}} \langle v, s \rangle = \frac{1}{m} \sum_{j=1}^m \|v \circ s(j)\|_1 = \frac{1}{m} \|v\|_1 \geq \frac{t}{2\mathcal{L}} \|v\|_1 \geq \frac{\delta^2 t}{2\|v\|_1} \geq \frac{\delta^2 t}{2\sqrt{k}}.$$

Moreover,

$$\mathbb{E}_{s \sim \mathcal{U}} \langle v, s \rangle = \frac{1}{m} \|v\|_1 \leq \frac{t}{\mathcal{L}} \|v\|_1 \leq \frac{t}{\|v\|_1}.$$

Let us write $\tilde{Y}^\top \tilde{Y}$ as follows:

$$\tilde{Y}^\top \tilde{Y} = \beta \|u\|^2 v v^\top + \mathbf{W}^\top \mathbf{W} + \sqrt{\beta} v u^\top \mathbf{W} + \sqrt{\beta} \mathbf{W}^\top u v^\top + E^\top E + \sqrt{\beta} v u^\top E + \sqrt{\beta} E^\top u v^\top + \mathbf{W}^\top E + E^\top \mathbf{W}.$$

We will bound each term separately.

Consider the term $E \mathbf{W}^\top$.

$$\mathbb{E}_{s \sim \mathcal{U}} |\langle E_i, \mathbf{W} s \rangle|^2 = \frac{1}{m} \sum_{j=1}^m |\langle E_i, \mathbf{W} s(j) \rangle|^2$$

Since for different j_1 and j_2 , $s(j_1)$ and $s(j_2)$ have disjoint supports, $\mathbf{W} s(j) \sim N(0, t)^n$ are independent. By the concentration of spectral norm of Gaussian matrices (fact 40),

$$\sum_{j=1}^m \langle E_i, \mathbf{W} s(j) \rangle^2 \leq 4t(m+n) \cdot \|E_i\|^2 \leq 10t n b^2$$

with probability at least $1 - \exp(-n)$. Hence

$$\mathbb{E}_{s \sim \mathcal{U}} |\langle E_i, \mathbf{W} s \rangle|^2 \leq \frac{10t b^2 n}{m}.$$

Therefore, with probability at least $1 - d \exp(-n)$,

$$\mathbb{E}_{s \sim \mathcal{U}} \sum_{i \in \mathcal{L}} |\langle E_i, \mathbf{W} s \rangle|^2 \leq |\mathcal{L}| \frac{10t b^2 n}{m}$$

and

$$\mathbb{E}_{s \sim \mathcal{U}} \sum_{i \in [d]} |\langle E_i, \mathbf{W} s \rangle|^2 \leq d \frac{10t b^2 n}{m}$$

18. If the last block has smaller size, we can add arbitrary entries to it.

Hence with probability at least $1 - d \exp(-n/2)$ there exists $s' = s(j)$ such that

$$\frac{10t}{\|v\|_1} \geq \langle v, s(j) \rangle \geq \frac{\delta^2 t}{10\|v\|_1} \geq \frac{\delta^2 t}{10\sqrt{k}},$$

and

$$\sum_{i \in \mathcal{L}} |\langle E_i, \mathbf{W} s(j) \rangle|^2 \leq |\mathcal{L}| \frac{40tb^2 n}{m} \leq 40t^2 \varepsilon^2 \beta^2 n^2 / k \leq \left(\frac{100\varepsilon}{\delta^2} \right)^2 \left\| \beta \|\mathbf{u}\|^2 v v^\top s(j) \right\|^2 \leq \delta^2 \left\| \beta \|\mathbf{u}\|^2 v v^\top s(j) \right\|^2,$$

where we used $|\mathcal{L}| \leq mt$, and

$$\sum_{i \in [d]} |\langle E_i, \mathbf{W} s' \rangle|^2 \leq \left(\frac{100\varepsilon}{\delta^2} \right)^2 \cdot \frac{d}{k} \cdot \left\| \beta \|\mathbf{u}\|^2 v v^\top s(j) \right\|^2.$$

Consider the term $\mathbf{W}^\top E$. By fact 40, with probability at least $1 - \exp(-n/2)$,

$$\|\mathbf{W}^\top E s(j)\| \leq t b \|\mathbf{W}^\top\| \leq 10tb\sqrt{n} \leq 10t\varepsilon\beta n / \sqrt{k} \leq \frac{400\varepsilon}{\delta^2} \|\beta \|\mathbf{u}\|^2 \langle v, s(j) \rangle \cdot v\| \leq \delta \left\| \beta \|\mathbf{u}\|^2 v v^\top s(j) \right\|.$$

Consider the term $\sqrt{\beta} E^\top \mathbf{u} v^\top$:

$$\left\| \left(\sqrt{\beta} E^\top \mathbf{u} v^\top s(j) \right)_{\mathcal{L} \setminus \text{supp}(s')} \right\| \leq \sqrt{k} \left\| \sqrt{\beta} E^\top \mathbf{u} v^\top s(j) \right\|_\infty \leq \sqrt{k} \beta \cdot b \cdot \|\mathbf{u}\| \cdot \langle v, s(j) \rangle \leq 10\varepsilon \left\| \beta \|\mathbf{u}\|^2 \langle v, s(j) \rangle \cdot v \right\|.$$

Consider the term $\sqrt{\beta} v \mathbf{u}^\top E$:

$$\left\| \sqrt{\beta} v \mathbf{u}^\top E s(j) \right\| \leq \sqrt{\beta} t b \|\mathbf{u}\| \leq \frac{\varepsilon}{\delta^2} \cdot \frac{\delta^2 t}{\sqrt{k}} \cdot \beta \sqrt{n} \|\mathbf{u}\| \leq \frac{10\varepsilon}{\delta^2} \left\| \beta \|\mathbf{u}\|^2 \langle v, s(j) \rangle \cdot v \right\|.$$

Consider the term $E^\top E$:

$$\left\| \left(E^\top E s(j) \right)_{\mathcal{L} \setminus \text{supp}(s')} \right\| \leq \sqrt{k} \cdot \left\| E^\top E s(j) \right\|_\infty \leq \sqrt{k} \cdot b^2 t \leq \sqrt{k} \cdot \frac{\varepsilon^2}{\delta^2} \cdot \frac{\delta^2 t}{k} \beta n \leq 10\varepsilon \left\| \beta \|\mathbf{u}\|^2 \langle v, s(j) \rangle \cdot v \right\|.$$

Moreover, note that from the bounds above we also get

$$\left\| \left(E^\top \mathbf{Y} + \mathbf{Y}^\top E \right) s(j) \right\| \leq 1000 \cdot \frac{\varepsilon}{\delta^2} \cdot \sqrt{\frac{d}{k}} \cdot \left\| \beta \|\mathbf{u}\|^2 v v^\top s(j) \right\| \leq 10^4 \frac{\varepsilon}{\delta^2} \cdot \sqrt{d} \cdot \frac{\beta n t}{\sqrt{k} \cdot \|v\|_1}.$$

Consider the term $\mathbf{W}^\top \mathbf{W}$. By the Chi-squared tail bound (fact 39), with probability at least $1 - \exp(-\tau/2)$,

$$\left\| \left(\mathbf{W}^\top \mathbf{W} s(j) \right)_{\mathcal{L} \setminus \text{supp}(s(j))} \right\| \leq 2\sqrt{tn \cdot (|\mathcal{L}| + \tau)}.$$

Since $|\mathcal{L}| \leq k$, with probability at least $1 - \exp(-k)$,

$$\left\| \left(\mathbf{W}^\top \mathbf{W} s(j) \right)_{\mathcal{L} \setminus \text{supp}(s(j))} \right\| \leq 10\sqrt{kt n} \leq \delta \left\| \beta \|\mathbf{u}\|^2 \langle v, s(j) \rangle \cdot v \right\|,$$

where we used $\beta \gtrsim \frac{k}{\delta^3 \sqrt{nt}}$.

