
On the Power of SVD in Clustering Problems

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A popular heuristic method for improving clustering results is to apply dimensionality
2 reduction before running clustering algorithms. It has been observed that
3 spectral-based dimensionality reduction tools, such as PCA or SVD, improve the
4 performance of clustering algorithms in many applications. This phenomenon
5 indicates that spectral method not only serves as a dimensionality reduction tool,
6 but also contributes to the clustering procedure in some sense. It is an interesting
7 question to understand the behavior of spectral steps in clustering problems.
8 As an initial step (but not the final step) in this direction, this paper studies the
9 power of vanilla-SVD algorithm in the stochastic block model (SBM). We show
10 that, in the symmetric setting, vanilla-SVD algorithm recovers all clusters correctly.
11 This result answers an open question posed by Van Vu (Combinatorics Probability
12 and Computing, 2018) in the symmetric setting.

13 1 Introduction

14 Clustering is a fundamental task in machine learning, with applications in many fields, such as
15 biology, data mining, and statistical physics. Given a set of objects, the goal is to partition them into
16 clusters according to their similarities. Objects and known relations can be represented in various
17 ways. In most cases, objects are represented as vectors in \mathbb{R}^d , forming a data set $\mathcal{D} \subset \mathbb{R}^d$; each
18 coordinate is called a feature, whose value is directly derived from raw data.

19 In many applications, the number of features could be very large. It has been observed that the
20 performance of classical clustering algorithms such as K-means may be worse on high-dimensional
21 datasets. Some people call this phenomenon *curse of dimensionality in machine learning* [SEK04]. A
22 popular heuristic method to address this issue is to apply dimensionality reduction before clustering.
23 Among tools for dimensionality reduction, it is noted in practice that spectral methods such as
24 *principal component analysis* (PCA) and *singular value decomposition* (SVD) significantly improve
25 clustering results, e.g., [SEK04, KAH19].

26 A natural question arises: *why do spectral steps help to cluster high-dimensional datasets?* Some prac-
27 titioners believe one reason is that the spectral method filters some noise from the high-dimensional
28 data [ARS⁺04, SEK04, ZLWZ09, KAH19]. Simultaneously, many theory works also (partially)
29 support this explanation [AFWZ20, EBW18, LZK22, MZ22a]. With this explanation in mind, people
30 analyzed the behavior of spectral-based algorithms with noise perturbation. Based on these analyses,
31 many algorithms were proposed to recover clusters in probabilistic generative models. Among them,
32 a well-studied model is the *signal-plus-noise model*.

33 **Signal-plus-noise model** In this model, we assume that each observed sample \hat{v}_i has the form
34 $\hat{v}_i = v_i + e_i$, where v_i is a ground-truth vector and e_i is a random noise vector. For any two vectors
35 \hat{v}_i, \hat{v}_j , if they are from the same cluster, their corresponding ground-truth vectors are identical, i.e.,
36 $v_i = v_j$. Signal-plus-noise model is very general; it has plentiful variants with different types of

37 ground-truth vectors and noise distribution. In this paper, we focus on an important instance known as
 38 the *stochastic block model* (SBM). Though SBM is not as broad as general signal-plus-noise model, it
 39 usually serves as a *benchmark* for clustering and provides *preliminary intuition* about random graphs.

40 **Stochastic block model** The SBM is first introduced by [HLL83] and is widely used as a theoretical
 41 benchmark for graph clustering algorithms. In the paper, we focus on the *symmetric version of*
 42 *stochastic block model* (SSBM), described as follows. Given a set of n vertices V , we uniformly
 43 partition them into k disjoint sets (clusters), denoted by V_1, \dots, V_k . Based on this partition, a random
 44 graph $\widehat{G} = (V, E)$ is sampled in the following way: for all pairs of vertices $u, v \in V$, an edge (u, v)
 45 is added independently with probability p , if $u, v \in V_\ell$ for some ℓ ; otherwise, an edge (u, v) is added
 46 independently with probability q .

47 We usually assume that $p > q$. The task is to recover the hidden partition V_1, \dots, V_k from the random
 48 graph \widehat{G} . We denote this model as $\text{SSBM}(V, n, k, p, q)$.

