

Second order accurate distributed eigenvector computation for extremely large matrices

Noureddine El Karoui*

Department of Statistics, U.C. Berkeley, Berkeley, CA 94720
e-mail: nkaroui@stat.berkeley.edu

and

Alexandre d’Aspremont[†]

ORFE, Princeton University, Princeton, NJ 08544
e-mail: aspremon@princeton.edu

Abstract: We propose a second-order accurate method to estimate the eigenvectors of extremely large matrices thereby addressing a problem of relevance to statisticians working in the analysis of very large datasets. More specifically, we show that averaging eigenvectors of randomly subsampled matrices efficiently approximates the true eigenvectors of the original matrix under certain conditions on the incoherence of the spectral decomposition. This incoherence assumption is typically milder than those made in matrix completion and allows eigenvectors to be sparse. We discuss applications to spectral methods in dimensionality reduction and information retrieval.

Received February 2010.

1. Introduction

Spectral methods have a long list of applications in statistics and machine learning. Beyond dimensionality reduction techniques such as PCA or CCA [And03, MKB79], they have been used in clustering [NJW02], ranking & information retrieval [PBMW98, HTF⁺01, LM05] or classification for example. Computationally, one of the most attractive features of these methods is their low numerical cost, in particular on problems where the data matrix is sparse (e.g. graph clustering or information retrieval). Computing a few leading eigenvalues and eigenvectors of a matrix, using the power or Lanczos methods for example, requires performing a sequence of matrix vector products and can be

*Support from an Alfred P. Sloan research Fellowship and NSF grants DMS-0605169 and DMS-0847647 (CAREER) is gratefully acknowledged.

[†]Support from NSF grants DMS-0625352, SES-0835550 (CDI), CMMI-0844795 (CAREER), a Peek junior faculty fellowship, a Howard B. Wentz Jr. award and a gift from Google is gratefully acknowledged.

processed very efficiently. This means that when the matrix is dense and has dimension n , the cost of each iteration is $O(n^2)$ in both storage and flops.

However, for extremely large scale problems arising in statistics or information retrieval for example, this cost quickly becomes prohibitively high and makes spectral methods impractical. In this paper, we propose a randomized, distributed algorithm to estimate eigenvectors (and eigenvalues) which makes spectral methods tractable on very large scale matrices. We show that our method is second order accurate and illustrate its performance on a few realistic datasets.

Going back to the numerical cost of spectral methods, we see that decomposing each matrix vector product in many smaller block operations partially alleviates the complexity problem, but makes the overall process very bandwidth intensive. Decomposition techniques thus improve the *granularity* of iterative eigenvalue methods (i.e. require many cheaper operations instead of a single very expensive one), but at the expense of significantly higher bandwidth requirements. Here, we focus on methods that improve the granularity of large-scale eigenvalue computations while having *very low bandwidth* requirements, meaning that they can be fully distributed over many loosely connected machines.

The idea of using subsampling to lower the complexity of spectral methods can be traced back at least to [GMKG91, PRTV00] who described algorithms based on subsampling and random projections respectively. Explicit error estimates followed in [FKV04, DKM06, AM07] which bounded the approximation error of either elementwise or columnwise matrix subsampling procedures. On the application side, a lot of work has been focused on the Pagerank vector, and [NZJ01] in particular study its stability under perturbations of the network matrix. Similar techniques are applied to spectral clustering in [HYJT08] and both works have close connections to ours. Following the *Netflix* competition on collaborative filtering, a more recent stream of works [RFP07, CR08, CT09, KMO09] has also been focused on *exactly* reconstructing a low rank matrix from a small, single incoherent set of observations. Finally, more recent “volume sampling” results provide relative error bounds [KV09], but so far, the sampling probabilities required to obtain these improved error bounds remain combinatorially hard to compute.

Our work here is focused on the impact of subsampling on eigenvector approximations. First we seek to understand how far we can reduce the granularity of eigenvalue methods using subsampling, before reconstructing eigenvectors becomes impossible. This question was partially answered in [CT09, KMO09] for matrices with low rank, incoherent spectrum, using a *single* subset of matrix coefficients, after solving a convex program with *high complexity*. Here we make much milder assumptions on matrix incoherence. In particular, we allow some eigenvectors to be *sparse* (while remaining incoherent on their support) and we approximate eigenvectors using *many* simple operations on subsampled matrices. Under certain conditions on the sampling rate which guarantee that we remain in a perturbative setting, we show that simply *averaging* many approximate eigenvectors obtained by subsampling reduces approximation error by an

order of magnitude. We also show on real data that this technique results in practice in significant improvement in the quality of the approximations.

At a technical level, our approach is to represent the subsampled matrices as additive perturbations of the original matrix. We present in Theorem 2 deterministic and non-asymptotic bounds that allow us to approximate the perturbed eigenvectors to any order by fairly explicit functions of the original matrix as well as the perturbation matrix. Precise use of these bounds then yields second order accurate approximation results in expectation, which is what we seek, given the averaging procedure we propose. The higher-order analysis we perform is necessary to show that even though the first-order error terms (in probability) have mean 0, we can take expectation and not have to worry about higher-order terms (in probability) exploding in expectation. These bounds also give us a very precise understanding of the perturbed eigenvectors, which is of independent interest.

Non-elementary random matrix theory plays a key role in allowing us to bound the norm of various random matrices appearing in our computations. Some non-trivial bounds on the norm of various random matrices appear in Appendix A, improving for instance (and in certain situations) on the results of [AM07]. Concentration of measure arguments also play an important role, allowing us to essentially shift the questions of bounding the norm a certain random matrix to controlling the mean (or median) of this norm. We make repeated use of Talagrand's inequality [Tal95] and of its consequences detailed for instance in [Led01], Chap. 4.

A simple take-away message from our analysis is that when the incoherence conditions we propose are met and when all eigenvectors have support of size n , one can sample the corresponding large matrix at rate $p = (\log n)/n$ (or larger) and still be able to approximate the eigenvectors of the corresponding matrix well. We also show that our approximations run into trouble if $p = (\log n)^{1-\delta}/n$, for some $\delta > 0$, so our results seem sharp in terms of sampling rates.

Notation. In what follows, we write \mathbf{S}_n the set of symmetric matrices of dimension n . For a matrix $X \in \mathbf{R}^{m \times n}$, we write $\|X\|_F$ its Frobenius norm, $\|X\|_2$ its spectral norm, $\sigma_i(X)$ its i -th largest singular value and let $\|X\|_\infty = \max_{ij} |X_{ij}|$, while $\mathbf{Card}(X)$ is the number of nonzero coefficients in X . We denote by $X(i, j)$ or X_{ij} its (i, j) -th element and by M_i the i -th column of M . Here, \circ denotes the Hadamard (i.e entrywise) product of matrices. When $x \in \mathbf{R}^n$ is a vector, we write its Euclidean norm $\|x\|_2$ and $\|x\|_\infty$ its ℓ_∞ norm. We write $\mathbf{1} \in \mathbf{R}^n$ the vector having all entries equal to 1. Finally, κ denotes a generic constant, whose value may change from display to display.

2. Subsampling

We first recall the subsampling procedure in [AM07] which approximates a symmetric matrix $M \in \mathbf{S}_n$ using a subset of its coefficients. The entries of M are

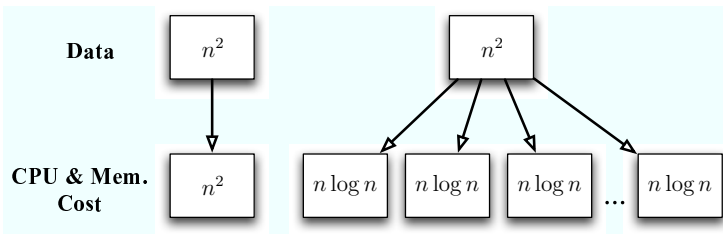


FIG 1. Our objective here is to approximate the spectral decomposition problem of size $O(n^2)$ by solving many independent problems of much smaller size.

independently sampled as

$$S_{ij} = \begin{cases} M_{ij}/p & \text{with probability } p \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

where $p \in [0, 1]$ is the sampling probability. Theorem 1.4 in [AM07] shows that when n is large enough

$$\|M - S\|_2 \leq 4\|M\|_\infty \sqrt{n/p}, \quad (2)$$

holds with high probability (if p is large enough). In what follows, we will prove a similar bound on $\|M - S\|_2$ using incoherence conditions on the spectral decomposition of M .

2.1. Computational benefits

Computing k leading eigenvectors and eigenvalues of a symmetric matrix of dimension n using iterative algorithms such as the power or Lanczos methods (see [GVL90, Chap. 8-9] for example) only requires matrix vector products, hence can be performed in $O(kn^2)$ flops when the matrix is dense. However, this cost is reduced to $O(k \text{Card}(M))$ flops for sparse matrices M . Because the matrix S defined in (1) has only pn^2 nonzero coefficients on average, the cost of computing k leading eigenvalues/eigenvectors of S will typically be $1/p$ times smaller than that of performing the same task on the full matrix M . Of course, sampling the matrix S still requires $O(n^2)$ flops, but can be done in a single pass over the data and be fully distributed. In what follows, we will show that, under incoherence conditions, averaging the eigenvectors of many independently subsampled matrices produces second order accurate approximations of the original spectral decomposition. While the global computational cost of this averaging procedure may not be globally lower, it is decomposed into many much smaller computations, and is thus particularly well adapted to large clusters of simple, loosely connected machines (Amazon EC2, Hadoop, etc.).

2.2. Sparse matrix approximations

Let us write the spectral decomposition of $M \in \mathbf{S}_n$ as

$$M = \sum_{i=1}^n \lambda_i u_i u_i^T$$

where $u_i \in \mathbf{R}^n$ for $i = 1, \dots, n$ and $\lambda \in \mathbf{R}^n$ are the eigenvalues of M with $\lambda_1 > \dots > \lambda_n$ (we assume they are all distinct). Let $\alpha \in [0, 1]^n$, we measure the *incoherence* of the matrix M as

$$\mu(M, \alpha) = \sum_{i=1}^n |\lambda_i| n^{\alpha_i} \|u_i\|_\infty^2 \tag{3}$$

Note that this definition is slightly different from that used in [CT09] because we do not seek to reconstruct the matrix M exactly, so the tail of the spectrum can be partially neglected in our case. In a uniformly bounded model where $n\|u_i\|_\infty^2 = O(1)$, the results of [CT09, §1.5.1] guarantee exact reconstruction of the matrix M given only a fraction of its entries by solving a semidefinite program. As we will see below, the fact that we only seek approximate eigenvectors here, instead of exact reconstructions, allows us to relax these requirements and handle sparse eigenvectors.

Let us define a matrix $Q \in \mathbf{S}_n$ with i.i.d. Bernoulli coefficients

$$Q_{ij} = \begin{cases} 1/p & \text{with probability } p \\ 0 & \text{otherwise.} \end{cases}$$

We can write

$$Q = \mathbf{1}\mathbf{1}^T + \sqrt{\frac{1-p}{p}} C$$

where C has i.i.d. entries with mean zero and variance one, defined as

$$C_{ij} = \begin{cases} \sqrt{(1-p)/p} & \text{with probability } p \\ -\sqrt{p/(1-p)} & \text{otherwise.} \end{cases}$$

We can now write the sampled matrix S in (1) as

$$S = M \circ Q = M + \sqrt{\frac{1-p}{p}} \left(\sum_{i=1}^n \lambda_i (u_i u_i^T) \circ C \right) \equiv M + E \tag{4}$$

and we now seek to bound the spectral norm of the residual matrix E as n goes to infinity. Naturally, if $\|E\|_2$ is small, S is a good approximation of M in spectral terms, because of Weyl's inequality and the Davis-Kahan $\sin(\theta)$ -theorem (see [Bha97]). So our aim now is to control $\|E\|_2$ so we can guarantee the quality of spectral approximations of M made using the sparse matrix S which is computationally easier to work with than the dense matrix M . We now make the following key assumptions on the incoherence of the matrix M .

Assumption 1. *There is a sequence of vectors $\alpha^{(n)} \in [0, 1]^n$ for which*

$$\mu(M, \alpha^{(n)}) \leq \mu \quad \text{and} \quad \mathbf{Card}(u_i) \leq n^{\alpha_i^{(n)}}, \quad i = 1, \dots, n$$

as n goes to infinity, where μ is an absolute constant.

In what follows, we will drop the dependence of α on n to make the notation less cumbersome, so instead of writing $\alpha^{(n)}$ we will just write α . We have the following theorem.

Theorem 1. *Suppose that Assumption 1 holds. Let us call $\alpha_{\min} = \min_{1 \leq i \leq n} \alpha_i$.*

1. *Assume that p and n are such that, $p < 1/2$, and for a given $\delta > 0$, $\alpha_{\min} > (\log n)^{(\delta-3)/4}$ and*

$$\frac{(\alpha_{\min} \log n)^4}{pn^{\alpha_{\min}}} \rightarrow 0, \quad \text{as } n \rightarrow \infty,$$

then we have

$$\limsup_{n \rightarrow \infty} \|E\|_2 \leq 2\mu (pn^{\alpha_{\min}})^{-1/2} \quad \text{a.s.} \quad (5)$$

2. *Further, if $pn^{\alpha_{\min}}/(\log n^{\alpha_{\min}})$ is bounded below by $\gamma > 0$, α_{\min} is such that $n^{\alpha_{\min}} \rightarrow \infty$ and $p < 1/2$, we have, for some finite $\mathfrak{K}(\gamma)$,*

$$\limsup_{n \rightarrow \infty} \|E\|_2 \leq \mathfrak{K}(\gamma)\mu (pn^{\alpha_{\min}})^{-1/2} \quad \text{a.s.} \quad (6)$$

Moreover, $\mathfrak{K}(\gamma)$ is of the form $(1+2/\sqrt{\gamma})\mathfrak{K}+8/\sqrt{\gamma\alpha_{\min}}$ for some universal \mathfrak{K} . Naturally, if

$$\liminf pn^{\alpha_{\min}}/(\log n^{\alpha_{\min}}) = \infty \quad \text{and} \quad \gamma\alpha_{\min} \rightarrow \infty,$$

then $\mathfrak{K}(\gamma)$ can be replaced by \mathfrak{K} .

Proof. Using [HJ91, Th. 5.5.19] or the fact that $uu^T \circ C = D_u C D_u$, where D_u is a diagonal matrix with the vector u on the diagonal (remember that $\|\cdot\|_2$ is a matrix norm and hence sub-multiplicative), we get

$$\|E\|_2 = \sqrt{\frac{1-p}{p}} \left\| \sum_{i=1}^n \lambda_i C \circ (u_i u_i^T) \right\|_2 \leq \sqrt{\frac{1-p}{p}} \sum_{i=1}^n |\lambda_i| n^{\alpha_i/2} \|u_i\|_\infty^2 \left\| \frac{C_{\alpha_i}}{n^{\alpha_i/2}} \right\|_2. \quad (7)$$

Since we assume that the vector u_i is sparse with $\mathbf{Card}(u_i) \leq n^{\alpha_i}$, C_{α_i} is a principal submatrix of C with dimension n^{α_i} . Now, we show in Theorem B.1 (this is the key element of the proof – see p.1375) that

$$\limsup_{n \rightarrow \infty} \max_{i \in \{1, \dots, n\}} \left\| \frac{C_{\alpha_i}}{n^{\alpha_i/2}} \right\|_2 \leq 2 \quad \text{a.s.} \quad ,$$

whenever $p = o\left(\frac{(\alpha_{\min} \log n)^4}{n^{\alpha_{\min}}}\right)$, and $\alpha_{\min} > (\log n)^{(\delta-3)/4}$ for some $\delta > 0$. (Our proof of Theorem B.1 relies on a result of Vu [Vu07] and Talagrand's inequality.) This yields Equation (5) and concludes that part of the proof.