Consider the term $\sqrt{\beta} \mathbf{W}^\top \mathbf{u} v^\top$. Since \mathbf{W} and \mathbf{u} are independent, with probability at least $1 - \exp(-k/2)$,

$$\left\| \left(\sqrt{\beta} \mathbf{W}^\top \mathbf{u} v^\top s(j) \right)_{\mathcal{L}} \right\| \leq 2\sqrt{\beta} \|\mathbf{u}\| \sqrt{k} \cdot \langle v, s(j) \rangle \leq \delta \|\beta\| \|\mathbf{u}\|^2 \langle v, s(j) \rangle \cdot \|v\|,$$

where we used $\beta n \gtrsim \frac{k\sqrt{n}}{\delta^2 \sqrt{t}}$ and the fact that $n \geq t$.

Consider the term $\sqrt{\beta} \mathbf{W}^\top \mathbf{u} v^\top$. With probability at least $1 - \exp(-0.1\sqrt{nt})$,

$$\left\| \sqrt{\beta} v \mathbf{u}^\top \mathbf{W} s(j) \right\| \leq 2\sqrt{\beta} \cdot \sqrt{t} \cdot \|\mathbf{u}\| \cdot (nt)^{1/4} \leq \delta \|\beta\| \|\mathbf{u}\|^2 \langle v, s(j) \rangle \cdot \|v\|,$$

where we used the fact that the distribution of $\mathbf{u}^\top \mathbf{W} s(j)$ given $\|\mathbf{u}\|$ is $\sim N(0, t \cdot \|\mathbf{u}\|^2)$. ■

For $s \in \mathcal{S}_t$ let $z_s(r)$ be the n -dimensional (random) vector defined as

$$z_{si}(r) = \begin{cases} \mathbf{1}[(\tilde{Y}^\top \tilde{Y} s)_i \geq r \cdot t \cdot n] & \text{if } s_i = 0 \\ 1 & \text{otherwise} \end{cases}$$

Lemma 15 Suppose that $k \gtrsim t \ln d$ and

$$\beta \gtrsim \frac{k}{\delta^2 \sqrt{tn}} \sqrt{\ln \left(2 + \frac{td}{k^2} \cdot \left(1 + \frac{d}{n} \right) \right)}.$$

Let r be such that

$$\frac{\beta \delta}{1000 \|v\|_1^2} \leq r \leq \frac{\beta \delta}{500 \|v\|_1^2}.$$

Let $s \in \mathcal{S}_t$ and

$$N(s) = \left(\tilde{Y}^\top \tilde{Y} - n \text{Id} - \beta \|\mathbf{u}\|^2 v v^\top \right) \circ \left(z_s(r) z_s^\top(r) \right).$$

Let

$$\mathcal{P}_k = \left\{ X \in \mathbb{R}^{d \times d} \mid X \succeq 0, \text{Tr } X = 1, \|X\|_1 \leq k \right\}.$$

Then with probability $1 - 10 \cdot d^{-10}$, for all $X \in \mathcal{P}_k$,

$$|\langle X, N \rangle| \leq \delta \beta n.$$

Proof To simplify the notation we write z instead of $z_s(r)$ and N instead of $N(s)$. Let

$$N_E = \left(\tilde{Y}^\top \tilde{Y} - (\tilde{Y} - E)^\top (\tilde{Y} - E) \right) \circ (z z^\top).$$

Note that

$$|\langle X, N \rangle| \leq \|N - N_E\| + |\langle X, N_E \rangle|.$$

By lemma 34,

$$|\langle X, N_E \rangle| \leq b^2 k + 2b \sqrt{k \left\| (\tilde{Y} - E)^\top (\tilde{Y} - E) \circ (z z^\top) \right\|}.$$

The first term can be bounded as follows:

$$b^2k \leq \varepsilon\beta^2n/k \leq \delta\beta\sqrt{n}/100 \leq \delta\beta n/100.$$

Note that with probability at least $1 - \exp(-n)$,

$$\left\| \left(\tilde{Y} - E \right)^\top \left(\tilde{Y} - E \right) \circ \left(zz^\top \right) \right\| \leq 10(\beta n + n).$$

Hence for the second term,

$$b\sqrt{k(\beta n + n)} \leq \varepsilon \cdot \frac{2\beta}{1 + \sqrt{\beta}} n \left(1 + \sqrt{\beta} \right) \leq \varepsilon\beta n \leq \delta\beta n/100.$$

Therefore,

$$|\langle X, N_E \rangle| \leq \delta\beta n/10.$$

By fact 25 and the bound on $\|(E^\top \mathbf{Y} + \mathbf{Y}^\top E)s'\|$ from lemma 14, at most $\tilde{\varepsilon}^2 d = (10^4 \varepsilon / \delta)^2 d$ entries of $(E^\top \tilde{Y} + \tilde{Y}^\top E)s'$ are larger (in absolute value) than $rt n/2$. By lemma 9, with probability at least $1 - d^{-10}$, number of entries i such that

$$|(\mathbf{W}^\top \mathbf{Y} s')_i| \geq t \cdot r \cdot n/2$$

is bounded by $pd + 10\sqrt{pd \ln d}$, where $p = \exp(-tnr^2/8)$. Hence number of nonzero entries of z is at most $10\tilde{\varepsilon}^2 d + pd + 10\sqrt{pd \ln d} + 2k \leq 10\tilde{\varepsilon}^2 d + 10pd + 10k$ with probability at least $1 - d^{-10}$.

Therefore, by lemma 42 and lemma 11, with probability at least $1 - 2d^{-10}$,

$$\|N - N_E\| \leq \delta\beta n/10 + \sqrt{(n + \beta n) \cdot \tilde{\varepsilon}^2 d \ln(1/\tilde{\varepsilon})} + \tilde{\varepsilon}^2 d \ln(1/\tilde{\varepsilon}).$$

Since the second and the third terms are bounded by $\delta\beta n/10$, we get the desired bound. \blacksquare

Lemma 16 *Let $0 < \delta < 0.1$ and let $\frac{\beta\delta^3}{1000\|v\|_1^2} \leq r \leq \frac{\beta\delta^3}{500\|v\|_1^2}$.*

For $s \in \mathcal{S}_t$ let $\hat{v}(s)$ be the top eigenvector of $X \in \mathcal{P}_k$ that maximizes $\langle X, \tilde{Y} \circ (\tilde{z}_s \tilde{z}_s^\top) \rangle$.

Suppose that $n \gtrsim k + t \ln^2 d$, $k \gtrsim t \ln d$ and

$$\beta \gtrsim \frac{k}{\delta^6 \sqrt{tn}} \sqrt{\ln \left(2 + \frac{td}{k^2} \cdot \left(1 + \frac{d}{n} \right) \right)}.$$

Then there exists $s' \in \mathcal{S}_t$ such that with probability $1 - 20d^{-10}$,

$$|\langle \hat{v}(s'), v \rangle| \geq 1 - \delta/10.$$

Proof Let s' be as in lemma 14 and let $\tilde{v} = v \circ z_{s'}(r)$. Then

$$\left(\tilde{Y}^\top \tilde{Y} - n\text{Id} \right) \circ \left(z_{s'}(r) z_{s'}^\top(r) \right) = \beta \|\mathbf{u}\|^2 \tilde{v} \tilde{v}^\top + N(s),$$

where $N(s)$ is the same as in lemma 15. By lemma 15, with probability $1 - 10 \cdot d^{-10}$, for all $X \in \mathcal{P}_k$,

$$|\langle X, N \rangle| \leq \delta \beta n.$$

Consider

$$\hat{\mathbf{X}}(s^*) \in \operatorname{argmax}_{X \in \mathcal{P}} \left\langle X, \tilde{Y} \circ \left(\tilde{z}_{s^*} \tilde{z}_{s^*}^\top \right) \right\rangle.$$

By lemma 32,

$$\left\langle \hat{\mathbf{X}}(s^*), vv^\top \circ \left(\tilde{z}_{s^*} \tilde{z}_{s^*}^\top \right) \right\rangle \geq 1 - 6\tilde{\delta}.$$