SBM as a signal-plus-noise model Though SBM was originally designed for graph clustering,
 we view it as a special form of vector clustering. Namely, given the adjacency matrix of a graph
 $\widehat{G} \in \{0, 1\}^{V \times V}$, the columns of \widehat{G} form a set of $n = |V|$ vectors. To see that SBM fits into the
 signal-plus-noise model, note that in SBM, the adjacency matrix $\widehat{G} \in \{0, 1\}^{V \times V}$ can be view as a
 fixed matrix G plus a random noise, i.e., $\widehat{G} = G + E$, where $G \stackrel{\text{def}}{=} \mathbf{E}[\widehat{G}]$ is the mean and E is a
 zero-mean random matrix. More precisely, in the case of SSBM,

$$G_{uv} = \begin{cases} p, & \text{if } u, v \in V_\ell \text{ for some } \ell; \\ q, & \text{otherwise;} \end{cases} \quad \text{and} \quad E_{uv} = \begin{cases} 1 - G_{uv}, & \text{with probability } G_{uv}; \\ -G_{uv}, & \text{with probability } 1 - G_{uv}. \end{cases}$$

49 1.1 Motivations: Analyzing Vanilla Spectral Algorithms

50 Since the seminal work by McSherry [McS01], many spectral-based algorithms have been proposed
 51 and studied in SBM [GM05, Vu18, LR15, EBW18, Col19, AFWZ20, MZ22b] and even more general
 52 signal-plus-noise models [AFWZ20, EBW18, CTP19, LZK22, MZ22a]. These algorithms are largely
 53 based on the spectral analysis of random matrices. The purpose of designing and analyzing such
 54 algorithms is twofold.

55 **Understand the limitation of spectral-based algorithms.** SBM is specified by parameters, such
 56 as n, k, p, q in the symmetric case. Clustering is usually getting harder for larger k and smaller
 57 gap $(p - q)$. Many existing works aim to understand in which *regimes* of these parameters it is
 58 possible to recover the hidden partition. In this regard, the state-of-the-art bound is given by Vu
 59 [Vu18]. Concretely, [Vu18] proved that, in the symmetric setting, there is an algorithm that recovers
 60 all clusters if $n \geq C \cdot k \left(\frac{\sigma\sqrt{k} + \sqrt{\log n}}{p - q} \right)^2$, where $\sigma^2 \stackrel{\text{def}}{=} \max\{p(1 - p), q(1 - q)\}$ and C is a constant.

61 **Understand spectral-based algorithms in practice.** Besides analyzing spectral algorithms in
 62 theory, the other purpose (*the primary purpose of this paper*) is to explain the usefulness of such
 63 algorithms in practice. Indeed, as we mentioned before, many spectral-based algorithms, as observed
 64 in practice, can filter the noise and address the curse of dimensionality [ARS⁺04, SEK04, ZLWZ09,
 65 KAH19]. Some representative algorithms are PCA and SVD. Furthermore, it has been observed that
 66 spectral algorithms used in practice, such as PCA or SVD, are usually *very simple*: they just project
 67 data points into some lower-dimension subspace, and no extra steps are conducted.

68 In stark contrast, most of the aforementioned theoretical algorithms have pre-processing or post-
 69 processing steps. For example, the idea in [LR15] is that one first applies SVD, and then runs a
 70 variant of K-means to clean up the clustering; the main algorithm in [Vu18] partitions the graph
 71 into several parts and uses these parts in different ways. As noted in [Vu18], these extra steps are
 72 only for the *purpose of theoretical analysis*. The author believed, from the perspective of algorithm
 73 design, these extra steps appear redundant. Later on, [AFWZ20] coined the phrase *vanilla spectral*
 74 *algorithms* to describe spectral algorithms that do not include any additional steps. Both [Vu18] and
 75 [AFWZ20] conjectured that vanilla spectral algorithms are themselves good clustering algorithms. In
 76 practice, this is a widely-used heuristic; however, in theory, the analysis of vanilla spectral algorithms

77 is not satisfactory due to the lack of techniques for analysis. We refer to [MZ22b] for a detailed
 78 discussion on barriers of the current analysis.

79 **Why do we study vanilla algorithms?** Our main focus is particularly on vanilla spectral algorithms
 80 for two reasons:

- 81 1. Vanilla spectral algorithms are the most popular in practice—no extra steps are widely used.
 82 Plus, their performance seems good enough. The lack of theoretical analysis is mostly due
 83 to technical obstacles.
- 84 2. A vanilla spectral algorithm is often simple and is not specifically designed for theoretical
 85 models such as SBM. In contrast, some complicated algorithms use extra steps which are
 86 designed for SBM. These steps made the analysis of SBM go through (as commented by
 87 [Vu18]). Meanwhile, these extra steps exploit specific structures and may cause ‘overfittings’
 88 on SBM, which makes these algorithms less powerful in practice.

89 The main purpose of this paper is to *theoretically understand* the power of *practically successful*
 90 vanilla spectral algorithms. To this end, we study SBM as a preliminary demonstration. We *do not*
 91 aim to design algorithms for SBM that outperforms existing algorithms.