The second part of the proof relies on non trivial results that allow us to control $\|C_{\alpha_i}\|_2$ even when $p \geq \gamma(\log n^{\alpha_i})/n^{\alpha_i}$. These results are given in full details in Appendix A. They allow us to conclude in Theorem B.1 that in the second setting,

$$\limsup_{n \rightarrow \infty} \max_{i \in \{1, \dots, n\}} \left\| \frac{C_{\alpha_i}}{n^{\alpha_i/2}} \right\|_2 \leq \mathfrak{K}(\gamma) \text{ a.s. .}$$

This yields the desired results. □

The proof of the theorem makes clear that the error term coming from the sparsest eigenvector will usually dominate all the others in the residual matrix E .

In these approximation methods, we naturally want to use a small p , so that S is very sparse and the computation of its spectral decomposition is numerically cheap. The result of Theorem B.2 guarantees that the subsampling approximation works whenever $p \gg (\alpha_{\min} \log n)^4/n^{\alpha_{\min}}$ (asymptotically, but we have in mind a very high-dimensional setting, so n will be large in practice) and in that setting, we get a bound of 2 for the norm of the relevant error matrix. When $p \geq \gamma(\log n^{\alpha_{\min}})/n^{\alpha_{\min}}$, where γ stays bounded away from 0, we just know that the norm of the relevant error matrix is bounded, but do not have an explicit value for the constant.

A natural question is therefore whether we could use p much smaller than this. Separate computations (see Subsection A.3) indicate that $\|C/n^{1/2}\|_2$ goes to infinity if $p \leq (\log n)^{1-\delta}/n$, which suggests that this subsampling approach to approximating eigenproperties of M might run into trouble if the sampling rate p gets much smaller than $\log n/n$. (Note also that our results are therefore sharp in terms of rates.) As a matter of fact, we could not control the quantities $\|C_{\alpha_i}/n^{\alpha_i/2}\|_2$ at this sampling rate, which is naturally problematic given the way we established the bound on $\|E\|_2$. Furthermore, if the sparsest eigenvector had support disjoint from the supports of all other eigenvectors, E would be the sum of two block diagonal matrices. Hence, its operator norm would be the maximum of the operator norms of the two blocks, at least one of which having potentially very large operator norm.

2.3. Tightness

Note that, in the limit case $\alpha = \mathbf{1}$ where the eigenvectors are fully dense and incoherent, our bound is similar to the original bound in [AM07, Theorem 1.4] or that of [KMO09, Th 1.1] (our model for M is completely different however). In fact, the bounds in (2) and (5) can be directly compared. In the fully dense case where $\alpha = \mathbf{1}$, we have

$$\sqrt{n}\|M\|_{\infty} = \sqrt{n} \left\| \sum_{i=1}^n \lambda_i u_i u_i^T \right\|_{\infty} \leq n^{-1/2} \sum_{i=1}^n |\lambda_i| n \|u_i\|_{\infty}^2 \leq n^{-1/2} \mu,$$

so in this limit case, the original bound in (2) is always tighter than our bound in (5). However, in the sparse incoherent case where $\alpha \neq \mathbf{1}$, the ratio of the

bound (2) in [AM07] over our bound (5) becomes

$$\frac{2 \left\| \sum_{i=1}^n \lambda_i n^{\frac{(\alpha_{\min}+1)}{2}} u_i u_i^T \right\|_{\infty}}{\sum_{i=1}^n |\lambda_i| n^{\alpha_i} \|u_i\|_{\infty}^2},$$

which can be large when $\alpha_{\min} < 1$. The results in [KMO09], which are focused on exact recovery of low rank incoherent matrices, do not apply when the eigenvectors are sparse (i.e. $\alpha \neq \mathbf{1}$).

2.4. Approximating eigenvectors

We now study the impact of subsampling on the eigenvectors and in particular on the one associated with the largest eigenvalue. We have the following theorem.

Theorem 2. *Assume that the eigenvalues of M are simple. When $S = M + E$, let us call $v_k \in \mathbf{R}^n$ and $\lambda_k(S)$ the k -th eigenpair of S , and $u_k \in \mathbf{R}^n$, λ_k the k -th eigenpair of M . We write R_k the reduced resolvent of M associated with u_k , defined as*

$$R_k = \sum_{j \neq k} \frac{1}{\lambda_j - \lambda_k} u_j u_j^T,$$

and let $\Delta_k = R_k(E - (\lambda_k(S) - \lambda_k)Id)$. We also call d_k the separation distance of λ_k , i.e the distance from λ_k to the nearest eigenvalue of M . If $\|E\|_2$ satisfies $\|E\|_2 < d_k/2$, then

$$\left\| v_k - u_k + \left[\sum_{m=0}^j (-1)^m \Delta_k^m \right] R_k E u_k \right\|_2 \leq \frac{1}{2} \left(\frac{2 \|E\|_2}{d} \right)^{j+2} \frac{1}{1 - \frac{2 \|E\|_2}{d}} \quad (8)$$

having normalized v_k so $v_k^T u_k = 1$.

As will be seen shortly, to prove this theorem we find an explicit and exact representation of $v_k - u_k$, and show that when we subtract from this quantity a $(j + 1)$ -term approximation, which is also explicit, we are left with an error term of order $j + 2$. The use of reduced resolvents in this setting is natural in the analytic perturbation theory of linear operators, which underlies our approach.

One virtue of this approximation is that in the settings corresponding to our original problem, the term of highest magnitude, let us call it A_0 , has mean 0. Hence, at least conceptually, when trying to bound $\|\mathbf{E}[v - u]\|_2$ we can simply use $\|\mathbf{E}[v - u]_2\| = \|\mathbf{E}[v - u - A_0]_2\|$ and finally

$$\|\mathbf{E}[v - u]\|_2 = \|\mathbf{E}[v - u - A_0]\|_2 \leq \mathbf{E}[\|v - u - A_0\|_2].$$

We will see in Theorem 3 that further technical problems arise (which force us to go an order higher in approximations), but that is essentially the idea and a main motivation for getting the present result.

We also note that the theorem gives us a very precise understanding of $v - u$, essentially to any order.

Proof. From now on we focus on u_k and drop the dependence on k in u_k, v_k, R_k, Δ_k etc. . . when this does not create confusion. We also use the notation λ_S and λ instead of $\lambda_k(S)$ and λ_k . If v is normalized so that $v^T u = 1$ (so $(v - u)^T u = 0$), we have the explicit formula [Kat95, Eq. 3.29]

$$v - u = -(\text{Id} + R(E - \gamma \text{Id}))^{-1} R E u ,$$

where $\gamma = \lambda_S - \lambda$. The formula is valid as soon as $(\text{Id} + R(E - \gamma \text{Id}))$ is invertible. Let us now call $\Delta = R(E - \gamma \text{Id})$ and assume that Δ has no eigenvalues equal to -1, i.e $\text{Id} + \Delta$ is invertible. Then we have

$$v - u + \left[\sum_{m=0}^j (-1)^m \Delta^m \right] R E u = (-1)^j \Delta^{j+1} (\text{Id} + \Delta)^{-1} R E u . \tag{9}$$

We also have by construction $Ru = 0$, so $R E u = \Delta u$. Hence, we can write

$$v - u + \left[\sum_{m=0}^j (-1)^m \Delta^m \right] R E u = (-1)^j \Delta^{j+2} (\text{Id} + \Delta)^{-1} u .$$

Now let us call d the separation distance of λ . Then $\|R\|_2 = 1/d$. Our assumptions guarantee that $\|E\|_2$ is such that $2 \|E\|_2 / d < 1$. We note that using Weyl's inequality, $|\lambda_S - \lambda| \leq \|S - M\|_2 = \|E\|_2$, hence $\|\Delta\| \leq 2 \|R\|_2 \|E\|_2 = 2 \|E\|_2 / d$ and

$$\|(\text{Id} + \Delta)^{-1}\|_2 \leq \frac{1}{1 - \frac{2\|E\|_2}{d}} .$$

Putting all the elements together and recalling that $\|u\|_2 = 1$, we get (8) from Equation (9). □

Spectral methods tend to focus on eigenvectors associated with extremal eigenvalues, so let us elaborate on the meaning of Theorem 2 for the eigenvector associated with the largest eigenvalue. If we suppose that the spectral norm of the residual matrix E is smaller than half the separation distance of the largest eigenvalue, i.e

$$\|E\|_2 < (\lambda_1 - \lambda_2) / 2 , \tag{10}$$

the previous result (and results such as [Kat95, Theorem II.3.9]) shows that we can use perturbation expansions to approximate the leading eigenvector of the subsampled matrix. Based on the bound in Equation (5), the condition stated in Equation (10) will be satisfied (asymptotically with high-probability) if, for some $\varepsilon > 0$,

$$\frac{\mu}{\sqrt{pn^{\alpha_{\min}}}} < (\lambda_1 - \lambda_2) / (4 + \varepsilon) .$$

We note that assumption (10) is likely reasonable if one eigenvalue is very large compared to the others, which is a natural setting for methods such as PCA. (Note however that our result is not limited to the largest eigenvalue but actually applies to any eigenvalue of the original matrix M , λ , for which $\|E\|_2$ is smaller

than half the distance from λ to any other eigenvalue of M . In particular, the result would apply to several separated eigenvalues.) We also note that the approximation

$$v = u - \left[\sum_{m=0}^j (-1)^m \Delta^m \right] REu$$

is accurate to order $j + 2$.

Let us now try to make our approximation slightly more explicit. If we write R the reduced resolvent of M (associated with u_1), and assume that $\lambda_1 - \lambda_2$ stays bounded away from 0, we have in this setting, using Equation (8) with $j = 1$,

$$v = u - REu + R(E - (\lambda_1(S) - \lambda_1) \text{Id})REu + O_P(\|E\|_2^3),$$

and therefore

$$v = u - REu + R(E - u^T E u \text{Id})REu + O_P(\|E\|_2^3), \tag{11}$$

after we account for the fact that $u^T E u$ is an order- $\|E\|_2^2$ accurate approximation of $\lambda_1(S) - \lambda_1$ [Kat95, Eq. 2.36 and 3.18]. This approximation makes clear that a key component in the accuracy of our approximations will be the size of the vector Eu . For simplicity here, we have normalized v so that $v^T u = 1$; a similar result holds if we set $v^T v = 1$ instead, if for instance $\|E\|_2 \rightarrow 0$ asymptotically.

2.5. Second order accuracy result for eigenvectors by averaging

In light of Equation (11), it is clear that v is a first order accurate approximation of u , because of the presence of the (first-order) term REu in the expansion. We now show that we can get a second order accurate approximation of the eigenvector u . Our results are based on an averaging procedure and hence are easy to implement in a distributed fashion. We have the following second-order accuracy result.

Theorem 3. *Let us call u_1 the eigenvector associated with the largest eigenvalue of M , and $v_1 = v_1 / \|v_1\|$ the eigenvector associated with the largest eigenvalue of $S = M + E$ and normalized so that $\|v_1\| = 1$ and $v_1^T u_1 \geq 0$. Let us call $\xi = \mu / (pn^{\alpha_{\min}})^{1/2}$. Suppose that the assumptions of Theorem 1 are satisfied (hence $\xi \rightarrow 0$). Suppose also that $d = (\lambda_1 - \lambda_2)$ satisfies*

$$\frac{d}{\xi} \rightarrow \infty. \tag{12}$$

Then we have

$$\mathbb{E}[\|v_1 - u_1\|_2] = O\left(\frac{1}{(\lambda_1 - \lambda_2)^2} \frac{\mu^2}{pn^{\alpha_{\min}}}\right) = O\left(\frac{\xi^2}{d^2}\right).$$

Practically, this means that if we average eigenvectors over many subsampled matrices (after removing indeterminacy by always making the first component positive), the residual error will be of order $\|E\|_2^2/d^2$ with

$$\limsup_{n \rightarrow \infty} \|E\|_2^2 \leq (\mathfrak{K}(\gamma))^2 \frac{\mu^2}{pn^{\alpha_{\min}}} \text{ a.s. ,}$$

where we recall that $\mathfrak{K}(\gamma)$ can be replaced by 2 in most cases (i.e as soon as $(\log n^{\alpha_{\min}})^4/(pn^{\alpha_{\min}}) \rightarrow 0$) and minimal conditions on α_{\min} mentioned above are satisfied. In other words, by averaging subsampled eigenvectors, we gain an order of accuracy (over the method that would just take one subsampled eigenvector) by canceling the effect of the first order residual term REu .

Technically, the previous theorem relies partly (and somewhat indirectly) on repeated use of a concentration inequality due to Talagrand, whose consequences are well explained in [Led01], Chapter 4. This inequality allows us to control the probability that $\|E\|_2$ exceeds a certain threshold, and hence tells us that with overwhelming probability, Theorem 2 is applicable in the setting we consider. However, this is not enough to be able to bound the expectations we care about. So we also make use of it to guarantee that the higher moments of $\|E\|_2$ can be controlled, which gives us explicit control on $\mathbf{E}[\|u - \tilde{v}_\varepsilon\|_2]$, where \tilde{v}_ε is a regularization of the perturbed eigenvector v which we need to use for technical reasons (we cannot directly take expectations in the bounds of Theorem 2).

Proof. To keep notations simple, we drop the index 1 in ν and u in the proof (so $\nu_1 = \nu$ and $u_1 = u$). In what follows, κ is a generic constant that may change from display to display. Before we start the proof per se, let us make a few remarks. (We recall that $\Delta = R(E - (\lambda_1(S) - \lambda_1(M))\text{Id})$.)

First, there is a technical difficulty when trying to work directly with v , namely the fact that it appears difficult to control $\mathbf{E}[\|(\text{Id} + \Delta)^{-1}\|_2]$ and hence to get a bound on $\mathbf{E}[\|v - u\|]$ (with the normalization $v^T u = 1$, $\|v\|$ could be very large; our bounds show that this can happen with only low probability but obviously $\mathbf{E}[\|v\|]$ could still be large). To go around this difficulty, we need two steps: first, we work with unit eigenvectors (so we go from v to ν), and second we need a “regularization” step and will replace v by a vector \tilde{v}_ε which is equal to v with high-probability and for which we can control $\mathbf{E}[\|\tilde{v}_\varepsilon - u\|]$. More precisely, for $\varepsilon > 0$, we call \tilde{v}_ε the vector such that

$$\tilde{v}_\varepsilon = \begin{cases} v & \text{if } \|(\text{Id} + \Delta)^{-1}\|_2 \leq \frac{1}{\varepsilon} \\ u - REu + \Delta REu & \text{otherwise.} \end{cases}$$

Its properties are studied in Theorem B.3. We call it below the ε -regularized version of v .

We note that under the assumptions of the current theorem we have $\frac{\xi}{d} \rightarrow 0$, so the results of Theorem B.3 apply. In particular, as shown in the proof of that Theorem, we have $\|M\|_\infty^2/p^2 = o(\xi^2)$. Also, Assumption 1 (which is made in Theorem 1), means μ is fixed so $\xi \rightarrow 0$, as $pn^{\alpha_{\min}} \rightarrow \infty$.