Hence by fact 33,

$$\left\langle \hat{\mathbf{X}}(s^*), vv^\top \right\rangle \geq 1 - 24\tilde{\delta}.$$

By lemma 29, $|\langle \hat{v}(s^*), \tilde{v} \rangle| \geq 1 - 100\tilde{\delta}$. By lemma 14 and lemma 26, with probability $1 - d^{-10}$,

$$\langle \tilde{v}, v \rangle \geq 1 - 100\tilde{\delta}.$$

Using fact 33 we get the desired bound. \blacksquare

Proof [Proof of theorem 13] Since $\sqrt{n/2} \leq \|u\| \leq \sqrt{2n}$ with probability at least $1 - \exp(-n/10)$, in this proof we assume that this bound on $\|u\|$ holds. Let $\tilde{\delta} = \delta/1000$. Note that we can apply lemma 15 if $\frac{\beta \tilde{\delta}^3}{1000 \|v\|_1^2} \leq r \leq \frac{\beta \tilde{\delta}^3}{500 \|v\|_1^2}$. Since we do not know β and $\|v\|_1$, we can create a list of candidates for r of size at most $2 \ln n$ (starting from $r = 1/n$ and finishing at $r = 1$).

For all $s \in \mathcal{S}_t$ and for all candidates for r we compute $\hat{v}(s)$ as in lemma 16. By lemma 16, we get a list of vectors \mathbf{L} of size $2 \ln(n) \cdot |\mathcal{S}_t|$ such that with probability $1 - 20 \ln(n) d^{-10} \geq 1 - d^{-9}$ there exists $v^* \in \mathbf{L}$ such that $|\langle v^*, v \rangle| \geq 1 - \tilde{\delta}/10$.

By fact 41, lemma 43, lemma 34, k' -sparse norm of $\mathbf{Y}^\top \mathbf{Y} - n \operatorname{Id} - \beta \|u\|^2 vv^\top$ is bounded by

$$10\sqrt{nk' \ln(ed/k')} + 10\sqrt{\beta nk' \ln(ed/k')} + b^2 k' + 10b\sqrt{k'(\beta n + n)}$$

with probability at least $1 - 3d^{-9}$. Let us show that if

$$k' \lesssim \min \left\{ \delta \beta n / b^2, \frac{\delta^2 \beta^2 n}{(1 + \beta) \ln d}, \frac{\delta^2 \beta^2 n}{(1 + \beta) b^2 \ln d} \right\},$$

then k' the k' -sparse norm is bounded by $\delta \beta n$. Indeed, if $\beta \geq 1$, then

$$\frac{\delta^2 \beta^2 n}{(1 + \beta) \ln d} \geq \frac{\delta^2 \beta n}{2 \ln d} \gtrsim k \sqrt{\frac{n}{t \ln^2 d}} \gtrsim k / \delta^2,$$

and if $\beta < 1$,

$$\frac{\delta^2 \beta^2 n}{(1 + \beta) \ln d} \geq \frac{\delta^2 \beta^2 n}{2 \ln d} \gtrsim \frac{k^2}{\delta^2 t \ln d} \gtrsim k / \delta^2,$$

and by our bound on b ,

$$k / \delta^2 \lesssim \min \left\{ \delta \beta n / b^2, \frac{\delta^2 \beta^2 n}{(1 + \beta) b^2 \ln d} \right\}.$$

By lemma 28, for $k' = \lceil 100k / \delta^2 \rceil$, we can compute a k' -sparse vector unit vector $\hat{v}(k')$ such that $|\langle \hat{v}(k'), v \rangle| \geq 1 - \delta$. \blacksquare

Appendix C. The Wigner Model

C.1. Classical Settings

In classical settings the input is an $d \times d$ matrix $\mathbf{Y} = \lambda v v^\top + \mathbf{W}$, where $v \in \mathbb{R}^d$ is a k -sparse unit vector, $\mathbf{W} \sim N(0, 1)^{d \times d}$. The goal is to compute a unit vector \hat{v} such that $|\langle \hat{v}, v \rangle| \geq 0.99$ with high probability (with respect to the randomness of \mathbf{W}).

In this section we prove the following theorem.

Theorem 17 *Let $d, k, t \in \mathbb{N}$, $\lambda > 0$, $0 < \delta < 0.1$. Let $\mathbf{Y} = \sqrt{\beta} \mathbf{u} \mathbf{u}^\top + \mathbf{W}$, where $v \in \mathbb{R}^d$ is a k -sparse unit vector, $\mathbf{W} \sim N(0, 1)^{d \times d}$.*

Suppose that $k \gtrsim \frac{t \ln d}{\delta^2}$, and

$$\lambda \gtrsim \frac{k}{\delta^2} \sqrt{\frac{\ln(2 + td/k^2)}{t}}.$$

Then there exists an algorithm that, given \mathbf{Y} , k , t and δ , in time $d^{O(t)}$ outputs a unit vector \hat{v} such that with probability $1 - o(1)$ as $d \rightarrow \infty$,

$$|\langle \hat{v}, v \rangle| \geq 1 - \delta.$$

As for Wishart model, we define vectors $\mathbf{z}_s(r)$ that we will use in the algorithm. Recall the definition of \mathcal{S}_t : for $t \in \mathbb{N}$ such that $1 \leq t \leq k$ we denote by \mathcal{S}_t the set of all d -dimensional vectors with values in $-\infty, t, \infty$ that have exactly t nonzero coordinates. For $s \in \mathcal{S}_t$ let \mathbf{z}_s be d -dimensional (random) vectors defined as

$$\mathbf{z}_{s_i}(r) = \begin{cases} \mathbf{1}_{[(\mathbf{Y}s)_i = \langle s, \mathbf{Y}_i \rangle \geq r \cdot t]} & \text{if } s_i = 0 \\ 1 & \text{otherwise} \end{cases}$$

for some $r > 0$ (here \mathbf{Y}_i denotes the i -th row of \mathbf{Y}).

Lemma 18 *If $i \notin \text{supp}(v)$, then for all $s \in \mathcal{S}_t$ and $r > 0$,*

$$\mathbb{P}[\langle s, \mathbf{Y}_i \rangle \geq r \cdot t] \leq \exp(-tr^2/2).$$

Proof *Since $v_i = 0$,*

$$\langle s, \mathbf{Y}_i \rangle = \langle s, \mathbf{W}_i \rangle \sim N(0, t).$$

The lemma follows from the tail bound for Gaussian distribution (fact 38). ■

Lemma 19 *Let $0 < \delta < 0.1$, $r \leq \frac{\lambda \delta}{100k}$ and suppose that $\lambda \gtrsim \frac{k}{\delta^2 \sqrt{t}}$ and $k \gtrsim \frac{\ln d}{\delta^2}$.*

Let $s^ \in \arg \max_{s \in \mathcal{S}_t} \langle s, v \rangle$. Then, with probability $1 - d^{-10}$*

$$\|v \circ \mathbf{z}_{s^*}\|^2 = \langle v \circ \mathbf{z}_{s^*}, v \rangle \geq 1 - \delta.$$

Proof *To simplify the notation we write \mathbf{z}_{s_i} instead of $\mathbf{z}_{s_i}(r)$. Let $\mathcal{L}_\delta = \left\{ i \in [d] \mid |v_i| \geq \frac{\sqrt{\delta}}{10\sqrt{k}} \right\}$. For all $i \in \mathcal{L}$, $\beta \|\mathbf{u}\|^2 \langle v, s^* \rangle |v_i| \geq 10rt$. Note that $\|v_{\mathcal{L}}\|^2 \geq 1 - \delta/10$.*