92 1.2 Our Results

93 The contribution of this paper is twofold. On the one hand, we show that vanilla algorithms (alg. 1) is
 94 indeed a clustering algorithm in SSBM for a wide range of parameters, breaking previous barrier on
 95 analyzing on only constant number of clusters. On the other hand, we provide a novel analysis on
 96 matrix perturbation with random noise. We discuss more details on this part in Section 1.4.

97 Recall that parameters of SBM is specified by $\text{SSBM}(V, n, k, p, q)$, where $n = |V|$. Let $\sigma^2 =$
 98 $\max\{p(1-p), q(1-q)\}$. Our main result is stated below.

99 **Theorem 1.1.** *There exists a constant $C > 0$. In the model $\text{SSBM}(V, n, k, p, q)$, if $\sigma^2 \geq C \log n/n$
 100 and $n \geq C \cdot k \left(\frac{\sqrt{kp \cdot \log^6 n + \sqrt{\log n}}}{p-q} \right)^2$, then alg. 1 recovers all clusters with probability $1 - O(n^{-1})$.*

101 Here we describe the vanilla-SVD algorithm in more detail. Algorithms in [McS01, Vu18, Col19]
 102 share a common idea: they both use SVD-based methods to find a clear-cut vector representation
 103 of vertices. That is, every node $v \in V$ is associated with a vector $\rho(v)$, and we say a vector
 104 representation ρ is *clear-cut* if the following holds for some threshold Δ : if $u, v \in V_\ell$ for some ℓ ,
 105 then $\|\rho(u) - \rho(v)\| \leq \Delta/4$; otherwise, $\|\rho(u) - \rho(v)\| \geq \Delta$.

106 Once a clear-cut representation is found, the clustering task is easy. If the parameters n, k, p, q
 107 are all known, we can calculate Δ and simply decide whether two vertices are in the same cluster
 108 based on their distance; in the case where Δ is unknown, we need one more step.¹ Following
 109 [Vu18], we denote by `ClusterByDistance` an algorithm that recovers the partition from a clear-cut
 110 representation. One natural representation is obtained by SVD as follows. Let $\widehat{G} \in \{0, 1\}^{V \times V}$ be the
 111 adjacent matrix of the input graph, and let $P_{\widehat{G}_k}$ be the orthogonal projection matrix onto the space
 112 spanned by the first k eigenvectors of \widehat{G} . Then set $\rho(u) \stackrel{\text{def}}{=} P_{\widehat{G}_k} \widehat{G}_u$, where \widehat{G}_u is the column index
 by $u \in V$. This yields alg. 1, the vanilla-SVD algorithm.

Algorithm 1: Vanilla-SVD algorithm for graph clustering

- 1 Input: adjacent matrix $\widehat{G} \in \{0, 1\}^{V \times V}$
 - 2 Output: a partition of V
 1. Compute $\rho(u) \stackrel{\text{def}}{=} P_{\widehat{G}_k} \widehat{G}_u$ for each $u \in V$.
 2. Run `ClusterByDistance` with representation ρ .
-

113

¹For example, one possible implementation is as follows: create a minimal spanning tree according to the distances under ρ , then remove the heaviest $(k - 1)$ edges, resulting in k connected components, and output these components as clusters.

1.3 Comparison with Existing Analysis for Vanilla Spectral Algorithms in SBM

To the best of our knowledge, there are very few works on the analysis of vanilla spectral algorithms [AFWZ20, EBW18, PPV⁺19]. All of them only apply to the case of $k = O(1)$. In this work, we obtain the first analysis for general parameters n, k, p, q , in the symmetric SBM setting.

Davis-Kahan approaches. To study spectral algorithms in signal-plus-noise models, a key step is to understand how random noise perturbs the eigenvectors of a matrix. A commonly-used technical ingredient is the Davis-Kahan $\sin \Theta$ theorem (or its variant). However, this type of approach faces two challenges in SBM.

- Davis-Kahan leads to *worst-case perturbation* bounds. For perturbations caused by random noises, such as signal-plus-noise models, Davis-Kahan $\sin \Theta$ theorem is sometimes *suboptimal*.
- These $\sin \Theta$ theorems only lead to bound on 2-norm. However, in SBM analysis, we may need $(2 \rightarrow \infty)$ -norm bounds. See [CTP19] for more discussions.

Previous works such as [AFWZ20, EBW18, PPV⁺19] mainly followed this approach. They proposed some novel ideas to (partially) address these two challenges, but only apply to the case of $k = O(1)$. In contrast, our approach, following the power-iteration-based analysis proposed by [MZ22b], completely avoids Davis-Kahan $\sin \Theta$ theorem and can handle the case of $k = \omega(1)$.

Comparison with [MZ22b]. Inspired by power iteration methods, Mukherjee and Zhang [MZ22b] proposed a new approach to analyze the perturbation of random matrices. The idea is to approximate the eigenvectors of a matrix by its power. In fact, this method has been widely used in practice as a fast algorithm to approximate eigenvectors. However, there are two limitations of [MZ22b].