If v is the eigenvector of S associated with its largest eigenvalue, using the fact that $(v - u)^T u = 0$ by construction, we have

$$\|v\|_2^2 = \|v - u\|_2^2 + \|u\|_2^2 = 1 + \|v - u\|_2^2$$

hence

$$\nu = \frac{v}{\sqrt{1 + \|v - u\|_2^2}}.$$

Turning our attention to \tilde{v}_ε , we see that, since $Ru = 0$ by construction and R is symmetric, $u^T \Delta = 0$, so $(\tilde{v}_\varepsilon - u)^T u = 0$, and hence

$$\|\tilde{v}_\varepsilon\|_2^2 = 1 + \|\tilde{v}_\varepsilon - u\|_2^2.$$

Now let us call

$$\beta = \frac{\tilde{v}_\varepsilon}{\sqrt{1 + \|\tilde{v}_\varepsilon - u\|_2^2}},$$

we see that $\beta = \nu$ as long as $\|(\text{Id} + \Delta)^{-1}\|_2 \leq 1/\varepsilon$, since when this happens, $v = \tilde{v}_\varepsilon$. Now we have

$$\begin{aligned} \mathbf{E}[\|u - \nu\|_2] &= \mathbf{E}[\|u - \nu\|_2 \mathbf{1}_{\nu=\beta}] + \mathbf{E}[\|u - \nu\|_2 \mathbf{1}_{\nu \neq \beta}] \\ &\leq \mathbf{E}[\|u - \beta\|_2 \mathbf{1}_{\nu=\beta}] + \mathbf{E}[\|u - \nu\|_2 \mathbf{1}_{\nu \neq \beta}] \\ &\leq \mathbf{E}[\|u - \beta\|_2] + 2P(\nu \neq \beta), \end{aligned}$$

since $\|u - \nu\|_2 \leq \|u\|_2 + \|\nu\|_2 = 2$ (note the importance of the change of normalization here, as this bound would not hold with v instead of ν). Let us now work on controlling both these quantities. For reasons that will be clear later, we now take $\varepsilon = \mathfrak{K}(\gamma)\xi/d$.

Control of $\mathbf{E}[\|u - \beta\|_2]$. Given that $u - \beta = (u - \tilde{v}_\varepsilon)/\sqrt{1 + \|u - \tilde{v}_\varepsilon\|_2^2} + u(1 - 1/\sqrt{1 + \|u - \tilde{v}_\varepsilon\|_2^2})$, we have

$$\begin{aligned} \|u - \beta\|_2 &\leq \frac{\|u - \tilde{v}_\varepsilon\|_2}{\sqrt{1 + \|u - \tilde{v}_\varepsilon\|_2^2}} + \|u\|_2 \left(1 - \frac{1}{\sqrt{1 + \|u - \tilde{v}_\varepsilon\|_2^2}}\right) \\ &\leq \|u - \tilde{v}_\varepsilon\|_2 + (\sqrt{1 + \|u - \tilde{v}_\varepsilon\|_2^2} - 1) \\ &\leq 2\|u - \tilde{v}_\varepsilon\|_2, \end{aligned}$$

since $\sqrt{1 + x^2} \leq 1 + x$ for $x \geq 0$. Let us call $\mu/(pn^{\alpha_{min}})^{1/2} = \xi$ and $d = \lambda_1 - \lambda_2$. We show in Theorem B.3 that, for some $\kappa > 0$, asymptotically

$$\mathbf{E}[\|u - \tilde{v}_\varepsilon\|_2] \leq \kappa \left(\frac{\xi^2}{d^2} + \frac{\xi^3}{d^3 \varepsilon}\right)$$

so when $\varepsilon > \xi/d$, we have $\mathbf{E}[\|u - \tilde{v}_\varepsilon\|_2] \leq \kappa \frac{\xi^2}{d^2}$ and therefore

$$\mathbf{E}[\|u - \beta\|_2] \leq \kappa \frac{\xi^2}{d^2}.$$

Control of $P(\nu \neq \beta)$. We have (essentially) seen in the proof of Theorem 2 above that if $2\|E\|_2/d < 1 - \varepsilon$, then $\|(\text{Id} + \Delta)^{-1}\|_2 \leq 1/\varepsilon$ (see also the proof of Theorem B.3). Hence

$$P\left(\|(\text{Id} + \Delta)^{-1}\|_2 > 1/\varepsilon\right) \leq P\left(\|E\|_2 > \frac{(1 - \varepsilon)d}{2}\right).$$

Recall that we have now chosen $\varepsilon = \mathfrak{K}(\gamma)\xi/d$. In that case, we have

$$\frac{(1 - \varepsilon)d}{2} = \frac{d}{2} - \frac{\mathfrak{K}(\gamma)}{2}\xi.$$

Now we show the following deviation inequality in Theorem B.2: if m_E is a median of $\|E\|_2$,

$$P\left(\left|\|E\|_2 - m_E\right| > t\right) \leq 4 \exp\left(-\frac{p^2}{8\|M\|_\infty^2}t^2\right).$$

Recall also that for n large enough $0 \leq m_E \leq (\mathfrak{K}(\gamma) + 1)\xi$ when the conditions of Theorem 1 apply (see Theorems 1 or arguments at the end of the proof of Theorem B.1). Suppose now that n is such that indeed $m_E \leq (\mathfrak{K}(\gamma) + 1)\xi$. Then if $\frac{d}{2} - (\frac{3}{2}\mathfrak{K}(\gamma) + 1)\xi > 0$, we have

$$\begin{aligned} P\left(\|E\|_2 > \frac{(1 - \varepsilon)d}{2}\right) &\leq P\left(\left|\|E\|_2 - m_E\right| > \frac{(1 - \varepsilon)d}{2} - m_E\right) \\ &\leq P\left(\left|\|E\|_2 - m_E\right| > \frac{d}{2} - \left(\frac{3}{2}\mathfrak{K}(\gamma) + 1\right)\xi\right). \end{aligned}$$

Now when $\xi/d \rightarrow 0$, and because $\mathfrak{K}(\gamma)$ stays bounded, $\frac{d}{2} - (\frac{3}{2}\mathfrak{K}(\gamma) + 1)\xi \geq \frac{d}{3}$ asymptotically. Note that by assumption, $\xi/d \rightarrow 0$. Therefore,

$$P\left(\|E\|_2 > \frac{(1 - \varepsilon)d}{2}\right) \leq 4 \exp\left(-\frac{p^2}{72\|M\|_\infty^2}d^2\right).$$

All we have to do now is to verify that the asymptotics we consider, the quantity on the right-hand side of the previous equation remains less than ξ^2/d^2 asymptotically. Elementary algebra shows that this is equivalent to saying that

$$1 \geq \frac{1}{d^2}72\frac{\|M\|_\infty^2}{p^2}(\ln(d^2/\xi^2) + \ln 4). \tag{13}$$

We have $\|M\|_\infty^2/p^2 = o(\xi^2)$, so the right-hand side is going to zero because it is $o((\xi/d)^2 \log(d/\xi))$ and $\xi/d \rightarrow 0$. So we have shown that under our assumptions,

$$P(\nu \neq \beta) \leq \frac{\xi^2}{d^2}.$$

We can finally conclude that

$$\mathbf{E}[\|\nu - u\|_2] \leq \kappa \frac{\xi^2}{d^2},$$

as announced in the theorem. □

This result applies to all eigenvectors corresponding to eigenvalues whose isolation distance (i.e distance to the nearest eigenvalue) satisfies the separation condition (12), which is a strong version of the separation condition (10). We note that we need the strong separation condition (Equation (12)) to be able to take expectations rigorously.

Finally, we note that theoretical as well as practical considerations seem to indicate that condition (10) (and hence (12)) is quite conservative. On the theoretical side, we see with Equation (9) that what really matters for the quality of the approximation is the norm of the vector

$$l_j = \Delta^{j+2}(\text{Id} + \Delta)^{-1}u ,$$

or its expectation. We used in our approximations the coarse bound $\|\Delta\|_2 \leq 2\|R\|_2\|E\|_2$, which is convenient because it does not require us to have information about the eigenvectors of Δ . However, we see that the norm of l_j could be small even when $\|R\|_2\|E\|_2$ is not very small, for instance if u belonged to a subspace spanned by eigenvectors of Δ associated with eigenvalues of this matrix that are small in absolute value. So it is quite possible that our method could work in a somewhat larger range of situations than the one for which we have theoretical guarantees. This is what our simulations below seem to indicate.

2.6. Variance

The expansion in Equation (11) also allows us to approximate the variance of the first-order residual REu after subsampling. This is useful in practice because it gives us an idea of how many independent computations we need to make to essentially void the effect of the first order term in the expansion of v . In terms of distributed computing, it therefore tells us how many machines we should involve in the computation. We have the following theorem.

Theorem 4. *Let u_1 be the eigenvector associated with λ_1 , the largest eigenvalue of M . Let us call $w_1 = u_1 \circ u_1$, and $\mathcal{M} = M \circ M$. Then*

$$\begin{aligned} & \mathbf{E}[\|REu_1\|_2^2] \\ & \leq \frac{1}{(\lambda_2 - \lambda_1)^2} \frac{1-p}{p} \left(\sum_{k=1}^n u_1(k)^2 \|M_k\|_2^2 - \left[2w_1^T \mathcal{M} w_1 - \sum_{k=1}^n w_1^2(k) \mathcal{M}_{kk} \right] \right) . \end{aligned}$$

Assuming w.l.o.g. that $\lambda_1 = \|M\|_2$, this bound yields in particular

$$\mathbf{E}[\|REu_1\|_2^2] \leq \frac{1}{(1 - \lambda_2/\lambda_1)^2} \|u_1\|_\infty^2 \frac{\mathbf{NumRank}(M)}{p} \tag{14}$$

where $\mathbf{NumRank}(M) = \|M\|_F^2 / \|M\|_2^2$ is the numerical rank of the matrix M and is a stable relaxation of the rank, satisfying $1 \leq \mathbf{NumRank}(M) \leq \mathbf{Rank}(M) \leq n$ (see [RV07] for a discussion).

Proof. By construction, $\mathbf{E}[E] = 0$ and

$$\mathbf{E}[\|REu_1\|_2^2] = \mathbf{E}[u_1^T ER^2Eu_1] = \sum_{j=2}^n \mathbf{E}\left[\frac{(u_1^T Eu_j)^2}{(\lambda_j - \lambda_1)^2}\right],$$

by definition of R . Now

$$\sum_{j=1}^n (u_1^T Eu_j)^2 = \|Eu_1\|_2^2 = u_1^T E^2 u_1,$$

because E is symmetric, the u_i 's form an orthonormal basis and $u_1^T Eu_j$ is the j -th coefficient of Eu_1 in this basis, so the sum of the squared coefficients is the squared norm of the vector. Hence

$$\mathbf{E}[\|REu_1\|_2^2] \leq \frac{1}{(\lambda_2 - \lambda_1)^2} (\mathbf{E}[u_1^T E^2 u_1] - \mathbf{var}(u_1^T Eu_1)).$$

The variance of $u_1^T Eu_1$ is easy to compute if we rewrite this quantity as a sum of independent random variables. Also, separate computations (see Appendix, Subsection B.3) show that $\mathbf{E}[E^2]$ is a diagonal matrix, whose i -th diagonal entry is $(1-p)\|M_i\|_2^2/p$, where M_i is the i -th column of M . Hence, in that case, having defined $w_1 = u_1 \circ u_1$ and $\mathcal{M} = M \circ M$, we get

$$\begin{aligned} & \mathbf{E}[\|REu_1\|_2^2] \\ & \leq \frac{1}{(\lambda_2 - \lambda_1)^2} \frac{1-p}{p} \left(\sum_{k=1}^n u_1(k)^2 \|M_k\|_2^2 - \left[2w_1^T \mathcal{M} w_1 - \sum_{k=1}^n w_1^2(k) \mathcal{M}_{kk} \right] \right). \end{aligned}$$

Assuming w.l.o.g. that $\lambda_1 = \|M\|_2$, we get (14). □

2.7. Nonsymmetric matrices

The results described above are easily extended to nonsymmetric matrices. Here $M \in \mathbf{R}^{m \times n}$, with $m \geq n$ and we write its singular value decomposition

$$M = \sum_{i=1}^n \sigma_i u_i v_i^T,$$

where $u_i \in \mathbf{R}^n$, $v_i \in \mathbf{R}^m$ and $\sigma_i > 0$. We can adapt the definition of incoherence to

$$\mu(M, \alpha, \beta) = \sum_{i=1}^n \sigma_i n^{\alpha_i/2} \|u_i\|_\infty m^{\beta_i/2} \|v_i\|_\infty$$

and reformulate our main assumption on M as follows.

Assumption 2. *There are vectors $\alpha \in [0, 1]^n$ and $\beta \in [0, 1]^n$ for which*

$$\mu(M, \alpha, \beta) \leq \mu \quad \text{and} \quad \mathbf{Card}(u_i) \leq n^{\alpha_i}, \quad \mathbf{Card}(v_i) \leq m^{\beta_i}, \quad i = 1, \dots, n$$

as m, n go to infinity with $m = \rho n$ for a given $\rho > 1$, where μ is an absolute constant.

In this setting, using again [HJ91, Th. 5.5.19], we get

$$\left\| \sum_{i=1}^n \sigma_i C \circ (u_i v_i^T) \right\|_2 \leq \sum_{i=1}^n \sigma_i n^{\alpha_i/2} \|u_i\|_\infty m^{\beta_i/2} \|v_i\|_\infty \left\| \frac{C_{\alpha_i, \beta_i}}{n^{\alpha_i/2} m^{\beta_i/2}} \right\|_2 \quad (15)$$

where we have assumed that u_i, v_i are sparse and C_{α_i, β_i} is a $n^{\alpha_i} \times m^{\beta_i}$ submatrix of C . As in (5), we can then bound the spectral norm of the residual and we have

$$\limsup_{n \rightarrow \infty} \|E\|_2 \leq \frac{\mathfrak{K}(\gamma)\mu}{\sqrt{p(n^{\alpha_{\min}} \wedge m^{\beta_{\min}})}}. \quad (16)$$

in probability, as soon as the sampling probability p is such that if $s_i = n^{\alpha_i} \vee m^{\beta_i}$, $p \geq \gamma(\log s_i)/s_i$, for all i . Perturbation results similar to (11) for left and right eigenvectors are detailed in [Ste98] for example.

We also note that our arguments go through if M is for instance a diagonalizable square matrix, after we replace all the potentially complex numbers appearing in the definition of incoherence by their modulus: if $M = \sum_i \lambda_i u_i v_i^T$, where $\lambda_i \in \mathbf{C}$ and $u_i, v_i \in \mathbf{C}^n$, $\mu(M, \alpha, \beta) = \sum_{i=1}^n |\lambda_i| n^{\alpha_i/2} \|u_i\|_\infty n^{\beta_i/2} \|v_i\|_\infty$, where for a vector $\nu \in \mathbf{C}^n$, $\|\nu\|_\infty = \max_k |\nu(k)|$.

3. Numerical experiments

In this section, we study the numerical performance of the subsampling/averaging results detailed above on both artificial and realistic data matrices

Dense matrices: PCA, SVD, etc. We first illustrate our results by approximating the leading eigenvector of a matrix M as the average of leading eigenvectors of subsampled matrices, for various values of the sampling probability p . To start with a naturally structured dense matrix, we form M as the covariance matrix of the 500 most active genes in the colon cancer data set in [ABN⁺99]. We let p vary from 10^{-4} to 1 and for each p , we compute the leading eigenvector of 1000 subsampled matrices, average these vectors and normalize the result. We call u the true leading eigenvector of M and v the approximate one. We now normalize v so that $\|v\|_2 = 1$ (which is standard, but different from the normalization we used in our theoretical investigations where we had $u^T v = 1$).