As in the proof of lemma 8, we can assume that

$$\langle s^*, v \rangle = \|v_T\|_1 \geq \sqrt{\delta t} / \sqrt{k}.$$

We need to bound the norm of $\mathbf{W} s^* \sim N(0, t \cdot \text{Id})$ restricted to the entries of \mathcal{L} . By fact 39, with probability $1 - \exp(-k/10)$,

$$\|(\mathbf{W} s^*)_{\mathcal{L}}\| \leq 2\sqrt{kt} \leq \frac{\delta}{100} \|\lambda \langle v, s^* \rangle v\|.$$

By lemma 26,

$$\|v_{\mathcal{L}} \circ z_{s^*}\|^2 \geq (1 - \delta/10) \|v_{\mathcal{L}}\|^2 \geq 1 - \delta.$$

■

Lemma 20 Let $p = \exp(-tr^2/2)$. Then, with probability $1 - d^{-10}$,

$$\max_{s \in \mathcal{S}_t} \left\| \mathbf{W} \circ \left(z_s(r) z_s^\top(r) \right) \right\| \leq 10 \sqrt{(pd + k) \cdot \ln\left(\frac{d}{pd + k}\right)}.$$

Proof Using the same argument as in the proof of lemma 10, we get that the number of nonzero $z_{si}(r)$ for every $s \in \mathcal{S}_t$ is bounded by $2pd + 2\sqrt{pdt \ln(n/t)} + 2k \leq 4pd + 4k$ with probability at least $1 - \exp(-pd - t \ln(d/t))$.

By fact 40, an $m \times m$ Gaussian matrix \mathbf{G} satisfies

$$\|\mathbf{G}\| \leq 2\sqrt{m} + \sqrt{\tau}$$

with probability $1 - \exp(-\tau/2)$ (for every $\tau > 0$). By union bound over all sets of size at most $4pd + 4k$ (corresponding to nonzero rows and columns of $\mathbf{W} \circ (z_s(r) z_s^\top(r))$), we get the desired bound. ■

Lemma 21 Let $0 < \delta < 0.1$ and $\frac{\delta\lambda}{200k} \leq r \leq \frac{\delta\lambda}{100k}$. For $s \in \mathcal{S}_t$ let $\hat{v}(s)$ be the top eigenvector of $\mathbf{Y} \circ (z_s(r) z_s^\top(r))$.

Suppose that $k \gtrsim \frac{t \ln d}{\delta^2}$, and

$$\lambda \gtrsim \frac{k}{\delta^2 \sqrt{t}} \sqrt{\ln\left(2 + \frac{td}{k^2}\right)}.$$

Then there exists $s' \in \mathcal{S}_t$ such that with probability $1 - 2d^{-10}$,

$$|\langle \hat{v}(s'), v \rangle| \geq 1 - 20\delta.$$

Proof Let $s' \in \operatorname{argmax}_{s \in \mathcal{S}_t} \langle s, v \rangle$ and let $\tilde{v} = v \circ z_{s'}(r)$. Then

$$\mathbf{Y} \circ \left(z_s(r) z_s^\top(r) \right) = \beta \|\mathbf{u}\|^2 \tilde{v} \tilde{v}^\top + \mathbf{W} \circ \left(z_s(r) z_s^\top(r) \right),$$

By lemma 20 and the same argument as in lemma 11 with $n = d$, with probability at least $1 - d^{-10}$,

$$\left\| \mathbf{W} \circ \left(\mathbf{z}_s(r) \mathbf{z}_s^\top(r) \right) \right\| \leq \delta \lambda.$$

By lemma 32, $|\langle \hat{\mathbf{v}}(u^*), \tilde{\mathbf{v}} \rangle| \geq 1 - 3\delta$. By lemma 19, with probability $1 - d^{-10}$, $\langle \tilde{\mathbf{v}}, v \rangle \geq 1 - \delta$. Therefore, by fact 33,

$$|\langle \hat{\mathbf{v}}(s'), v \rangle| \geq 1 - 20\delta.$$

■

Proof [Proof of theorem 17] For all $s \in \mathcal{S}_t$ we compute $\hat{\mathbf{v}}(s)$ as in lemma 21 with $\frac{\delta' \lambda}{200k} \leq r \leq \frac{\delta' \lambda}{100k}$, where $\delta' = \delta/1000$. Since we do not know λ , we can create a list of candidates for r of size at most $2 \ln d$ (from $r = 1/d$ to $r = 1$).

By lemma 21, we get a list of vectors \mathbf{L} of size $2 \ln(d) |\mathcal{S}_t|$ such that with probability $1 - d^{-9}$ there exists $\mathbf{v}^* \in \mathbf{L}$ such that $|\langle \mathbf{v}^*, v \rangle| \geq 1 - \delta/10$. By fact 41, k' -sparse norm of \mathbf{W} is bounded by

$$10 \sqrt{k' \ln(ed/k')}$$

with probability at least $1 - d^{-10}$. Hence by lemma 28 with $k' = \lceil 100k/\delta^2 \rceil$, we get the desired estimator. ■

C.2. Adversarial Perturbations

The following theorem is the restatement of theorem 4.

Theorem 22 *Let $d, k, t \in \mathbb{N}$, $\lambda > 0$, $0 < \delta < 0.1$. Let $\tilde{\mathbf{Y}} = \sqrt{\beta} \mathbf{u} \mathbf{v}^\top + \mathbf{W} + E$, where $v \in \mathbb{R}^d$ is a k -sparse unit vector, $\mathbf{W} \sim N(0, 1)^{n \times d}$, and $E \in \mathbb{R}^{d \times d}$ is a matrix with entries*

$$\|E\|_\infty = \varepsilon \lambda / k \lesssim \delta^3 \lambda / k.$$

Suppose that $k \gtrsim \frac{t \ln d}{\delta^2}$ and

$$\lambda \gtrsim \frac{k}{\delta^2} \sqrt{\frac{\ln(2 + td/k^2)}{t}}.$$

Then there exists an algorithm that, given $\tilde{\mathbf{Y}}$, k , t and δ , in time $d^{O(t)}$ outputs a unit vector $\hat{\mathbf{v}}$ such that with probability $1 - o(1)$ as $d \rightarrow \infty$,

$$|\langle \hat{\mathbf{v}}, v \rangle| \geq 1 - \delta.$$

Proof Let $\tilde{\delta} = \delta/1000$ and let $r = \frac{\tilde{\delta}}{100k}$. For $s \in \mathcal{S}_t$ let \tilde{z}_s be n -dimensional (random) vectors defined as

$$\tilde{z}_{si}(r) = \begin{cases} \mathbf{1}_{[\langle s, \tilde{Y}_i \rangle \geq r \cdot t]} & \text{if } s_i = 0 \\ 1 & \text{otherwise} \end{cases}$$

Let $\mathbf{Y} = \tilde{\mathbf{Y}} - E$, and \mathbf{z}_{si} be the same as in the non-adversarial case defined for \mathbf{Y} . Note that

$$\mathbf{z}_{si}(r - \varepsilon \cdot \lambda / k) \leq \tilde{z}_{si}(r) \leq \mathbf{z}_{si}(r + \varepsilon \cdot \lambda / k).$$

Let $s^* \in \operatorname{argmax}_{s \in \mathcal{S}_t} \langle s, v \rangle$. By lemma 19, with probability $1 - d^{-10}$,

$$\|v \circ \tilde{z}_{s^*}\|^2 = \langle v \circ \tilde{z}_{s^*}, v \rangle \geq z_{si}(r - \varepsilon \cdot \lambda/k) \geq 1 - 2\tilde{\delta}.$$

By the argument from the proof of lemma 10, with probability $1 - d^{-10}$, number of nonzero entries of \tilde{z}_s is at most $4pd + 4k$ and

$$\left\| \mathbf{W} \circ \left(\tilde{z}_{s^*} \tilde{z}_{s^*}^\top \right) \right\| \leq 10 \sqrt{(pd + k) \cdot \ln \left(\frac{d}{pd + k} \right)} \lesssim \tilde{\delta}^2 \lambda.$$