In Figure 2, we plot $u^T v$ as a function of p together with the median of $u^T v$ sampled over all individual subsampled matrices, with dotted lines at plus and minus one standard deviation. We also record the proportion of samples where $\|E\|$ satisfies the perturbation condition (10).

We repeat this experiment on a (nonsymmetric) term-document matrix formed using press release data from PRnewswire, to test the impact of subsampling on Latent Semantic Indexing results. Once again, we let p vary from 10^{-2} to 1 and for each p , we compute the leading eigenvector of 1000 subsampled matrices, average these vectors and normalize the result. We call u the true leading eigenvector of M and v the approximate one. In Figure 3 on the left, we plot

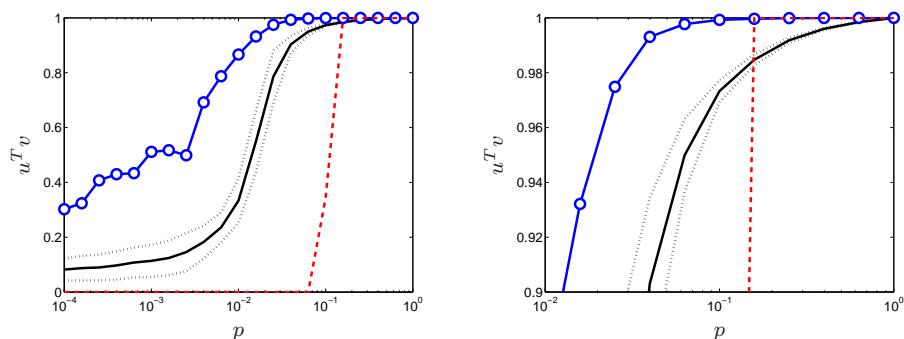


FIG 2. Left: Alignment $u^T v$ between the true and the normalized average of 1000 subsampled eigenvectors (blue circles), median value of $u^T v$ over all sampled matrices (solid black line), with dotted lines at plus and minus one standard deviation and proportion of samples satisfying condition (10) (dashed red line), for various values of the sampling probability p on a gene expression covariance matrix. Right: Zoom on the the interval $p \in [10^{-2}, 1]$.

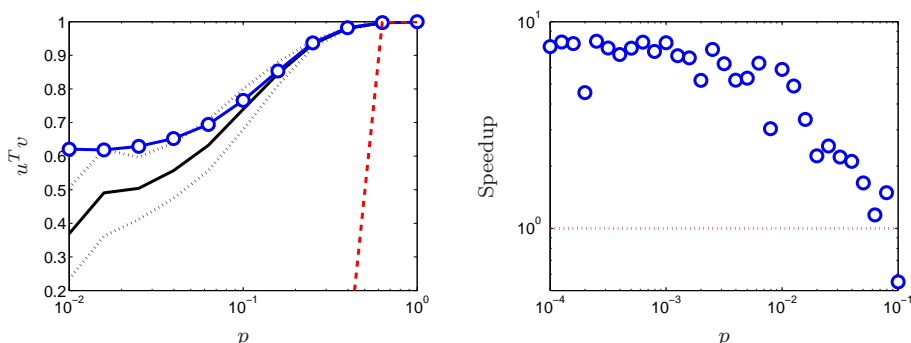


FIG 3. Left: Alignment $u^T v$ between the true and the normalized average of 1000 subsampled left eigenvectors (blue circles), median value (solid black line), dotted lines at plus and minus one standard deviation and proportion of samples satisfying condition (10) (dashed red line), for various values of the sampling probability p on a term document matrix with dimensions 6779×11171 . Right: Speedup in computing leading eigenvectors on gene expression data, for various values of the sampling probability p .

$u^T v$ as a function of p together with the median of $u^T v$ sampled over all individual subsampled matrices, with dotted lines at plus and minus one standard deviation. The matrix M is 6779×11171 with spectral gap $\sigma_2/\sigma_1 = 0.66$.

In Figure 3 on the right, we plot the ratio of CPU time for subsampling a gene expression matrix of dimension 2000 and computing the leading eigenvector of the subsampled matrix (on a single machine), over CPU time for computing the leading eigenvector of the original matrix. Two regimes appear, one where the eigenvalue computation dominates with computation cost scaling with p , another where the sampling cost dominates and the speedup is simply the ratio

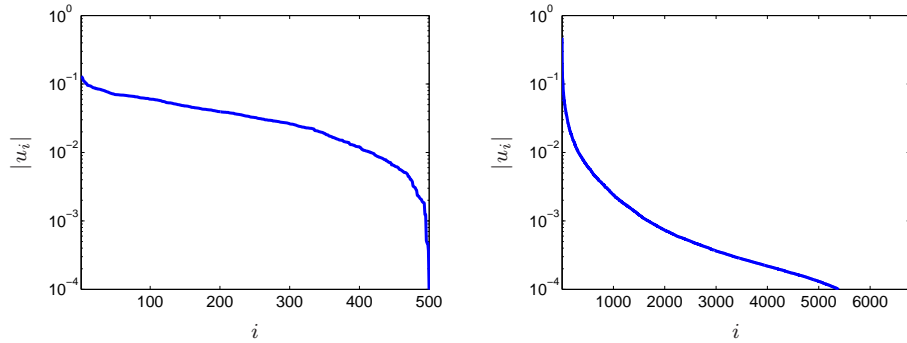


FIG 4. Magnitude of eigenvector coefficients $|u_i|$ in decreasing order for both the leading eigenvector of the gene expression covariance matrix (left) and leading left eigenvector of the 6779×11171 term document matrix (right).

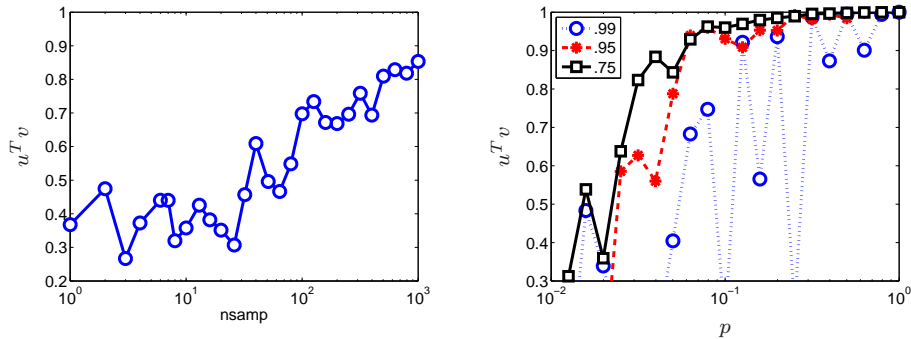


FIG 5. Left: Alignment $u^T v$ between the true leading eigenvector u and the normalized average leading eigenvector versus number of samples, on the gene expression covariance matrix with subsampling probability $p = 10^{-2}$. Right: Alignment $u^T v$ for various values of the spectral gap $\lambda_2/\lambda_1 \in \{0.75, 0.95, 0.99\}$.

between sampling time and the CPU cost of a full eigenvector computation. Of course, the principal computational benefit of subsampling is the fact that memory usage is directly proportional to p .

A key difference between the experiments of Figure 2 and those of 3 is that the leading eigenvector of the gene expression data set is much more incoherent than the leading left eigenvector of the term-document matrix, which explains part of the difference in performance. We compare both eigenvectors in Figure 4.

We then study the impact of the number of samples on precision. We use again the colon cancer data set in [ABN+99]. In Figure 5 on the left, we fix the sampling rate at $p = 10^{-2}$ and plot $u^T v$ as a function of the number of samples used in averaging. We also measure the impact of the eigenvalue gap λ_2/λ_1 on precision. We scale the spectrum of the gene expression covariance matrix so that its first eigenvalue is $\lambda_1 = 1$ and plot the alignment $u^T v$ be-

tween the true and the normalized average of 100 subsampled eigenvectors over subsampling probabilities $p \in [10^{-2}, 1]$ for various values of the spectral gap $\lambda_2/\lambda_1 \in \{0.75, 0.95, 0.99\}$.

Graph matrices: ranking. Here, we test the performance of the methods described above on graph matrices used in ranking algorithms such as pagerank [PBMW98] (because of its susceptibility to manipulations however, this is only one of many features used by search engines). Suppose we are given the adjacency matrix of a web graph, with

$$\begin{cases} A_{ij} = 1, & \text{if there is a link from } i \text{ to } j \\ A_{ij} = 0, & \text{otherwise,} \end{cases}$$

where $A \in \mathbf{R}^{n \times n}$ (one such matrix is displayed in Figure 6). Whenever a node has no out-links, we link it with every other node in the graph, so that $B = A + \delta \mathbf{1}^T/n$, with $\delta_i = 1$ if and only if $\text{deg}_i = 0$, where deg_i is the degree of node i . We then normalize B into a stochastic matrix $P_{ij}^g = B_{ij}/\text{deg}_i$. The matrix P^g is the transition matrix of a Markov chain on the graph modeling the behavior of a web surfer randomly clicking on links at every page. For most web graphs, this Markov chain is usually not irreducible but if we set

$$P = cP^g + (1 - c)\mathbf{1}\mathbf{1}^T/n$$

for some $c \in (0, 1]$, the Markov chain with transition matrix P will be irreducible. An additional benefit of this modification is that the spectral gap of P is at least c [HK03]. The leading (Perron-Frobenius) eigenvector u of this matrix is called the *Pagerank* vector [PBMW98], its coefficients u_i measure the stationary probability of page i being visited by a random surfer driven by the transition matrix P , hence reflect the importance of page i according to this model.

The coefficients of pagerank vectors typically follow a power law for classic values of the damping factor [PRU06, BC06] which means that the bounds in

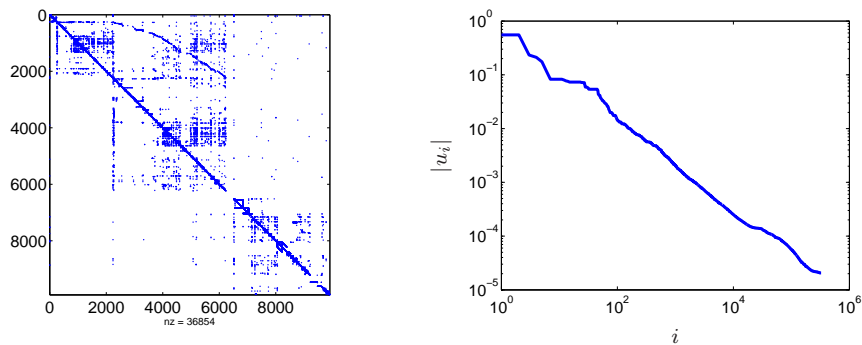


FIG 6. Left: The `wb-cs.stanford` graph. Right: Loglog plot of the Pagerank vector coefficients for the `cnr-2000` graph.

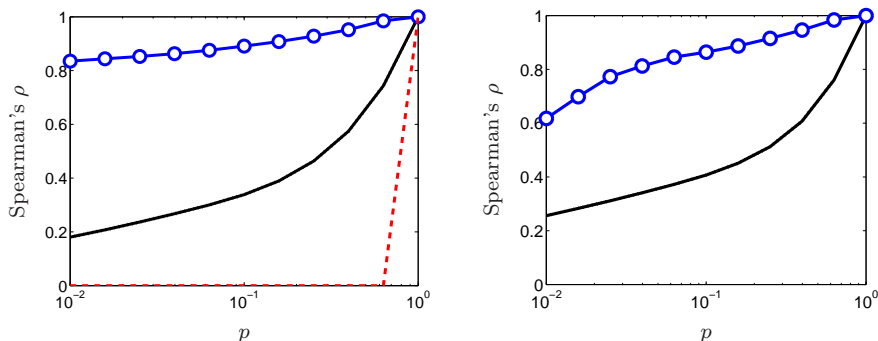


FIG 7. Ranking correlation (Spearman’s ρ) between true and averaged pagerank vector (blue circles), median value of the correlation over all subsampled matrices (solid black line), dotted lines at plus and minus one standard deviation and proportion of samples satisfying the perturbation condition (10) (dashed red line), for various values of the sampling probability p . Left: On the CNR-2000 graph. Right: On the UK-2002 graph.

assumption 1 do not hold. Empirically however, while the distance between true and averaged eigenvectors quickly gets large, the ranking correlation (measured using Spearman’s ρ [Mel07]) is surprisingly robust to subsampling.

We use two graphs from the Webgraph database [BV04], `cnr-2000` which has 3×10^5 nodes and 3×10^6 edges, and `UK-2002` with 2×10^7 nodes and 3×10^8 edges. For each graph, we form the transition matrix P as in [GZB04] with uniform teleportation probability and set the teleportation coefficient $c = 0.85$. In Figure 6 we plot the `wb-cs.stanford` graph and the Pagerank vector for `cnr-2000` in loglog scale. In Figure 7 we plot the ranking correlation (Spearman’s ρ) between true and averaged Pagerank vector (over 1000 samples), the median value of the correlation over all subsampled matrices and the proportion of samples satisfying the perturbation condition (10), for various values of the sampling probability p . We notice that averaging very significantly improves ranking correlation, far outside the perturbation regime.

4. Conclusion

We have proposed a method to compute the eigenvectors of very large matrices in a distributed fashion:

1. To each node in a computer cluster of size N , we send a subsampled version S_i of the matrix of interest, M .
2. Node i computes the relevant eigenvectors of S_i .
3. The N eigenvectors are averaged together and normalized to produce our final estimator.

The key to the algorithm is that Step 2 is numerically cheap (because S_i is very sparse), and hence can be executed fast even on small machines. Therefore a

cluster or cloud of small machines could be used to approximate the eigenvectors of M , a difficult problem in general when M is extremely large.

We have shown that under carefully stated conditions, the algorithm described above will yield a second-order accurate approximation of the eigenvectors of M . This gain in accuracy comes from the averaging step of our algorithm. We note that arguments similar to the ones we used in this paper could be made to compute second-order accurate approximations of the eigenvalues of M . (We restricted ourselves to eigenvectors here because in methods such as PCA, the eigenvectors are in some sense more important than the eigenvalues.) Our results depend on a measure of incoherence for M that we propose in this paper. This notion of incoherence makes the algorithm very suitable for matrices whose eigenvectors are not very sparse (though somewhat sparse eigenvectors can also be handled). We hence expect that the method we propose could be useful in working with, for instance, kernel matrices. Our work also shows that subsampling will work if the sampling probability is small, but is likely to fail if that probability is too small.

At a more general level, we note that a similar averaging step could be used for other randomized algorithms in numerical linear algebra, provided that these methods yield unbiased estimates of the entries of M . Then, similar techniques to the ones we employ could be used to investigate second-order accuracy of the corresponding algorithms (after the averaging step), provided the original algorithm also results, with high-probability, in a bounded and reasonably small perturbation of the matrix M . (Much of our random matrix analysis is devoted to showing just this for the particular subsampling algorithm we are concerned with.)

Finally, our simulations show that we gain significantly in accuracy by averaging subsampled eigenvectors (which suggests that our theoretical passage from first-order to second-order accuracy is also relevant in practice) and that the performance of our method seems to degrade for very incoherent matrices, a result that is also in line with our theoretical predictions.

Appendix A: On the median of $\|C\|_2$

A key quantity in the problems we are considering is the random variable

$$\left\| \frac{C}{\sqrt{n}} \right\|_2,$$

where C is a $n \times n$ random matrix (possibly symmetric) whose entries are i.i.d with

$$C_{i,j} = \begin{cases} \sqrt{\frac{1-p}{p}} & \text{with probability } p \\ -\sqrt{\frac{p}{1-p}} & \text{with probability } 1-p \end{cases}.$$

While if $p \geq \gamma(\log n)^4/n$, for γ bounded away from 0, a result of Vu [Vu07] allows us to bound $\mathbf{E}[\| \frac{C}{\sqrt{n}} \|_2]$, we are not aware of existing results in the case $p \ll (\log n)^4/n$.

We show in this Appendix that there is a phase transition for p of order $(\log n)/n$: if $p = \gamma(\log n)/n$, where γ is bounded below, $\mathbf{E}[\|\frac{C}{\sqrt{n}}\|_2]$ and $\|\frac{C}{\sqrt{n}}\|_2$ stay bounded (with high-probability in the latter case). When $p = (\log n)^{1-\delta}/n$, for some $\delta > 0$, then $\|\frac{C}{\sqrt{n}}\|_2 \rightarrow \infty$ in probability.

So in this Appendix, unless otherwise noted, we suppose that $p = \gamma(\log n)^{1+\varepsilon}/n$, where $\varepsilon \geq 0$ and $p \leq 1/2$. We assume throughout this appendix that $\gamma(1-p)$ is bounded below. This is without loss of generality, for otherwise we would change the rate at which p is going to 0 with n .

A.1. The case of non-symmetric C

We have the following theorem.

Theorem A.1. *Suppose that the entries of C are i.i.d with the distribution mentioned above, so C is not symmetric. C is $n \times n$ and $p = \gamma(\log n)^{1+\varepsilon}/n$, $\gamma > 0$. Then there exists a K (independent of n or p and finite) such that*

$$\left\| \frac{C}{\sqrt{n}} \right\|_2 \leq K \left(1 + \frac{2(\log n)^{-\varepsilon/2}}{\sqrt{\gamma(1-p)}} \right)$$

with very high-probability as n gets large.

We also have the same result when C is symmetric but this requires a separate argument, given below. Note that using Vu's result, we get $K = 2$ when $\varepsilon \geq 3+\delta$.

We also note that the previous result implies (through elementary linear algebraic considerations) that the same bound holds if C is $n \times m$ with $m \leq n$.

Proof. The proof is in several steps. We use a result of Seginer [Seg00] in connection with Talagrand's inequality and some careful manipulations. Note that the entries of C are supported on an interval $[u, v]$ with $v-u = 1/\sqrt{p(1-p)}$. So Talagrand's inequality gives us that if F is a convex 1-Lipschitz function (with respect to Euclidian norm) of the $C_{i,j}$, we have

$$P(|F - m_F| \geq t) \leq 4 \exp(-p(1-p)t^2/4),$$

according to [Led01], Corollary 4.10, where m_F is a median of the random variable $F(\{C_{i,j}\})$.

Standard results also give us control of the deviation from the mean, through, according to Proposition 1.9 in [Led01],

$$|\mu_F - m_F| \leq 4 \sqrt{\frac{\pi}{p(1-p)}}. \tag{A.1}$$

Now, it is well known that $\|C\|_2$ is a convex 1-Lipschitz function of the entries of C , so

$$\boxed{P(|\|C\|_2 - m_C| \geq t) \leq 4 \exp(-p(1-p)t^2/4)}. \tag{A.2}$$

Therefore, when $p = \gamma(\log n)^{1+\varepsilon}/n$, we have

$$P\left(\frac{|\|C\|_2 - m_C|}{\sqrt{n}} \geq t\right) \leq 4 \exp(-\gamma(\log n)^{1+\varepsilon}(1-p)t^2/4).$$

So if we can establish that $m_C/\sqrt{n} \leq K_1$, at least asymptotically, then we will have, for any given $l > 0$,

$$\|C\|_2/\sqrt{n} \leq K_1 + \frac{(\log n)^{-\frac{1+\varepsilon}{2(l+1)}}}{\sqrt{\gamma(1-p)}},$$

with very high-probability (by just picking $t = [\gamma(1-p)(\log n)^{\frac{1+\varepsilon}{1+t}}]^{-1/2}$).

Furthermore, the estimate above (Equation (A.1)) gives us, in connection with Proposition 1.9 in [Led01], that

$$\boxed{\left| \frac{\mathbf{E}[\|C\|_2]}{\sqrt{n}} - \frac{m_C}{\sqrt{n}} \right| \leq 4 \frac{\sqrt{\pi}}{\sqrt{\gamma(1-p)(\log n)^{1+\varepsilon}}}.}$$

So control of $\mathbf{E}[\|C\|_2]$ is all we need.

Now Seginer's result states that if A is a $n \times m$ matrix with i.i.d entries and mean 0, we have, for K a universal constant (i.e independent of A, n, m),

$$\mathbf{E}[\|A\|_2] \leq K \left(\mathbf{E} \left[\max_{1 \leq i \leq n} \|A(i, \cdot)\|_2 \right] + \mathbf{E} \left[\max_{1 \leq j \leq m} \|A(\cdot, j)\|_2 \right] \right).$$

Here $\|A(i, \cdot)\|_2$ is the Euclidian norm of the i -th row of A , and $\|A(\cdot, j)\|_2$ is the Euclidian norm of the j -th column of A .

Let us focus on $\mathbf{E}[\|A(i, \cdot)\|_2]$. Let us call m_i the median of $\|A(i, \cdot)\|_2$. Clearly,

$$\max_{1 \leq i \leq n} \|A(i, \cdot)\|_2 \leq \max_{1 \leq i \leq n} |\|A(i, \cdot)\|_2 - m_i| + \max_{1 \leq i \leq n} m_i.$$

Since the $A(i, \cdot)$ are identically distributed, $m_i = m_j$ for all i, j .

In particular,

$$\mathbf{E} \left[\max_{1 \leq i \leq n} \|A(i, \cdot)\|_2 \right] \leq \mathbf{E} \left[\max_{1 \leq i \leq n} |\|A(i, \cdot)\|_2 - m_i| \right] + m_1.$$

Two tasks remain: bounding m_1 and controlling $\mathbf{E}[\max_{1 \leq i \leq n} |\|A(i, \cdot)\|_2 - m_i|]$, when A is replaced by C . We start by the latter.

Controlling $\mathbf{E} \left[\frac{\max_{1 \leq i \leq n} |\|C(i, \cdot)\|_2 - m_i|}{\sqrt{n}} \right]$

Note that $\|C(i, \cdot)\|_2 = \sqrt{\sum_{j=1}^n C_{i,j}^2}$. This is clearly a convex 1-Lipschitz function of the $C_{i,j}$. So Talagrand's inequality guarantees that

$$P(|\|C(i, \cdot)\|_2 - m_i| \geq t) \leq 4 \exp(-p(1-p)t^2/4).$$

Using a union bound, we have

$$P(\max_{1 \leq i \leq n} \|C(i, \cdot)\|_2 - m_i \geq t) \leq (4n \exp(-p(1-p)t^2/4)) \wedge 1 .$$

Also, replacing p by its value, we have

$$P\left(\frac{\max_{1 \leq i \leq n} \|C(i, \cdot)\|_2 - m_i}{\sqrt{n}} \geq t\right) \leq (4n \exp(-\gamma(\log n)^{1+\varepsilon}(1-p)t^2/4)) \wedge 1 .$$

We now recall that for a non-negative random variable Y , we have

$$\mathbf{E}[Y] = \int_0^\infty P(Y \geq t) dt .$$

Hence, for any $x > 0$,

$$\begin{aligned} \mathbf{E} \left[\frac{\max_{1 \leq i \leq n} \|C(i, \cdot)\|_2 - m_i}{\sqrt{n}} \right] &\leq \int_0^\infty P\left(\frac{\max_{1 \leq i \leq n} \|C(i, \cdot)\|_2 - m_i}{\sqrt{n}} \geq t\right) dt \\ &\leq \int_0^x 1 dt + \int_x^\infty 4n \exp(-\gamma(\log n)^{1+\varepsilon}(1-p)t^2/4) dt \\ &\leq x + \frac{8n}{\gamma(\log n)^{1+\varepsilon}(1-p)x} \exp(-\gamma(\log n)^{1+\varepsilon}(1-p)x^2/4) , \end{aligned}$$

since

$$\int_x^\infty \exp(-\sigma t^2) dt \leq \frac{1}{2\sigma x} \exp(-\sigma x^2) .$$

If we pick $x = 2(\log n)^{-\varepsilon/2}/\sqrt{\gamma(1-p)}$, we have

$$\gamma(\log n)^{1+\varepsilon}(1-p)x^2/4 = (\log n) ,$$

so

$$\mathbf{E} \left[\frac{\max_{1 \leq i \leq n} \|C(i, \cdot)\|_2 - m_i}{\sqrt{n}} \right] \leq \frac{4}{\sqrt{\gamma(1-p)}(\log n)^{1+\varepsilon/2}} + \frac{2(\log n)^{-\varepsilon/2}}{\sqrt{1-p}\sqrt{\gamma}} .$$

Controlling $\frac{m_i}{\sqrt{n}}$

Using the fact that $\|C(i, \cdot)\|_2$ is a convex 1-Lipschitz function of the $C_{i,j}$, we have

$$\left| \frac{\mathbf{E}[\|C(i, \cdot)\|_2]}{\sqrt{n}} - \frac{m_i}{\sqrt{n}} \right| \leq 4 \frac{\sqrt{\pi}}{\sqrt{\gamma(1-p)}(\log n)^{1+\varepsilon}} .$$

Now, by concavity of $x \rightarrow \sqrt{x}$,

$$\mathbf{E}[\|C(i, \cdot)\|_2] = \mathbf{E} \left[\sqrt{\sum_{j=1}^n C_{i,j}^2} \right] \leq \sqrt{\mathbf{E} \left[\sum_{j=1}^n C_{i,j}^2 \right]} = \sqrt{n} .$$

Hence, for all i ,

$$\frac{m_i}{\sqrt{n}} \leq 1 + 4 \frac{\sqrt{\pi}}{\sqrt{\gamma(1-p)(\log n)^{1+\varepsilon}}}.$$

Putting everything together

Using our bounds we conclude that, for K_2 a universal constant (e.g $K_2 = 2K$, where K is the constant appearing in Seginer's result), and for $\varepsilon \geq 0$,

$$\mathbf{E} \left[\frac{\|C\|_2}{\sqrt{n}} \right] \leq K_2 \left[1 + \frac{2}{\sqrt{\gamma(1-p)(\log n)^{\varepsilon/2}}} \left(2 \frac{\sqrt{\pi}}{(\log n)^{1/2}} + \frac{2}{(\log n)} + 1 \right) \right]. \tag{A.3}$$

□

A.2. The case of symmetric C

Our aim is now to show that

Theorem A.2. *Suppose that $p = \gamma(\log n)^{1+\varepsilon}/n$, where $\varepsilon \geq 0$. Suppose that C is $n \times n$ and symmetric and its entries above and on the diagonal are i.i.d with the distribution mentioned at the beginning. Then there exists a K (independent of n or p and finite) such that*

$$\left\| \frac{C}{\sqrt{n}} \right\|_2 \leq K \left(1 + \frac{2(\log n)^{-\varepsilon/2}}{\sqrt{\gamma(1-p)}} \right).$$

with very high-probability as n gets large.

The proof now relies on a series of lemmas that will allow us to use Theorem A.1

Lemma A.1. *Suppose A and B are two matrices whose entries are independent and for all (i, j) and k integer,*

$$0 \leq \mathbf{E} [A_{i,j}^k] \leq \mathbf{E} [B_{i,j}^k].$$

Then, for all k integer,

$$\mathbf{E} [\text{trace}((A^T A)^k)] \leq \mathbf{E} [\text{trace}((B^T B)^k)].$$

In particular, for all k integer, if A and B have n non-zero singular values,

$$\left(\mathbf{E} [\|A\|_2^{2k}] \right)^{1/2k} \leq n^{1/2k} \left(\mathbf{E} [\|B\|_2^{2k}] \right)^{1/2k}.$$

Proof. Let us consider trace $((A^T A)^k)$. If we expand it in terms of $A_{i,j}$, this is the sum of many products of $A_{i,j}$'s. By independence of the $A_{i,j}$'s, the expectation of each of these products is the product of the $r_{i,j}$ -th moments of $A_{i,j}$, where $r_{i,j}$ is the number of times $A_{i,j}$ appear in the product.

Now for each such term, the corresponding term involving B is greater because it is greater term by term, as our conditions on the moments clearly show (here it is important that all the moments of $A_{i,j}$ be non-negative and less than all the moments of $B_{i,j}$).

So we have

$$\mathbf{E} [\text{trace} ((A^T A)^k)] \leq \mathbf{E} [\text{trace} ((B^T B)^k)] .$$

Now, $\text{trace} ((A^T A)^k) = \sum_{j=1}^n (\sigma_j(A))^{2k}$, where σ_j are the decreasingly ordered singular values of A , from which we simply deduce the second result. As a matter of fact, $\|A\|_2 = \sigma_1(A)$, so $\|A\|_2^{2k} \leq \text{trace} ((A^T A)^k)$, and $\text{trace} ((B^T B)^k) \leq n\sigma_1(B)^{2k}$. \square

Lemma A.2. *Suppose C is symmetric, $n \times n$ with entries i.i.d with the distribution given at the beginning. Then, for some universal K_3 ,*

$$\frac{1}{\sqrt{n}} \mathbf{E} [\|C\|_2] \leq K_3 \left(1 + \frac{2(\log n)^{-\varepsilon/2}}{\sqrt{\gamma(1-p)}} \right) .$$

The same result applies if $\mathbf{E} [\|C\|_2]$ is replaced by a median of $\|C\|_2$.

Proof. The first thing to note is that for all k , $\mathbf{E} [C_{i,j}^k] \geq 0$, when $p \leq 1/2$. Note that the maximal entry of C on the diagonal is at most of order $\sqrt{n/(\log n)^{1+\varepsilon}} \ll \sqrt{n}$, so we can replace the entries of C on the diagonal by zeroes without affecting the final result. Call \tilde{C} the corresponding matrix.

Now, we can write $\tilde{C} = C_0 + C_0^T$, where C_0 has i.i.d entries above the diagonal and 0 on and underneath it. Note that the entries of C_0 are independent. Clearly,

$$\|\tilde{C}\|_2 \leq 2 \|C_0\|_2 .$$

Call D a $n \times n$ matrix whose entries are i.i.d with the distribution given at the beginning. Note that for all i, j , and k integer, $0 \leq \mathbf{E} [C_0^k(i, j)] \leq \mathbf{E} [D^k(i, j)]$. Therefore, the previous lemma gives us

$$\mathbf{E} [\text{trace} ((C_0^T C_0)^{2k})] \leq \mathbf{E} [\text{trace} ((D^T D)^{2k})] .$$

Hence,

$$\mathbf{E} [\|C_0\|_2] \leq \left[\mathbf{E} [\text{trace} ((C_0^T C_0)^{2k})] \right]^{1/2k} \leq n^{1/2k} \left[\mathbf{E} [\|D\|_2^{2k}] \right]^{1/2k} .$$

By convexity, we have, if a and b are non-negative, $(a + b)^{2k} \leq 2^{2k-1}(a^{2k} + b^{2k})$. We also have trivially that $(a + b)^{1/2k} \leq a^{1/2k} + b^{1/2k}$.