Let $\mathcal{P}_k = \{X \in \mathbb{R}^{d \times d} \mid X \succeq 0, \operatorname{Tr} X = 1, \|X\|_1 \leq k\}$. For all $X \in \mathcal{P}_k$,

$$\left| \langle X, \mathbf{W} \circ \left(\tilde{z}_{s^*} \tilde{z}_{s^*}^\top \right) \rangle \right| \leq \left\| \mathbf{W} \circ \left(\tilde{z}_{s^*} \tilde{z}_{s^*}^\top \right) \right\| \cdot \|X\| \lesssim \tilde{\delta} \lambda.$$

and

$$\left| \langle X, E \circ \left(\tilde{z}_{s^*} \tilde{z}_{s^*}^\top \right) \rangle \right| \leq \left\| E \circ \left(\tilde{z}_{s^*} \tilde{z}_{s^*}^\top \right) \right\|_\infty \cdot \|X\|_1 \leq \tilde{\delta} \lambda.$$

Consider

$$\hat{X}(s^*) \in \operatorname{argmax}_{X \in \mathcal{P}} \langle X, \tilde{Y} \circ \left(\tilde{z}_{s^*} \tilde{z}_{s^*}^\top \right) \rangle.$$

By lemma 32,

$$\langle \hat{X}(s^*), vv^\top \circ \left(\tilde{z}_{s^*} \tilde{z}_{s^*}^\top \right) \rangle \geq 1 - 6\tilde{\delta}.$$

Hence by fact 33,

$$\langle \hat{X}(s^*), vv^\top \rangle \geq 1 - 24\tilde{\delta}.$$

Let $\hat{v}(s^*)$ be the top eigenvector of $\hat{X}(s^*)$. By lemma 29,

$$|\langle \hat{v}(s^*), v \rangle| \geq 1 - 100\tilde{\delta}.$$

By fact 41, k' -sparse norm of \mathbf{W} is bounded by

$$10 \sqrt{k' \ln(ed/k')}$$

with probability at least $1 - d^{-10}$. And k' -sparse norm of E is bounded by

$$k' \|E\|_\infty \leq \varepsilon \lambda \frac{k'}{k}.$$

Hence by lemma 28 with $k' = \lceil 100k/\tilde{\delta}^2 \rceil$ and the list of $\mathbf{X}(s)$ for all $s \in \mathcal{S}_t$, we get the desired estimator. \blacksquare

Appendix D. Heavy-tailed Symmetric Noise

The following theorem is a restatement of theorem 5.

Theorem 23 *Let $k, d, t \in \mathbb{N}$, $\lambda > 0$, $A \geq 1$, $0 < \alpha < 1$, $0 < \delta < 0.1$. Let*

$$\mathbf{Y} = \lambda v v^\top + \mathbf{N},$$

where $v \in \mathbb{R}^d$ is a k -sparse unit vector such that $\|v\|_\infty \leq A/\sqrt{k}$ and $\mathbf{N} \sim N(0, 1)^{d \times d}$ is a random matrix with independent (but not necessarily identically distributed) symmetric about zero entries such that

$$\mathbb{P}[|\mathbf{N}_{ij}| \leq 1] \geq \alpha.$$

Suppose that $k \gtrsim \frac{t \ln d}{\delta^4 A^4 \alpha^2}$,

$$t \gtrsim \frac{\ln(2 + td/k^2)}{\alpha^2 A^4 \delta^6},$$

and

$$\lambda \geq k.$$

Then there exists an algorithm that, given \mathbf{Y} , k , t and λ , in time $d^{O(t)}$ finds a unit vector \hat{v} such that with probability $1 - o(1)$ as $d \rightarrow \infty$,

$$|\langle \hat{v}, v \rangle| \geq 1 - O(\delta).$$

Before proving this theorem, we state here a theorem from [d'Orsi et al. \(2022\)](#).

Theorem 24 (d'Orsi et al. (2022)) *Let $\delta, \alpha \in (0, 1)$ and $\zeta \geq 0$. Let $\tilde{\Omega} \subseteq \mathbb{R}^m$ be a compact convex set. Let $b, r, \gamma \in \mathbb{R}$ be such that*

$$\max_{X \in \tilde{\Omega}} \|X\|_\infty \leq b,$$

$$\max_{X \in \tilde{\Omega}} \|X\|_2 \leq r,$$

and

$$\mathbb{E}_{\mathbf{W} \sim N(0, \text{Id})} \left[\sup_{X \in \tilde{\Omega}} \langle X, \mathbf{W} \rangle \right] \leq \gamma.$$

Consider

$$\mathbf{Y} = X^* + \mathbf{N},$$

where $X^* \in \tilde{\Omega}$ and \mathbf{N} is a random m -dimensional vector with independent (but not necessarily identically distributed) symmetric about zero entries satisfying $\mathbb{P}[|\mathbf{N}_i| \leq \zeta] \geq \alpha$.

Then the minimizer $\hat{X} = \operatorname{argmin}_{X \in \tilde{\Omega}} F_h(\mathbf{Y} - X)$ of the Huber loss with parameter $h \geq 2b + \zeta$ satisfies

$$\|\hat{X} - X^*\|_2 \leq O\left(\sqrt{\frac{h}{\alpha}(\gamma + r\sqrt{\log(1/\delta)})}\right)$$

with probability at least $1 - \delta$ over the randomness of \mathbf{N} .

Using this result, we will prove theorem 23

Proof [Proof of theorem 23] Consider

$$\tau_h(x) := \begin{cases} h & \text{if } x > h \\ x & \text{if } |x| \leq h \\ -h & \text{if } x < -h \end{cases}$$

where $h = 3\lambda\|v\|_\infty^2 \leq 3\frac{\lambda}{k}A^2$. Let $\mathbf{T} = \tau_h(\mathbf{Y})$ (i.e. the matrix obtained from Y by applying τ_h to each entry).

For $s \in \mathcal{S}_t$ let \mathbf{z}_s be n -dimensional (random) vectors defined as

$$\mathbf{z}_{s_i}(r) = \begin{cases} \mathbf{1}_{[(\mathbf{T}s)_i = \langle s, \mathbf{T}_i \rangle \geq r \cdot t]} & \text{if } s_i = 0 \\ 1 & \text{otherwise} \end{cases}$$

Let $\frac{\lambda\delta^2\alpha}{10k} \leq r \leq \frac{\lambda\delta^2\alpha}{5k}$.

Let $\mathcal{L} = \left\{ i \in [d] \mid |v_i| \geq \frac{\delta}{\sqrt{k}} \right\}$. Note that $\|v_{\mathcal{L}}\|^2 \geq 1 - \delta^2$ and hence $|\mathcal{L}| \geq \frac{k}{2A^2}$.

Let $s^* \in \operatorname{argmax}_{s \in \mathcal{S}_t} \langle s, v \rangle$. By the same argument as in the proof of lemma 8,

$$\langle s^*, v \rangle = \|v_{\mathcal{L}}\|_1 \geq \delta t / \sqrt{k}.$$

By Chernoff bound, for every $i \in \mathcal{L}$ with probability at least $1 - \exp(-t\alpha/2)$, for at least $\alpha t/10$ of $j \in \operatorname{supp}(s^*(j))$, $|\mathbf{N}_{ij}| \leq 1$. Let \mathcal{C}_i be the set of such entries. Then for $j \in \mathcal{C}_i$,

$$\tau_h(\mathbf{Y}_{ij}) = \lambda v_i v_j + |\mathbf{N}_{ij}|.$$

By Hoeffding's inequality,

$$\sum_{j \in \mathcal{C}_i} \tau_h(\mathbf{Y}_{ij}) \geq \lambda \frac{\delta |\mathcal{C}_i|}{\sqrt{k}} v_i - \sqrt{|\mathcal{C}_i| q}$$

with probability at least $1 - \exp(-q/2)$ Note that for all $j \in \operatorname{supp}(j)$,

$$\mathbb{E} \tau_h(\mathbf{Y}_{ij}) \geq 0.$$

By Hoeffding's inequality,

$$\sum_{j \in \operatorname{supp}(s^*(j)) \setminus \mathcal{C}_i} \tau_h(\mathbf{Y}_{ij}) \geq -h\sqrt{tq}$$

with probability at least $1 - \exp(-q/2)$. Hence for $i \in \mathcal{L}$, with probability at least $1 - \delta^2/A^2$.