Hence, if M_D is a median of $\|D^T D\|_2$, we have

$$\left[\mathbf{E} \left[\|D\|_2^{2k} \right] \right]^{1/2k} \leq 2^{1-1/2k} \left(\left[\mathbf{E} \left[\left| \|D\|_2 - M_D \right|^{2k} \right] \right]^{1/2k} + M_D \right).$$

When the random variable X is such that $P(|X - m| > t) \leq C \exp(-ct^2)$, arguments similar to those of Proposition 1.10 in [Led01] show that, for any $q \geq 1$,

$$\mathbf{E} [|X - m|^q] \leq C \frac{q}{2} \Gamma(q/2) c^{-q/2}.$$

Therefore,

$$\left[\mathbf{E} [|X - m|^q] \right]^{1/q} \leq C^{1/q} \left(\frac{q}{2} \right)^{1/q} (\Gamma(q/2))^{1/q} c^{-1/2}.$$

Applying this result in our context with $k = \lceil \log n \rceil$, and using the fact that $\Gamma(x) \sim ((x-1)/e)^{(x-1)} \sqrt{2\pi x}$ as $x \rightarrow \infty$, we have, after realizing that $(\Gamma(k))^{1/2k} \sim \sqrt{k/e}$, for a certain constant K ,

$$\left[\mathbf{E} \left[\left| \|D\|_2 - M_D \right|^{2k} \right] \right]^{1/2k} \leq K 2 \sqrt{\frac{n}{\gamma(1-p)(\log n)^{1+\varepsilon}}} \sqrt{k/e}.$$

Note that K can be taken arbitrarily close to 1, if k is large enough. In any case, we see that, for a certain constant K , picking $k = \lceil \log n \rceil$,

$$\frac{1}{\sqrt{n}} \left[\mathbf{E} \left[\left| \|D\|_2 - M_D \right|^{2k} \right] \right]^{1/2k} \leq K \frac{(\log n)^{-\varepsilon/2}}{\sqrt{\gamma(1-p)}}.$$

Also, for this choice of k , $n^{1/2k} \leq e^{1/2}$, so we have, for yet another constant K ,

$$\frac{1}{\sqrt{n}} \mathbf{E} [\|C_0\|_2] \leq \frac{M_D}{\sqrt{n}} + K \frac{(\log n)^{-\varepsilon/2}}{\sqrt{\gamma(1-p)}}.$$

Our work in the previous subsection guarantees that $\frac{M_D}{\sqrt{n}}$ remains bounded by something of the order $K(1 + 2\frac{(\log n)^{-\varepsilon/2}}{\sqrt{\gamma(1-p)}})$, so we have established that, for yet another K ,

$$\frac{1}{\sqrt{n}} \mathbf{E} [\|C_0\|_2] \leq K \left(1 + 2\frac{(\log n)^{-\varepsilon/2}}{\sqrt{\gamma(1-p)}} \right).$$

And the result concerning $\mathbf{E} [\|C_0\|_2]$ is shown.

The result concerning the median comes out of the fact that $\|C\|_2$ is a convex $\sqrt{2}$ -Lipschitz function of its entries on or above the diagonal. Hence Talagrand's inequality applies and shows that median and mean are arbitrarily close in the setting we are considering. (See the proof of Theorem A.2 for more details.) \square

Proof of Theorem A.2. For C symmetric, $\|C\|_2$ is a convex $\sqrt{2}$ -Lipschitz function of its elements on or above the diagonal. Hence, by Talagrand's inequality,

$$P(\|C\|_2 - m_C > t) \leq 4 \exp(-\gamma(\log n)^{1+\varepsilon}(1-p)t^2/2).$$

Because we have established control of $\mathbf{E} [\|C\|_2]$ above, we also control m_C and conclude as in the proof of Theorem A.1. \square

A.3. On $\|C\|_2$ when $1/n \leq p \ll (\log n)/n$

A.3.1. Case of symmetric C

At the end of Subsection 2.2, we mentioned a corollary (see below) of the following theorem:

Theorem A.3. *Suppose that $p = (\log n)^{1-\delta}u_n/n$, for a fixed δ in $(0, 1)$ and for a fixed κ , $0 < u_n \leq \kappa$. Suppose also that $np \geq 1$. Suppose further that we can find $v_n > 0$ such that $v_n \rightarrow \infty$, while $v_n = o(\log n, [u_n^{-1}(\log n)^\delta]^{1/4})$. Then*

$$\|C/\sqrt{n}\|_2 \rightarrow \infty \text{ in probability.}$$

Recall that practically, this theorem suggests that if we don't sample enough the matrix M (i.e p is too small), a subsampling approximation to its eigen-properties is not likely to work. Let us now prove it.

Proof. We first remark that the diagonal entries of C do not matter for the result we are trying to show. Let us call D_C the diagonal of C . There are two situations. If $np \rightarrow \infty$, then $\max_i |C_{i,i}/\sqrt{n}| \leq (np)^{-1/2} \rightarrow 0$. Similarly, if np is bounded away from 0 (recall that $np \geq 1$ for us), $\max_i |C_{i,i}/\sqrt{n}| \leq (np)^{-1/2}$, so $\|D_C/\sqrt{n}\|_2$ remains bounded. Since by the triangle inequality,

$$\|C - D_C\|_2 - \|D_C\|_2 \leq \|C\|_2 \leq \|C - D_C\|_2 + \|D_C\|_2 ,$$

we conclude that it is sufficient to show that $\|(C - D_C)/\sqrt{n}\|_2 \rightarrow \infty$ in probability to get the result we seek. Let us call $\tilde{C} = C - D_C$. We note that \tilde{C} is just C where we have replaced the diagonal by zeroes. (As an aside we note that if $np \rightarrow 0$, then $P(\exists i : C_{i,i} = \sqrt{(1-p)/p}) = 1 - (1-p)^n \rightarrow 0$, so $\|D_C\|_2 = \sqrt{p/(1-p)}$ with probability one.)

Our strategy is to show that the largest diagonal entry of $\tilde{C}^T \tilde{C}/n$ goes to infinity. To do so, we will rely on results in random graph theory. Let us examine more closely this diagonal. Using the definition of C , we see that, if $T = \tilde{C}^T \tilde{C}$, and d_i is the number of times $\sqrt{(1-p)/p}$ appears in the i -th column of \tilde{C} ,

$$T(i, i) = \frac{(n-1)p}{1-p} + d_i \left(\frac{1-p}{p} - \frac{p}{1-p} \right) .$$

Now $\{d_i\}$ is the degree sequence of an Erdős-Renyi random graph. According to [Bol01], Theorem 3.1, if k is such that $n \binom{n-1}{k} p^k (1-p)^{n-1-k} \rightarrow \infty$, then, if X_k is the number of vertices with degree greater than k ,

$$\lim_{n \rightarrow \infty} P(X_k \geq t) = 1 ,$$

for any t . So if we can exhibit such a k , then $\max d_i \geq k$ with probability going to 1. We now note that for small p ,

$$\left(\frac{1-p}{p} - \frac{p}{1-p} \right) \geq \frac{1}{2p} .$$

Hence, if our k is also such that $k/pn \rightarrow \infty$, we will indeed have

$$\max_i \frac{T(i, i)}{n} \rightarrow \infty$$

and the theorem will be proved.

We propose to take $k = np(1 + v_n)$, where we choose v_n such that k is integer (which can be done without problems, as the arguments below rely only on the order of magnitude of v_n). According to [Bol01], Theorem 1.5, if $h = k - np$, $np \geq 1$, and $q = 1 - p$,

$$\binom{n}{k} p^k (1-p)^{n-k} \geq \frac{1}{\sqrt{2\pi pqn}} \exp\left(-\frac{h^2}{2pqn} - \frac{h^3}{2q^2n^2} - \frac{h^4}{3p^3n^3} - \frac{h}{pn} - \beta\right), \tag{A.4}$$

where $\beta = 1/(12k) + 1/(12(n-k))$. In our case, $h = npv_n$. Let us show that all the terms in the exponential are negligible compared to $\log n$ as $n \rightarrow \infty$:

- $\beta \rightarrow 0$ because $k \rightarrow \infty$ and $npv_n = o((\log n)^{2-\delta})$, given that $v_n = o(\log n)$. Hence $n - k \rightarrow \infty$.
- $h/(pn) = v_n = o(\log n)$ by assumption.
- $h^4/(pn)^3 = npv_n^4 = o(u_n(\log n)^{1-\delta}(\log n)^\delta/u_n) = o(\log n)$, since $v_n = o((u_n^{-1}(\log n)^\delta)^{1/4})$.
- $h^3/n^2 = npv_n^3p^2 = o(npv_n^4p^2) = o(p^2 \log n)$, since $v_n^3 = o(v_n^4)$ ($v_n \rightarrow \infty$ by assumption).
- $h^2/np = npv_n^2 = o(npv_n^4) = o(\log n)$.

In light of these estimates, we have as $n \rightarrow \infty$,

$$\sqrt{n} \exp\left(-\frac{h^2}{2pqn} - \frac{h^3}{2q^2n^2} - \frac{h^4}{3p^3n^3} - \frac{h}{pn} - \beta\right) \rightarrow \infty.$$

Therefore, with this choice of k ,

$$n \binom{n-1}{k} p^k (1-p)^{n-1-k} \rightarrow \infty.$$

We can finally conclude that

$$\max_i T(i, i)/n \geq \frac{k}{2np} \text{ with probability going to } 1.$$

But because $v_n \rightarrow \infty$, we have $k/(2np) \rightarrow \infty$ and the theorem is proved. \square

We have the following corollary to which we appealed in Subsection 2.2.

Corollary A.4. *When $p \sim (\log n)^{1-\delta}/n$ for some fixed $\delta \in (0, 1)$,*

$$\|C/\sqrt{n}\|_2 \rightarrow \infty \text{ in probability.}$$

The previous corollary follows immediately from Theorem A.3, by noticing that u_n is lower bounded under our assumptions and by taking $v_n = (\log n)^{\delta/5}$.

A.3.2. Case of non-symmetric C

In this situation, we can use the same approach as before, namely showing divergence of $\|C/\sqrt{n}\|_2$ by showing that the diagonal of $C^T C/n$ explodes with high-probability. We have the following theorem.

Theorem A.5. *Suppose we are now in the setting where C is a $n \times n$ non-symmetric matrix, with i.i.d entries defined at the beginning of this Appendix. Suppose that $p = (\log n)^{1-\delta} u_n/n$, for a fixed δ in $(0, 1)$ and for a fixed κ , $0 < u_n \leq \kappa$. Suppose also that $np \geq 1$. Suppose further that we can find $v_n > 0$ such that $v_n \rightarrow \infty$, while $v_n = o(\log n, [u_n^{-1}(\log n)^\delta]^{1/4})$. Then*

$$\|C/\sqrt{n}\|_2 \rightarrow \infty \text{ in probability.}$$

Elementary linear algebra shows that the same result holds true if C is $n \times m$, with $m \geq n$.

Proof. We have the same representation as above for the diagonal entries of $T = C^T C$,

$$T(i, i) = \frac{np}{1-p} + d_i \left(\frac{1-p}{p} - \frac{p}{1-p} \right),$$

where now d_i are i.i.d Binomial(n, p). Now, for $t_n \in \mathbf{R}$,

$$P\left(\max_i \frac{d_i}{np} \leq t_n\right) = \prod_{i=1}^n P\left(\frac{d_i}{np} \leq t_n\right) = \left[1 - P\left(\frac{d_i}{np} > t_n\right)\right]^n.$$

So if we can find u_n such $P\left(\frac{d_i}{np} > t_n\right) \geq u_n$, where $nu_n \rightarrow \infty$, we will have $P(\max_i \frac{d_i}{np} \leq t_n) \rightarrow 0$ and therefore

$$P\left(\max_i \frac{d_i}{np} > t_n\right) \rightarrow 1.$$

Therefore $\max_i T_{i,i}/n > t_n$ in this situation asymptotically in probability, and if t_n can be chosen to go to ∞ we are done.

Note that $P\left(\frac{d_i}{np} > t_n\right) \geq P(d_i = \lceil t_n np + 1 \rceil)$, assuming $t_n np + 1 \leq n$ (which will not be a problem below).

If we take t_n such that $k = t_n np + 1$, where $k = np(v_n + 1) \in \mathbb{N}$ as above, we have $P(d_i = t_n np + 1) = \binom{n}{k} p^k q^{n-k} = b(k; n, p)$. Note also that $t_n \sim v_n \rightarrow \infty$. We have seen in the previous proof that under our conditions, $nb(k; n, p) \rightarrow \infty$. So we conclude that

$$P(\max_i T_{i,i}/n > t_n) \rightarrow 1,$$

and hence

$$\|C/\sqrt{n}\|_2 \rightarrow \infty \text{ in probability.}$$

□

Appendix B: Other results

B.1. On $\|C\|_2$ for moderate p

Let us consider the symmetric random matrix C with entries distributed as, for $i \geq j$,

$$C_{i,j} = \begin{cases} \sqrt{\frac{1-p}{p}} & \text{with probability } p \\ -\sqrt{\frac{p}{1-p}} & \text{with probability } 1-p \end{cases} . \quad (\text{B.1})$$

We assume that C is $n \times n$. Our aim is to show that we can control $\|C\|_2$ and in particular its deviation around its median. We do so by using Talagrand's inequality.

We have the following theorem.

Theorem B.1. *Suppose that we observe n matrices C_{α_i} , for $1 \leq i \leq n$ with entries distributed as those of the matrix C just described. Suppose these matrices are of size n^{α_i} , where α_i are positive numbers. Call $\alpha_{\min} = \min_{1 \leq i \leq n} \alpha_i$.*

1. *Assume that, for some fixed $\delta > 0$, $\alpha_{\min} > (\log n)^{(\delta-3)/4}$. Suppose further that p is such that $\lim_{n \rightarrow \infty} (\alpha_{\min} \log n)^4 / (n^{\alpha_{\min}} p) = 0$. Then*

$$\limsup_{n \rightarrow \infty} \max_{i \in \{1, \dots, n\}} \left\| \frac{C_{\alpha_i}}{n^{\alpha_i/2}} \right\|_2 \leq 2 \text{ a.s.} . \quad (\text{B.2})$$

2. *If $pn^{\alpha_{\min}} / (\log n^{\alpha_{\min}})$ is bounded below by $\gamma > 0$, α_{\min} is such that $n^{\alpha_{\min}} \rightarrow \infty$ and $p < 1/2$, we have, for some finite $\mathfrak{K}(\gamma)$,*

$$\limsup_{n \rightarrow \infty} \max_{i \in \{1, \dots, n\}} \left\| \frac{C_{\alpha_i}}{n^{\alpha_i/2}} \right\|_2 \leq \mathfrak{K}(\gamma) \text{ a.s.} . \quad (\text{B.3})$$

Moreover, $\mathfrak{K}(\gamma)$ is of the form $(1 + 2/\sqrt{\gamma})\mathfrak{K} + 8/\sqrt{\gamma\alpha_{\min}}$ for some universal \mathfrak{K} . Naturally, if

$$\liminf pn^{\alpha_{\min}} / (\log n^{\alpha_{\min}}) = \infty \text{ and } \gamma\alpha_{\min} \rightarrow \infty ,$$

then $\mathfrak{K}(\gamma)$ can be replaced by \mathfrak{K} .

Proof. The proof is in two steps. We first show that we can control the deviation of the random quantity $\|C_{\alpha_i}\|_2$ around its median, uniformly in i . Then we show that we can control the corresponding medians.