$$\langle s^*, \mathbf{T}_i \rangle \geq \frac{\lambda\delta^2\alpha t}{2k} - 10h\sqrt{t \log(A/\delta)} \geq \frac{\lambda\delta^2\alpha t}{4k} \geq r \cdot t.$$

Let $\mu = \delta^2/A^2$. By Chernoff bound, with probability at least $1 - \exp(-\mu k/10)$, for at most $2\mu k$ entries $i \in \mathcal{L}$, $\mathbf{z}_{s^*i} = 0$. Hence with probability at least $1 - \exp(-\mu|\mathcal{L}|/100)$,

$$\|v_{\mathcal{L}} \circ \mathbf{z}_{s^*}\|^2 \geq 1 - 2A^2\mu \geq 1 - \delta^2.$$

By Hoeffding's inequality, for all $i \notin \text{supp}(v)$,

$$|\langle s^*, \mathbf{T}_i \rangle| \leq h\sqrt{tq}$$

with probability at least $1 - \exp(-q/2)$. Hence

$$|\langle s^*, \mathbf{T}_i \rangle| < r \cdot t$$

with probability at least $1 - \exp\left(-\frac{r^2 t}{10h^2}\right) = 1 - p$.

By the same argument as in lemma 10, number of nonzero $z_{s^*i}(r)$ is at most $4pd + 4k$ with probability at least $1 - \exp(-pd - t \ln(d/t))$. For $s \in \mathcal{S}_t$, denote by $\mathbf{Z}(s)$ the set of $i \in [d]$ such that $z_{si}(r) = 1$.

For $Q \subset [d]$ let

$$\mathcal{P}_Q = \{X \in \mathbb{R}^{Q \times Q} \mid X \succeq 0, \text{Tr } X \leq \lambda, \|X\|_1 \leq \lambda k\},$$

Note that

$$\gamma(Q) := \mathbb{E}_{\mathbf{G} \sim N(0,1)^{Q \times Q}} \left[\sup_{X \in \mathcal{P}_Q} \langle X, \mathbf{G} \rangle \right] \leq \lambda \cdot \mathbb{E}_{\mathbf{G} \sim N(0,1)^{Q \times Q}} \|\mathbf{G}\| \leq 10\lambda\sqrt{|Q|}.$$

Hence by theorem 24 and a union bound over all sets Q of size at most $4pd + 4k$, we get with probability at least $1 - d^{-10}$, for all $s \in \mathcal{S}_t$,

$$\left\| \hat{\mathbf{X}}(s) - \lambda \tilde{\mathbf{v}}(s) \tilde{\mathbf{v}}(s)^\top \right\|_{\text{F}}^2 \leq O\left(\frac{h}{\alpha} \cdot \lambda \sqrt{(pd+k) \cdot \ln\left(\frac{d}{pd+k}\right)}\right) \leq O\left(\frac{A^2}{\alpha} \cdot \frac{\lambda^2}{k} \cdot \sqrt{(pd+k) \cdot \ln\left(\frac{d}{pd+k}\right)}\right),$$

where $\hat{\mathbf{X}}(s)$ is the minimizer of the Huber loss with parameter h over $\mathcal{P}_{\mathbf{Z}(s)}$ and $\tilde{\mathbf{v}} = v \circ z_s(r)$.

Note that

$$\sqrt{(pd+k) \cdot \ln\left(\frac{d}{pd+k}\right)} \lesssim \sqrt{pd \ln(1/p)} + \sqrt{k \ln d}.$$

The second term can be bounded as follows

$$\sqrt{k \ln d} \lesssim \delta^2 k \frac{\alpha}{A^2}.$$

Note that

$$r^2 t / h^2 \geq \frac{\delta^4 \alpha^2}{A^4} t \gtrsim \ln(2 + td/k^2).$$

Hence the first term can be bounded as follows

$$\sqrt{pd \ln(1/p)} \lesssim \frac{k}{\sqrt{d}} \cdot \frac{\alpha \delta^2}{A^2} \cdot \sqrt{d} \leq \delta^2 k \frac{\alpha}{A^2}.$$

Hence for all $s \in \mathcal{S}_t$,

$$\left\| \hat{\mathbf{X}}(s) - \lambda \tilde{\mathbf{v}}(s) \tilde{\mathbf{v}}(s)^\top \right\|_{\text{F}}^2 \leq \delta^2 \lambda^2.$$

Since $\|v \circ z_{s^*}\|^2 \geq 1 - 2\delta^2$,

$$\left\| \hat{\mathbf{X}}(s^*) \right\|_{\text{F}} \geq 1 - 10\delta\lambda.$$

Consider some s' such that $\left\| \hat{\mathbf{X}}(s') \right\|_{\text{F}} \geq 1 - 10\delta\lambda$. For such s' ,

$$\langle v \circ z_{s'}(r), v \rangle = \|v \circ z_{s'}(r)\|^2 \geq 1 - 100\delta.$$

Moreover, since $\left\| \hat{\mathbf{X}}(s') - \lambda \tilde{\mathbf{v}}(s') \tilde{\mathbf{v}}(s')^\top \right\|_{\text{F}} \leq \delta\lambda$, the top eigenvector $\hat{\mathbf{v}}(s')$ of $\hat{\mathbf{X}}(s')$ satisfies

$$|\langle \hat{\mathbf{v}}(s'), v \rangle| \geq 1 - O(\delta).$$

■

Appendix E. Properties of sparse vectors

This section contain tools used throughout the rest of the paper.

Fact 25 *Let $r, \delta > 0$ and let $x \in \mathbb{R}^m$ such that $\|x\| \leq R$. Let $\mathcal{S} = \{i \in [m] \mid |x_i| \geq \delta\}$. Then*

$$|\mathcal{S}| \leq R^2/\delta^2.$$

Proof

$$\delta^2 \cdot |\mathcal{S}| \leq \|x_{\mathcal{S}}\|^2 \leq \|x\|^2 \leq r^2.$$

■

Lemma 26 *Let $\delta, \delta' \in (0, 1)$. Let $x, y \in \mathbb{R}^m$ such that $\|y\| \leq \delta\|x\|$. Let $\mathcal{S} = \{i \in [m] \mid |y_i| \geq \delta'|x_i|\}$. Then*

$$\|v_{\mathcal{S}}\| \geq (1 - \delta/\delta')\|v\|.$$

Proof Consider the vector y' such that $y'_i = -x_i$ for all $i \in \mathcal{S}$ and $y'_i = 0$ for all $i \notin \mathcal{S}$. It follows that

$$\|y'\| \leq \frac{1}{\delta'}\|y\| \leq \frac{\delta}{\delta'}\|x\|.$$

Hence

$$\|x_{\mathcal{S}}\| = \|x + y'\| \geq \|x\|_2 - \|y'\| \geq (1 - \delta/\delta')\|x\|.$$

■

Lemma 27 *Let $v \in \mathbb{R}^d$ be a k -sparse unit vector, and suppose that for some unit vector $v' \in \mathbb{R}^d$, $|\langle v, v' \rangle| \geq 1 - \delta$. For $k' \geq k$, let \mathcal{K}' be the set of k' largest (in absolute value) entries of v' . Then*

$$|\langle v, v'_{\mathcal{K}'} \rangle| \geq 1 - \delta - \sqrt{k/k'}.$$

Proof Since $\|v\| = 1$, $\|v'_{\mathcal{K}'} - v'\|_\infty \leq 1/k'$. Hence

$$|\langle v, v'_{\mathcal{K}'} - v' \rangle| \leq \|v'_{\mathcal{K}'} - v'\|_\infty \cdot \|v\|_1 \leq \sqrt{k/k'}.$$

Therefore,

$$|\langle v, v'_{\mathcal{K}'} \rangle| \geq |\langle v, v' \rangle| - |\langle v, v'_{\mathcal{K}'} - v' \rangle| \geq 1 - \delta - \sqrt{k/k'}.$$

■

Lemma 28 *Let $\lambda, \kappa, \delta > 0$ and let $v \in \mathbb{R}^d$ be a k -sparse unit vector. Let $L \subset \mathbb{R}^d$ be a finite set of unit vectors such that*

$$\max_{x \in L} |\langle v, x \rangle| \geq 1 - \delta.$$

Let $k' \geq 2\lceil k/\delta^2 \rceil$ and let $N \in \mathbb{R}^{d \times d}$ be a matrix such that for every k' -sparse unit vector $u \in \mathbb{R}^d$

$$u^\top N u \leq \kappa.$$

Then, there exists an algorithm running in time $O(d^2 \cdot L)$ that, given $Y = \lambda v v^\top + N, L, k$ and δ as input, finds a unit vector \hat{v} such that

$$|\langle v, \hat{v} \rangle| \geq 1 - 4\delta - \frac{2\kappa}{\lambda}.$$

Proof For each $x \in L$ we can compute a k' -sparse vector $s(x)$ that coincides with x on the top k' largest (in absolute value) entries. Let $x^* \in \operatorname{argmax}_{x \in L} |\langle v, s(x) \rangle|$. By lemma 27,

$$|\langle v, s(x^*) \rangle| \geq 1 - 2\delta.$$

Hence

$$s(x^*)^\top (\lambda v v^\top + N) s(x^*) \geq (1 - 2\delta)^2 \lambda - \kappa \geq (1 - 4\delta) \lambda - \kappa.$$

Let

$$\tilde{v} \in \operatorname{argmax}_{s(x), x \in L} s(x)^\top (\lambda v v^\top + N) s(x).$$

Then

$$\tilde{v}^\top (\lambda v v^\top + N) \tilde{v} \geq s(x^*)^\top (\lambda v v^\top + N) s(x^*) \geq (1 - 4\delta) \lambda - \kappa.$$

Since $\tilde{v}^\top N \tilde{v} \leq \kappa$, we get

$$|\langle v, \tilde{v} \rangle| \geq |\langle v, \tilde{v} \rangle|^2 \geq 1 - 4\delta - 2\kappa/\lambda.$$

Hence $\hat{v} = \frac{1}{\|\tilde{v}\|} \tilde{v}$ is the desired estimator.

■

Appendix F. Linear Algebra

Lemma 29 *Let $M \in \mathbb{R}^{d \times d}$, $M \succeq 0$, $\text{Tr } M = 1$ and let $z \in \mathbb{R}^d$ be a unit vector such that $z^\top M z \geq 1 - \varepsilon$. Then the top eigenvector v_1 of M satisfies $\langle v_1, z \rangle^2 \geq 1 - 2\varepsilon$.*

Proof Write $z = \alpha v_1 + \sqrt{1 - \alpha^2} v_\perp$ where v_\perp is a unit vector orthogonal to v_1 .

$$\begin{aligned} z^\top M z &= \alpha^2 v_1^\top M v_1 + (1 - \alpha^2) v_\perp^\top M v_\perp \\ &= \alpha^2 (\lambda_1 - v_\perp^\top M v_\perp) + v_\perp^\top M v_\perp \\ &\geq 1 - \varepsilon \end{aligned}$$

As $v_1^\top M v_1 \geq z^\top M z$ and $v_\perp^\top M v_\perp \leq 1 - v_1^\top M v_1 \leq 1 - z^\top M z \leq \varepsilon$, rearranging

$$\alpha^2 \geq \frac{1 - \varepsilon - v_\perp^\top M v_\perp}{\lambda_1 - v_\perp^\top M v_\perp} \geq 1 - 2\varepsilon. \quad \blacksquare$$

Fact 30 *Let $A, B \in \mathbb{R}^{d \times d}$, $A, B \succeq 0$. Then $\langle A, B \rangle \geq 0$.*

Fact 31 *Let $X \in \mathbb{R}^{d \times d}$ be a positive semidefinite matrix. Then for all $A \in \mathbb{R}^{d \times d}$,*

$$|\langle A, X \rangle| \leq \|A\| \cdot \text{Tr } X.$$

Lemma 32 *Let $\Omega \subset \mathbb{R}^m$. Let $Y = S + N$, where $S \in \Omega$ and $N \in \mathbb{R}^m$ satisfies*

$$\sup_{X \in \Omega} |\langle X, N \rangle| \leq \delta.$$

Then $\hat{X} \in \text{argmax}_{X \in \Omega} \langle X, Y \rangle$ satisfies

$$\langle \hat{X}, S \rangle \geq \|S\|^2 - 2\delta.$$

Proof

$$\langle \hat{X}, S \rangle = \langle \hat{X}, Y \rangle - \langle \hat{X}, N \rangle \geq \langle \hat{X}, Y \rangle - \delta \geq \langle S, Y \rangle - \delta = \langle S, S \rangle + \langle S, N \rangle - \delta \geq \|S\|^2 - 2\delta. \quad \blacksquare$$

Fact 33 *Let $a, b, c \in \mathbb{R}^m$ such that $\|a\| = \|c\| = 1$ and $\|b\| \leq 1$. Suppose that $\langle a, b \rangle \geq 1 - \delta$ and $\langle c, b \rangle \geq 1 - \delta$. Then $\langle a, c \rangle \geq 1 - 4\delta$.*

Proof Note that $\|a - b\|^2 \leq 2 - 2\langle a, b \rangle \leq 2\delta$ and similarly $\|c - b\|^2 \leq 2\delta$. Hence

$$2 - 2\langle a, c \rangle = \|a - c\|^2 \leq (\|a - b\| + \|c - b\|)^2 \leq 8\delta. \quad \blacksquare$$

Lemma 34 Let $\mathcal{P} \subset \mathbb{R}^{d \times d}$ be some set of PSD matrices. For matrix $M \in \mathbb{R}^{n \times d}$ let

$$\mathfrak{s}_{\mathcal{P}}(M) := \sup_{X \in \mathcal{P}} \left| \langle M^{\top} M, X \rangle \right|.$$

Then for arbitrary matrices $A, B \in \mathbb{R}^{n \times d}$,

$$\sup_{X \in \mathcal{P}} \left| \langle B^{\top} A + A^{\top} B, X \rangle \right| \leq 2\sqrt{\mathfrak{s}_{\mathcal{P}}(A) \cdot \mathfrak{s}_{\mathcal{P}}(B)}.$$

Proof Let X' be an arbitrary element of \mathcal{S} . For some $c \in \mathbb{R}$ (we will choose the value of c later), let $C = (A - cB)^{\top} (A - cB)$. Since C and X' are PSD, $\langle C, X' \rangle \geq 0$. Hence

$$\left| c \cdot \langle B^{\top} A + A^{\top} B, X' \rangle \right| \leq \langle A^{\top} A, X' \rangle + c^2 \langle B^{\top} B, X' \rangle \leq \mathfrak{s}_{\mathcal{P}}(A) + c^2 \mathfrak{s}_{\mathcal{P}}(B).$$

Therefore,

$$\left| \langle B^{\top} A + A^{\top} B, X' \rangle \right| \leq \frac{\mathfrak{s}_{\mathcal{P}}(A)}{|c|} + |c| \cdot \mathfrak{s}_{\mathcal{P}}(B).$$

To minimize this expression, we can take $|c| = \sqrt{\frac{\mathfrak{s}_{\mathcal{P}}(A)}{\mathfrak{s}_{\mathcal{P}}(B)}}$. Hence

$$\left| \langle B^{\top} A + A^{\top} B, X' \rangle \right| \leq 2\sqrt{\mathfrak{s}_{\mathcal{P}}(A) \cdot \mathfrak{s}_{\mathcal{P}}(B)}.$$