- **Control of the deviation of $\|C_{\alpha_i}\|_2$ around its median.**

We note that the application $C \rightarrow \|C\|_2$ is a convex, $\sqrt{2}$ -Lipschitz (with respect to Euclidian/Frobenius norm) function of the entries of C that are on or above the main diagonal. As a matter of fact, since $\|\cdot\|$ is a norm, it is convex. Furthermore, if A and B are two symmetric matrices,

$$\|A - B\|_2 \leq \|A - B\|_F = \sqrt{\sum_{i,j} (a_{i,j} - b_{i,j})^2} \leq \sqrt{2} \sqrt{\sum_{i \leq j} (a_{i,j} - b_{i,j})^2}$$

Now recall the consequence of Talagrand's inequality [Tal95] spelled out in [Led01], Corollary 4.10 and Equation (4.10): if F is a convex, 1-Lipschitz function (with respect to Euclidian norm) on \mathbf{R}^n , of n independent random variables (X_1, \dots, X_n) that take value in $[u, v]$, and if m_F is a median of $F(X_1, \dots, X_n)$, then

$$P(|F - m_F| > t) \leq 4 \exp(-t^2/[4(u - v)^2]). \tag{B.4}$$

The random variables that are above the main diagonal of C are bounded, and take value in $[-\sqrt{\frac{p}{1-p}}, \sqrt{\frac{1-p}{p}}]$. We note that

$$\left(\sqrt{\frac{1-p}{p}} + \sqrt{\frac{p}{1-p}}\right)^2 = \frac{1}{p(1-p)}.$$

Therefore, calling m_n the median of $\|n^{-1/2}C\|_2$, we have, in light of Equation (B.4),

$$P\left(\left\|\frac{C}{n^{1/2}}\right\|_2 - m_n > t\right) \leq 4 \exp\left(-\frac{nt^2}{8/(p(1-p))}\right) = 4 \exp\left(-\frac{t^2}{8}p(1-p)n\right). \tag{B.5}$$

Suppose now that we have a collection C_{α_i} of matrices of size n^{α_i} with entries distributed as in Equation (B.1). (Note that the matrices could be dependent.) Let us call $m_{n^{\alpha_i}}$ the median of $\|C_{\alpha_i}/n^{\alpha_i/2}\|_2$. Then we have, by a simple union bound argument, for any k ,

$$\begin{aligned} P\left(\max_{1 \leq i \leq k} \left\|\frac{C_{\alpha_i}}{n^{\alpha_i/2}}\right\|_2 - m_{n^{\alpha_i}} > t\right) &\leq 4 \sum_{i=1}^k \exp\left(-\frac{t^2}{8}p(1-p)n^{\alpha_i}\right) \\ &\leq 4k \exp\left(-\frac{t^2}{8}p(1-p)n^{\alpha_{\min}}\right), \end{aligned}$$

where $\alpha_{\min} = \min_{1 \leq i \leq k} \alpha_i$.

Suppose now that $k = n$, $p \leq 1/2$, $pn^{\alpha_{\min}} \geq \gamma \alpha_{\min}^\zeta (\log n)^\zeta$, for some $\zeta \geq 1$. We assume that γ is bounded below. (Also, implicitly, γ is indexed by n and is allowed to potentially go to ∞ .) For any $\eta > 0$, let

$$t_n = 4 \frac{\sqrt{2+\eta}}{\sqrt{\gamma \alpha_{\min}^\zeta}} (\log n)^{(1-\zeta)/2}.$$

Since $pn^{\alpha_{\min}}(1-p) \geq \gamma \alpha_{\min}^\zeta (\log n)^\zeta / 2$, we have

$$\log n - \frac{t_n^2}{8}pn^{\alpha_{\min}}(1-p) \leq (\log n)(1 - (2 + \eta)).$$

Hence,

$$n \exp\left(-\frac{t_n^2}{8}pn^{\alpha_{\min}}(1-p)\right) \leq \frac{1}{n^{1+\eta}},$$

and since this is the general term of a converging series, we have, when $p \leq 1/2$ and $pn^{\alpha_{\min}} \geq \gamma\alpha_{\min}^\zeta(\log n)^\zeta$

$$\max_{1 \leq i \leq n} \left\| \frac{C_{\alpha_i}}{n^{\alpha_i/2}} \right\|_2 - m_{n^{\alpha_i}} < t_n \text{ a.s. ,}$$

by a simple application of the Borel-Cantelli lemma. Hence, we have, for any $\eta > 0$,

$$\max_{1 \leq i \leq n} \left\| \frac{C_{\alpha_i}}{n^{\alpha_i/2}} \right\|_2 \leq \max_{1 \leq i \leq n} m_{n^{\alpha_i}} + 4 \frac{\sqrt{2+\eta}}{\sqrt{\gamma\alpha_{\min}^\zeta}} (\log n)^{(1-\zeta)/2} \text{ a.s. .} \tag{B.6}$$

We note that when $p \rightarrow 0$, which is the setting that is of interest to us, the quantity $4\sqrt{2+\eta}$ can be replaced by $2^{3/2}\sqrt{2+\eta}$.

* Setting of Part 1:

In the setting of Part 1 of the theorem, we have $\gamma = \gamma_n \rightarrow \infty$, $\zeta = 4$ and $(\log n)^{1-\zeta}/\alpha_{\min}^\zeta \leq (\log n)^{-\delta}$. Hence, $t_n \rightarrow 0$ and so that

$$\max_{1 \leq i \leq n} \left\| \frac{C_{\alpha_i}}{n^{\alpha_i/2}} \right\|_2 - m_{n^{\alpha_i}} \rightarrow 0 \text{ a.s. ,}$$

and somewhat heuristically,

$$\max_{1 \leq i \leq n} \left\| \frac{C_{\alpha_i}}{n^{\alpha_i/2}} \right\|_2 \simeq \max_{1 \leq i \leq n} m_{n^{\alpha_i}} \text{ a.s. .}$$

* Setting of Part 2:

In that setting, we have $\zeta = 1$, so we conclude that

$$\max_{1 \leq i \leq n} \left\| \frac{C_{\alpha_i}}{n^{\alpha_i/2}} \right\|_2 \leq \max_{1 \leq i \leq n} m_{n^{\alpha_i}} + 4 \frac{\sqrt{2+\eta}}{\sqrt{\gamma\alpha_{\min}}} \text{ a.s. .}$$

• **Controlling the medians**

* Setting of Part 1:

Recall Vu's Theorem 1.4 in [Vu07], applied to our situation where we are dealing with bounded random variables with mean 0 and variance 1: if the matrix C has entries as above and is $n \times n$, then almost surely,

$$\left\| \frac{C}{n^{1/2}} \right\|_2 \leq 2 + \kappa_0 \left(\frac{1-p}{p} \right)^{1/4} n^{-1/4} \log(n) ,$$

for some constant κ_0 . So as soon as $(\log n)^4/(pn)$ remains bounded, so does m_n , the median of $\left\| \frac{C}{n^{1/2}} \right\|_2$. In particular, if $(\log n)^4/(pn) \rightarrow 0$, we have

$$\limsup_{n \rightarrow \infty} m_n \leq 2 .$$

Using elementary properties of the function f such that $f(t) = (\log t)^4/t$, we can therefore conclude that if α_{\min} is such that

$$\frac{(\alpha_{\min} \log n)^4}{n^{\alpha_{\min}} p} \rightarrow 0 ,$$

we have

$$\limsup_{n \rightarrow \infty} \max_{1 \leq i \leq n} m_{n^{\alpha_i}} \leq 2.$$

(Note that this is true because we are taking the maximum of elements of a fixed deterministic sequence that is asymptotically less than or equal to $2 + \varepsilon$, for any ε and the smallest argument is going to infinity. All the work using Talagrand's inequality was done to allow us to switch from having to control the maximum of a random sequence to that of a deterministic sequence.)

Now when $(\alpha_{\min} \log n)^4 / (pn^{\alpha_{\min}}) \rightarrow 0$, we have a fortiori $pn^{\alpha_{\min}} > (\log n)^{1+\delta}$ when $\alpha_{\min} > (\log n)^{(\delta-3)/4}$. So we conclude that when $(\alpha_{\min} \log n)^4 / (pn^{\alpha_{\min}}) \rightarrow 0$ and $\alpha_{\min} > (\log n)^{(\delta-3)/4}$,

$$\limsup_{n \rightarrow \infty} \max_{1 \leq i \leq n} \left\| \frac{C_{\alpha_i}}{n^{\alpha_i/2}} \right\|_2 \leq 2 \text{ a.s. .}$$

This shows Part 1 of the theorem.

* Setting of Part 2:

We have established in Lemma A.2 that as $n \rightarrow \infty$, if $p \geq \gamma(\log n)/n$, we have, for a universal constant \mathfrak{K} ,

$$m_n \leq (1 + 2\gamma^{-1/2})\mathfrak{K}.$$

Hence, when $n^{\alpha_{\min}} \rightarrow \infty$ and $p \geq \gamma(\log n^{\alpha_{\min}})/n^{\alpha_{\min}}$, we also have

$$\limsup_{n \rightarrow \infty} \max_{1 \leq i \leq n} m_{n^{\alpha_i}} \leq (1 + 2\gamma^{-1/2})\mathfrak{K}.$$

And we conclude as before to show Part 2. □

Let us now consider the related issue of understanding the matrix $E = r_p M \circ C$, where $r_p = \sqrt{(1-p)/p}$, M is a deterministic matrix and C is a random matrix as above.

Theorem B.2. *Suppose $E = r_p M \circ C$, where C is a symmetric random matrix distributed as above, M is a deterministic matrix and $r_p = \sqrt{(1-p)/p}$. Let us call m_E a median of $\|E\|_2$. Then we have*

$$P(|\|E\|_2 - m_E| > t) \leq 4 \exp\left(-\frac{p^2}{8\|M\|_\infty^2} t^2\right).$$

Hence, in particular,

$$\mathbf{E} \left[\|E\|_2^2 \right] \leq m_E^2 + 32 \frac{\|M\|_\infty^2}{p^2} + 8m_E \sqrt{\frac{2\pi\|M\|_\infty^2}{p^2}}. \tag{B.7}$$

and

$$\mathbf{E}[\|E\|_2^3] \leq 4m_E^3 + 12\sqrt{\pi} \left(\frac{8\|M\|_\infty^2}{p^2} \right)^{3/2}. \tag{B.8}$$

We note in connection with this theorem that if M is symmetric (and $n \times n$) and can be written in spectral form as $M = \sum \lambda_i u_i u_i^T$, we clearly have $\|M\|_\infty \leq \sum_{i=1}^n |\lambda_i| \|u_i\|_\infty^2$, and hence, if μ is the incoherence of the matrix M , we also have

$$\|M\|_\infty \leq \sum_{i=1}^n |\lambda_i| \|u_i\|_\infty^2 \leq \mu n^{-\alpha_{\min}}. \tag{B.9}$$

Proof of Theorem B.2. The crux of the proof is quite similar to that of Theorem B.1: we will rely on Talagrand's concentration inequality for convex 1-Lipschitz functions of bounded random variables. To do so let us consider the map: $C \rightarrow f(C) = \|M \circ C\|_2$. This map f is convex as the composition of a norm with an affine mapping. Let us now show that it is $(\sqrt{2} \|M\|_\infty)$ -Lipschitz with respect to Euclidian norm: if we denote by $c_{i,j}^{(k)}$ the (i, j) -th entry of the matrix C_k , we have

$$\begin{aligned} |f(C_1) - f(C_2)| &= \left| \|M \circ C_1\|_2 - \|M \circ C_2\|_2 \right| \leq \|M \circ (C_1 - C_2)\|_2 \\ &\leq \|M \circ (C_1 - C_2)\|_F = \sqrt{\sum_{i,j} M_{i,j}^2 (c_{i,j}^{(1)} - c_{i,j}^{(2)})^2} \\ &\leq \max_{i,j} |M_{i,j}| \sqrt{\sum_{i,j} (c_{i,j}^{(1)} - c_{i,j}^{(2)})^2} \leq \|M\|_\infty \sqrt{2} \sqrt{\sum_{i \leq j} (c_{i,j}^{(1)} - c_{i,j}^{(2)})^2} \end{aligned}$$

Hence, f is indeed a $(\sqrt{2} \|M\|_\infty)$ -Lipschitz function of the entries of C that are above or on the diagonal. Now the function of C we care about is $g(\cdot) = r_p f(\cdot)$, which is convex and $\sqrt{2} \|M\|_\infty r_p$ - Lipschitz. Given that the entries of C are bounded, we have, as in the proof of Theorem B.1,

$$P(|\|E\|_2 - m_E| > t) \leq 4 \exp\left(-\frac{p(1-p)}{8r_p^2 \|M\|_\infty^2} t^2\right) = 4 \exp\left(-\frac{p^2}{8 \|M\|_\infty^2} t^2\right).$$

Now using the proof of Proposition 1.9 in [Led01] (see p.12 of this book), we conclude that

$$\begin{aligned} \mathbf{E} [|\|E\|_2 - m_E|] &\leq 4 \sqrt{\frac{2\pi \|M\|_\infty^2}{p^2}}, \text{ and} \\ \mathbf{E} [|\|E\|_2 - m_E|^2] &\leq 32 \frac{\|M\|_\infty^2}{p^2}. \end{aligned}$$

Therefore,

$$\mathbf{E} [\|E\|_2^2] \leq m_E^2 + 32 \frac{\|M\|_\infty^2}{p^2} + 8m_E \sqrt{\frac{2\pi \|M\|_\infty^2}{p^2}},$$

since for a and b positive, $a^2 \leq b^2 + (a - b)^2 + 2b|a - b|$.

More generally, we see, using essentially Proposition 1.10 in [Led01] and elementary properties of the Gamma function, that if the random variable F is such that for a deterministic number a_F , $P(|F - a_F| > t) \leq C \exp(-cr^2)$, then

$$\mathbf{E}[|F - a_F|^k] \leq CT \left(\frac{k}{2} + 1\right) c^{-k/2}.$$

Applying this result with $k = 3$, we get

$$\mathbf{E}[|\|E\|_2 - m_E|^3] \leq 3\sqrt{\pi} \left(\frac{8\|M\|_\infty^2}{p^2}\right)^{3/2}.$$

In our context, using the fact that, for positive a and b , $(a + b)^3 \leq 4(a^3 + b^3)$ by convexity, we also have

$$\mathbf{E}[\|E\|_2^3] \leq 4 \left(m_E^3 + 3\sqrt{\pi} \left(\frac{8\|M\|_\infty^2}{p^2}\right)^{3/2}\right).$$

□

B.2. Regularized eigenvector considerations

We now have the following (regularized) second order accuracy result, which is a critical component of the proof of Theorem 3, one of the main results of the paper.

The results above allowed us to control $\mathbf{E}[\|E\|_2^2]$, $\mathbf{E}[\|E\|_2^3]$ and these moment bounds, together with $\mathbf{E}[E] = 0$ of course and the bound on $\|M\|_\infty$ in (B.9), are all we need to show that averaging produces second order accurate eigenvector approximations.

Theorem B.3. *Suppose that the assumptions of Theorem 1 are satisfied. We consider the approximation of u the eigenvector associated with the largest eigenvalue of M . Recall that v is the eigenvector corresponding to the leading eigenvalue of the subsampled matrix S . For $\varepsilon > 0$, we call \tilde{v}_ε the vector such that*

$$\tilde{v}_\varepsilon = \begin{cases} v & \text{if } \|(Id + \Delta)^{-1}\|_2 \leq \frac{1}{\varepsilon} \\ u - REu + \Delta REu & \text{otherwise} \end{cases}.$$

Then, we have

$$\|\mathbf{E}[\tilde{v}_\varepsilon - u]\|_2 \leq \mathbf{E} \left[\frac{4\|R\|_2^3\|E\|_2^3}{\varepsilon} + 2\|R\|_2^2\|E\|_2^2 \right].$$

Also, for any $\eta > 0$, we have asymptotically,

$$\|\mathbf{E}[u - \tilde{v}_\varepsilon]\|_2 \leq \frac{2(\mathfrak{K}(\gamma))^2 + \eta}{(\lambda_1 - \lambda_2)^2} \frac{\mu^2}{pn^{\alpha_{\min}}} + \frac{16(\mathfrak{K}(\gamma))^3 + \eta}{\varepsilon(\lambda_1 - \lambda_2)^3} \left(\frac{\mu^2}{pn^{\alpha_{\min}}}\right)^{3/2}.$$

Suppose further that we are in an asymptotic setting where $\frac{1}{\lambda_1 - \lambda_2} \frac{\mu}{(pn^{\alpha_{\min}})^{1/2}} \rightarrow 0$. Then, $v - \tilde{v}_\varepsilon = 0$ with high-probability.