Since it holds for arbitrary $X' \in \mathcal{S}$, we get the desired bound. ■

Appendix G. Concentration Inequalities

Fact 35 (Chernoff's inequality, Vershynin (2018)) Let ζ_1, \dots, ζ_n be independent Bernoulli random variables such that $\mathbb{P}(\zeta_i = 1) = \mathbb{P}(\zeta_i = 0) = p$. Then for every $\Delta > 0$,

$$\mathbb{P}\left(\sum_{i=1}^n \zeta_i \geq pn(1 + \Delta)\right) \leq \left(\frac{e^{-\Delta}}{(1 + \Delta)^{1+\Delta}}\right)^{pn}.$$

and for every $\Delta \in (0, 1)$,

$$\mathbb{P}\left(\sum_{i=1}^n \zeta_i \leq pn(1 - \Delta)\right) \leq \left(\frac{e^{-\Delta}}{(1 - \Delta)^{1-\Delta}}\right)^{pn}.$$

Fact 36 (Hoeffding's inequality, Wainwright (2019)) Let z_1, \dots, z_n be mutually independent random variables such that for each $i \in [n]$, z_i is supported on $[-c_i, c_i]$ for some $c_i \geq 0$. Then for all $t \geq 0$,

$$\mathbb{P}\left(\left|\sum_{i=1}^n (z_i - \mathbb{E} z_i)\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2 \sum_{i=1}^n c_i^2}\right).$$

Fact 37 (Bernstein's inequality Wainwright (2019)) Let z_1, \dots, z_n be mutually independent random variables such that for each $i \in [n]$, z_i is supported on $[-B, B]$ for some $B \geq 0$. Then for all $t \geq 0$,

$$\mathbb{P}\left(\sum_{i=1}^n (z_i - \mathbb{E} z_i) \geq t\right) \leq \exp\left(-\frac{t^2}{2 \sum_{i=1}^n \mathbb{E} z_i^2 + \frac{2Bt}{3}}\right).$$

Fact 38 *Wainwright (2019)* Let $X \sim N(0, \sigma^2)$, then for all $t > 0$,

$$\mathbb{P}(X \geq \sigma \cdot t) \leq e^{-t^2/2}.$$

Fact 39 *Laurent and Massart (2000)* Let $X \sim \chi_m^2$, then for all $x > 0$,

$$\begin{aligned} \mathbb{P}(X - m \geq 2x + 2\sqrt{mx}) &\leq e^{-x} \\ \mathbb{P}(m - X \geq x) &\leq e^{-\frac{x^2}{4m}} \end{aligned}$$

Fact 40 *Wainwright (2019)* Let $W \sim N(0, 1)^{n \times d}$. Then with probability $1 - \exp(-t/2)$,

$$\|W\| \leq \sqrt{n} + \sqrt{d} + \sqrt{t}$$

and

$$\left\| W^\top W - n\text{Id} \right\| \leq d + 2\sqrt{dn} + t + 4\sqrt{t(n+d)}.$$

Fact 41 *d'Orsi et al. (2020)* Let $\mathbf{W} \sim N(0, 1)^{n \times d}$ be a Gaussian matrix. Let $1 \leq k \leq d$. Then with probability at least $1 - \left(\frac{k}{ed}\right)^k$

$$\max_{\substack{u \in \mathbb{R}^n \\ \|u\|=1}} \max_{\substack{k\text{-sparse } v \in \mathbb{R}^d \\ \|v\|=1}} u^\top \mathbf{W} v \leq \sqrt{n} + 3\sqrt{k \ln\left(\frac{ed}{k}\right)}$$

and

$$\max_{\substack{k\text{-sparse } v \in \mathbb{R}^d \\ \|v\|=1}} v^\top \mathbf{W}^\top \mathbf{W} v - n \leq 10\sqrt{kn \ln(ed/k)} + 10k \ln(ed/k).$$

Lemma 42 Let \mathbf{Y} be an instance of sparse PCA in Wishart model. For $m \in \mathbb{N}$ let $Z_m = \{z \in \{0, 1\}^d \mid \|z\|_1 \leq m\}$. Suppose that $m \geq 100 \ln d$ and $n \geq 0.1 \cdot m \ln(ed/m)$. Then, with probability at least $1 - d^{-10}$,

$$\max_{z \in Z_m} \left\| \left(\mathbf{Y}^\top \mathbf{Y} - n\text{Id} - \beta \|\mathbf{u}\|^2 v v^\top \right) \circ (z z^\top) \right\| \leq 10\sqrt{(n + \beta n) \cdot m \ln(ed/m)} + 10m \ln(ed/m).$$

Proof We can write $(\mathbf{Y}^\top \mathbf{Y} - n\text{Id} - \beta \|\mathbf{u}\|^2 v v^\top) \circ (z z^\top)$ as

$$\left(\mathbf{W}^\top \mathbf{W} - n \cdot \text{Id} + \sqrt{\beta} \mathbf{W}^\top \mathbf{u} v^\top + \sqrt{\beta} v \mathbf{u}^\top \mathbf{W} \right) \circ (z z^\top).$$

Note that

$$\left\| \left(\mathbf{W}^\top \mathbf{u} v^\top + v \mathbf{u}^\top \mathbf{W} \right) \circ (z z^\top) \right\| \leq 2 \left\| \left(\frac{1}{\|\mathbf{u}\|} \mathbf{u}^\top \mathbf{W} \right) \circ z \right\| \cdot \|\mathbf{u}\|.$$

With probability at least $1 - \exp(-n/10)$, $\|\mathbf{u}\| \leq \sqrt{2n}$. By lemma 43, with probability at least $1 - \exp(-m)$,

$$\left\| \left(\mathbf{W}^\top \mathbf{u} v^\top + v \mathbf{u}^\top \mathbf{W} \right) \circ (z z^\top) \right\| \leq 5\sqrt{\beta n} \cdot \sqrt{m \ln(em/k)}.$$

By fact 40, for every $m' \in \mathbb{N}$, $\mathbf{G} \sim N(0, 1)^{n \times m'}$ satisfies

$$\left\| \mathbf{G}^\top \mathbf{G} - n \cdot \text{Id} \right\| \leq m' + 2\sqrt{m'n} + \tau + 4\sqrt{\tau(m' + n)}.$$

with probability $1 - \exp(-\tau/2)$ (for every $\tau > 0$). By union bound over all sets of size $m' \leq m$ we get the desired bound. \blacksquare

Lemma 43 *Let $\mathbf{W} \sim N(0, 1)^{n \times d}$ and let $a \in \mathbb{R}^n$ and $b \in \mathbb{R}^d$ be vectors independent of \mathbf{W} . For $m \in \mathbb{N}$ let $Z_m = \left\{ z \in \{0, 1\}^d \mid \|z\|_1 \leq m \right\}$. Suppose that $m \geq 100 \ln d$. Then with probability at least $1 - d^{-10}$,*

$$\max_{z \in Z_m} \left\| \left(b a^\top \mathbf{W} \right) \circ \left(z z^\top \right) \right\| \leq 3 \cdot \|a\| \cdot \|b\| \cdot \sqrt{m \ln(em/k)}.$$

Proof For fixed set $J \subset [d]$, let $\mathbf{1}_J \in \{0, 1\}^d$ be an indicator vector of this set. Random vector $\left(\frac{1}{\|a\|} a^\top \mathbf{W} \right) \circ \mathbf{1}_J$ has standard $|J|$ -dimensional Gaussian distribution, hence by fact 40, for all $\tau > 0$, with probability at least $1 - \exp(-\tau/2)$,

$$\left\| \left(\frac{1}{\|a\|} a^\top \mathbf{W} \right) \circ \mathbf{1}_J \right\| \leq \sqrt{|J|} + 1 + \sqrt{\tau}.$$

By union bound over all sets of size at most m , we get

$$\max_{z \in Z_m} \left\| \left(\frac{1}{\|a\|} a^\top \mathbf{W} \right) \circ z \right\| \leq 3\sqrt{m \cdot \ln(ed/m)}$$

with probability at least $1 - \exp(-m)$. \blacksquare