Proof. Let us first explain that our regularization does not change the vector we are dealing with with high-probability. $\tilde{v}_\varepsilon = v$ as long as $\|(\text{Id} + \Delta)^{-1}\|_2 \leq 1/\varepsilon$, which is guaranteed if $2\|E\|_2/d \leq 1 - \varepsilon$, where $d = \lambda_1 - \lambda_2$. When we assume that $\frac{1}{\lambda_1 - \lambda_2} \frac{\mu}{(pn^{\alpha_{\min}})^{1/2}} \rightarrow 0$, since we have according to Theorem B.2 $\|E\|_2 \leq \mathfrak{K}(\gamma) \frac{\mu}{(pn^{\alpha_{\min}})^{1/2}}$ with high-probability, we conclude that with high-probability, $\tilde{v}_\varepsilon = v$ (provided $\mathfrak{K}(\gamma)$, whose form is made explicit there, stays bounded).

Using Equation (9) with $j = 1$, we see that, since $\|\Delta\|_2 \leq 2\|R\|_2\|E\|_2$,

$$\|\tilde{v}_\varepsilon - (u - REu + \Delta REu)\|_2 \leq \frac{1}{\varepsilon} \|\Delta\|_2^2 \|RE\|_2 \leq \frac{4\|R\|_2^3 \|E\|_2^3}{\varepsilon}.$$

Recall that by construction $\mathbf{E}[E] = 0$. Hence, since R is a fixed deterministic matrix and u is a deterministic vector,

$$\mathbf{E}[\tilde{v}_\varepsilon - u] = \mathbf{E}[\tilde{v}_\varepsilon - u + REu].$$

So, if we now use the fact that $\|u\| = 1$, we have

$$\begin{aligned} \|\mathbf{E}[\tilde{v}_\varepsilon - u]\|_2 &= \|\mathbf{E}[\tilde{v}_\varepsilon - u + REu]\|_2 \\ &\leq \|\mathbf{E}[\tilde{v}_\varepsilon - u + REu - \Delta REu]\|_2 + \|\mathbf{E}[\Delta REu]\|_2 \\ &\leq \mathbf{E}[\|\tilde{v}_\varepsilon - u + REu - \Delta REu\|_2] + \mathbf{E}[\|\Delta REu\|_2] \\ &\leq \mathbf{E}\left[\frac{4\|R\|_2^3 \|E\|_2^3}{\varepsilon} + 2\|R\|_2^2 \|E\|_2^2\right]. \end{aligned}$$

This proves the first result of the theorem. Let us now show that we can control the right-hand side of the previous equation.

We prove in Theorem B.2 that

$$\mathbf{E}\left[\|E\|_2^2\right] \leq m_E^2 + 32 \frac{\|M\|_\infty^2}{p^2} + 8m_E \sqrt{\frac{2\pi \|M\|_\infty^2}{p^2}},$$

where m_E is a median of the random variable $\|E\|_2$. Our asymptotic control of $\|E\|_2$ in (5) and (6) allows us to control m_E , namely,

$$\limsup_{n \rightarrow \infty} m_E^2 \leq (\mathfrak{K}(\gamma))^2 \frac{\mu^2}{pn^{\alpha_{\min}}}.$$

In other respects, following (B.9), we clearly have $\|M\|_\infty \leq \sum_{i=1}^n |\lambda_i| \|u_i\|_\infty^2$, and hence

$$\|M\|_\infty \leq n^{-\alpha_{\min}} \mu.$$

Hence,

$$\frac{\|M\|_\infty^2}{p^2} \leq \frac{\mu^2}{(pn^{\alpha_{\min}})^2} = o\left(\frac{\mu^2}{pn^{\alpha_{\min}}}\right),$$

since we are in a setting where $pn^{\alpha_{\min}} \rightarrow \infty$. Similarly, $m_E \sqrt{\frac{\|M\|_{\infty}^2}{p^2}} = o\left(\frac{\mu^2}{pn^{\alpha_{\min}}}\right)$, so we have for $\eta > 0$,

$$2 \|R\|_2^2 \mathbf{E} \left[\|E\|_2^2 \right] \leq \frac{2(\mathfrak{K}(\gamma))^2 + \eta}{(\lambda_1 - \lambda_2)^2} \frac{\mu^2}{pn^{\alpha_{\min}}}$$

asymptotically.

Furthermore, we prove in Theorem B.2 that

$$\mathbf{E}[\|E\|_2^3] \leq 4m_E^3 + 12\sqrt{\pi} \left(\frac{8 \|M\|_{\infty}^2}{p^2} \right)^{3/2} \leq 4m_E^3 + o\left(\left(\frac{\mu^2}{pn^{\alpha_{\min}}} \right)^{3/2} \right).$$

Hence, for $\eta > 0$,

$$4 \|R\|_2^3 \mathbf{E}[\|E\|_2^3] \leq \frac{16(\mathfrak{K}(\gamma))^3 + \eta}{(\lambda_1 - \lambda_2)^3} \left(\frac{\mu^2}{pn^{\alpha_{\min}}} \right)^{3/2}.$$

□

B.3. Variance computations

We provide some details here to complement the explanations we gave in the proof of Theorem 4 in Subsection 2.6.

On $\mathbf{E}[E^2]$ Let us explain why this matrix is diagonal and compute the coefficients on the diagonal. Recall that $E = \sqrt{(1-p)/p} M \circ C$, where C is a random matrix whose above-diagonal elements are independent, have mean 0 and variance 1. E is naturally symmetric and we call E_i its i -th column. Naturally, $E^2(i, j) = E_i^T E_j$. Suppose first that $i \neq j$. The elements of E_i and E_j are independent, except for E_{ij} and E_{ji} , which are equal. In particular, E_{ki} and E_{kj} are independent for all $1 \leq k \leq n$. Recall also that $\mathbf{E}[C] = 0$, so $\mathbf{E}[E] = 0$. Combining all these elements, we conclude that, if $i \neq j$,

$$\mathbf{E}[E_i^T E_j] = \sum_{k=1}^n \mathbf{E}[E_{ki} E_{kj}] = \sum_{k=1}^n \mathbf{E}[E_{ki}] \mathbf{E}[E_{kj}] = 0.$$

Therefore $\mathbf{E}[E^2]$ is diagonal. Let us now turn our attention to computing the elements of the diagonal. This is simple since

$$\mathbf{E}[E_i^T E_i] = \frac{1-p}{p} \sum_{k=1}^n M_{ki}^2 \mathbf{E}[E_{ki}^2] = \frac{1-p}{p} \sum_{k=1}^n M_{ki}^2 = \frac{1-p}{p} \|M_i\|_2^2.$$

This is the result we announced in the proof of Theorem 4 in Subsection 2.6.

On $\text{var}(u^T Eu)$ Rewriting this quantity as a sum of independent quantities greatly simplifies the computation. If we pursue this route, we have

$$u^T Eu = \sum_{i,j} u(i)u(j)E_{ij} = 2 \sum_{i>j} u(i)u(j)E_{ij} + \sum_i u(i)^2 E_{ii} .$$

Because the previous expression is a sum of independent random variables, we immediately conclude that

$$\begin{aligned} \frac{p}{1-p} \text{var}(u^T Eu) &= 4 \sum_{i>j} u(i)^2 u(j)^2 M_{ij}^2 + \sum_i u(i)^4 M_{ii}^2 \\ &= 2 \left(2 \sum_{i>j} u(i)^2 u(j)^2 M_{ij}^2 + \sum_i u(i)^4 M_{ii}^2 \right) - \sum_i u(i)^4 M_{ii}^2 . \end{aligned}$$

Calling $w = u \circ u$ and $\mathcal{M} = M \circ M$, we immediately recognize in the last expression the quantity

$$2(w^T \mathcal{M} w) - \sum_k w(k)^2 \mathcal{M}_{kk} ,$$

as announced in the proof of Theorem 4.

References

- [ABN⁺99] A. ALON, N. BARKAI, D.A. NOTTERMAN, K. GISH, S. YBARRA, D. MACK, and A.J. LEVINE. Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Cell Biology*, 96:6745–6750, 1999.
- [AM07] D. ACHLIOPTAS and F. MCSHERRY. Fast computation of low-rank matrix approximations. *Journal of the ACM*, 54(2), 2007. [MR2295993](#)
- [And03] T.W. ANDERSON. *An introduction to multivariate statistical analysis*. Wiley Series in Probability and Statistics. Wiley-Interscience [John Wiley & Sons], Hoboken, NJ, third edition, 2003. [MR1990662](#)
- [BC06] L. BECCHETTI and C. CASTILLO. The distribution of PageRank follows a power-law only for particular values of the damping factor. In *World Wide Web Conference*, pages 941–942. ACM New York, NY, USA, 2006.
- [Bha97] RAJENDRA BHATIA. *Matrix analysis*, volume 169 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, 1997. [MR1477662](#)
- [Bol01] BÉLA BOLLOBÁS. *Random graphs*, volume 73 of *Cambridge Studies in Advanced Mathematics*. Cambridge University Press, Cambridge, second edition, 2001. [MR1864966](#)
- [BV04] PAOLO BOLDI and SEBASTIANO VIGNA. The WebGraph framework I: Compression techniques. In *Proc. of the Thirteenth International World Wide Web Conference (WWW 2004)*, pages 595–601, Manhattan, USA, 2004. ACM Press.

- [CR08] E.J. CANDÉS and B. RECHT. Exact matrix completion via convex optimization. *preprint*, 2008. [MR2565240](#)
- [CT09] E.J. CANDÉS and T. TAO. The Power of Convex Relaxation: Near-Optimal Matrix Completion. [arXiv:0903.1476](#), 2009.
- [DKM06] P. DRINEAS, R. KANNAN, and M.W. MAHONEY. Fast Monte Carlo Algorithms for Matrices II: Computing a Low-Rank Approximation to a Matrix. *SIAM Journal on Computing*, 36:158, 2006. [MR2231644](#)
- [FKV04] A. FRIEZE, R. KANNAN, and S. VEMPALA. Fast monte-carlo algorithms for finding low-rank approximations. *Journal of the ACM (JACM)*, 51(6):1025–1041, 2004. [MR2145262](#)
- [GMKG91] D.J. GROH, R.A. MARSHALL, A.B. KUNZ, and C.R. GIVENS. An approximation method for eigenvectors of very large matrices. *Journal of Scientific Computing*, 6(3):251–267, 1991. [MR1154901](#)
- [GVL90] G.H. GOLUB and C.F. VAN LOAN. Matrix computation. *North Oxford Academic*, 1990.
- [GZB04] D. GLEICH, L. ZHUKOV, and P. BERKHIN. Fast parallel PageRank: A linear system approach. *Yahoo! Research Technical Report YRL-2004-038*, 2004.
- [HJ91] R.A. HORN and C.R. JOHNSON. *Topics in matrix analysis*. Cambridge university press, 1991. [MR1091716](#)
- [HK03] T.H. HAVELIWALA and S.D. KAMVAR. The Second Eigenvalue of the Google Matrix. *Stanford CS Tech report*, 2003.
- [HTF⁺01] T. HASTIE, R. TIBSHIRANI, J. FRIEDMAN, et al. *The elements of statistical learning: data mining, inference, and prediction*. Springer, 2001. [MR1851606](#)
- [HYJT08] L. HUANG, D. YAN, M.I. JORDAN, and N. TAFT. Spectral Clustering with Perturbed Data. *Advances in Neural Information Processing Systems (NIPS)*, 2008.
- [Kat95] T. KATO. *Perturbation theory for linear operators*. Springer, 1995. [MR1335452](#)
- [KMO09] R.H. KESHAVAN, A. MONTANARI, and S. OH. Matrix Completion from a Few Entries. [arXiv:0901.3150](#), 2009.
- [KV09] R. KANNAN and S. VEMPALA. *Spectral algorithms*. 2009. [MR2558901](#)
- [Led01] M. LEDOUX. *The concentration of measure phenomenon*, volume 89 of *Mathematical Surveys and Monographs*. American Mathematical Society, Providence, RI, 2001. [MR1849347](#)
- [LM05] A.N. LANGVILLE and C.D. MEYER. A survey of eigenvector methods for web information retrieval. *SIAM Review*, 47(1):135–161, 2005. [MR2149104](#)
- [Mel07] MASSIMO MELUCCI. On rank correlation in information retrieval evaluation. *SIGIR Forum*, 41(1):18–33, 2007.
- [MKB79] KANTILAL VARICHAND MARDIA, JOHN T. KENT, and JOHN M. BIBBY. *Multivariate analysis*. Academic Press [Harcourt Brace Jo-

- vanovich Publishers], London, 1979. Probability and Mathematical Statistics: A Series of Monographs and Textbooks. [MR0560319](#)
- [NJW02] A. NG, M. JORDAN, and Y. WEISS. On spectral clustering: Analysis and an algorithm. In *Advances in Neural Information Processing Systems 14*, page 849. MIT Press, 2002.
- [NZJ01] A.Y. NG, A.X. ZHENG, and M.I. JORDAN. Stable algorithms for link analysis. In *ACM SIGIR*, pages 258–266. ACM New York, NY, USA, 2001.
- [PBMW98] L. PAGE, S. BRIN, R. MOTWANI, and T. WINOGRAD. The pagerank citation ranking: Bringing order to the web. *Stanford CS Technical Report*, 1998.
- [PRTV00] C.H. PAPADIMITRIOU, P. RAGHAVAN, H. TAMAKI, and S. VEMPALA. Latent semantic indexing: a probabilistic analysis. *Journal of Computer and System Sciences*, 61(2):217–235, 2000. [MR1802556](#)
- [PRU06] G. PANDURANGAN, P. RAGHAVAN, and E. UPFAL. Using pagerank to characterize web structure. *Internet Mathematics*, 3(1):1–20, 2006. [MR2283881](#)
- [RFP07] B. RECHT, M. FAZEL, and P.A. PARRILO. Guaranteed Minimum-Rank Solutions of Linear Matrix Equations via Nuclear Norm Minimization. *Arxiv preprint arXiv:0706.4138*, 2007.
- [RV07] MARK RUDELSON and ROMAN VERSHYNIN. Sampling from large matrices: An approach through geometric functional analysis. *J. ACM*, 54(4):21, 2007. [MR2351844](#)
- [Seg00] YOAV SEGNER. The expected norm of random matrices. *Combin. Probab. Comput.*, 9(2):149–166, 2000. [MR1762786](#)
- [Ste98] G.W. STEWART. *Matrix algorithms*. Society for Industrial and Applied Mathematics, 1998. [MR1653546](#)
- [Tal95] MICHEL TALAGRAND. Concentration of measure and isoperimetric inequalities in product spaces. *Inst. Hautes Études Sci. Publ. Math.*, (81):73–205, 1995. [MR1361756](#)
- [Vu07] V.H. VU. Spectral norm of random matrices. *Combinatorica*, 27(6):721–736, 2007. [MR2384414](#)