- Their analysis requires a nice structure of the mean matrix, i.e., all large eigenvalues are more or less the same.
- Their algorithm is not vanilla as it has a ‘centering step’. Moreover, their algorithm requires the knowledge of parameters p, q, k , and particularly, the centering step alone requires the knowledge of q . In comparison, we only need to know k ; further, we can also guess k (by checking the number of large eigenvalues) and then make alg. 1 fully parameter-free.

To overcome these limitations, we introduce a novel ‘polynomial approximation + entrywise analysis’ method, which makes this analysis more robust and requires less structure. More details will be discussed in Section 1.4.

1.4 Proof Outline and Technical Contributions

Let s_u denote the size of the cluster to which u belongs. Assume for now that all V_i ’s are of size roughly n/k . Indeed, this happens with high probability inasmuch as the partition is uniformly sampled.

Our goal is to show that there exists some threshold $\Delta > 0$ such that for every $u, v \in V$: if $u, v \in V_\ell$ for some ℓ , then $\|P_{\hat{G}_k} \hat{G}_u - P_{\hat{G}_k} \hat{G}_v\| \leq \Delta/4$; otherwise, $\|P_{\hat{G}_k} \hat{G}_u - P_{\hat{G}_k} \hat{G}_v\| \geq \Delta$.

Write $\varepsilon(u) \stackrel{\text{def}}{=} \|P_{\hat{G}_k} \hat{G}_u - G_u\|$. Then $\left| \|P_{\hat{G}_k} \hat{G}_u - P_{\hat{G}_k} \hat{G}_v\| - \|G_v - G_u\| \right| \leq \varepsilon(u) + \varepsilon(v)$. Note that $\|G_v - G_u\| = 0$ if $u, v \in V_\ell$ for some ℓ , otherwise, $\|G_v - G_u\| = (p - q) \cdot \sqrt{s_u + s_v} > (p - q) \sqrt{n/k}$. Therefore, setting $\Delta = 0.8(p - q) \sqrt{n/k}$, it suffices to show that $\varepsilon(u) \leq 0.1(p - q) \sqrt{n/k}$ for every $u \in V$.

We decompose $\varepsilon(u)$ into two terms:

$$\varepsilon(u) \leq \left\| P_{\hat{G}_k} (\hat{G}_u - G_u) \right\| + \left\| (P_{\hat{G}_k} - I) G_u \right\| = \underbrace{\left\| P_{\hat{G}_k} E_u \right\|}_{\text{"noise term"}} + \underbrace{\left\| (P_{\hat{G}_k} - I) G_u \right\|}_{\text{"deviation term"}}. \quad (1)$$

We shall bound the two terms from above separately. Intuitively, the noise term is small means $P_{\hat{G}_k}$ reduces the noise, while the deviation term is small means $P_{\hat{G}_k}$ preserves the data.

157 **Upper bound of the noise term** It is known that $P_{\widehat{G}_k}$ (resp., P_{G_k}) can be write as a polynomial of
 158 \widehat{G} (resp., G). By Weyl’s inequality, the eigenvalues of \widehat{G} are not too far from those of G . Therefore, in
 159 our case, one can find a simple polynomial φ which only depends on G , such that $\varphi(\widehat{G})$ (resp., $\varphi(G)$)
 160 is a good approximation of $P_{\widehat{G}_k}$ (resp., P_{G_k}); this is formalized in Lemma 3.2. Then we have the
 161 following decomposition: $\|P_{\widehat{G}_k} E_u\| \leq 2 \|\varphi(\widehat{G}) E_u\| \leq 2 \|\varphi(G) E_u\| + 2 \left\| \left(\varphi(\widehat{G}) - \varphi(G) \right) E_u \right\|$,
 162 where the first inequality follows from Lemma 3.2, which roughly says $\varphi(\widehat{G})$ is a good approximation
 163 of $P_{\widehat{G}_k}$.

- 164 1. The first term, $\|\varphi(G) E_u\|$, is small with high probability. To see this, we use Lemma 3.2
 165 again: $\|\varphi(G) E_u\| \leq \frac{3}{2} \|P_{G_k} E_u\|$. According to a known result (c.f. Proposition 2.4),
 166 $\|P_{G_k} E_u\|$ is small with high probability, largely because the projection P_{G_k} and the vector
 167 E_u are independent.
- 168 2. The second term is the tricky part, and we draw on an entrywise analysis. Namely, we study
 169 every entry of $(\varphi(\widehat{G}) - \varphi(G)) E_u$, using the new inequality from [MZ22b]. See Lemma 3.3
 170 for more details.

171 The upper bound for the noise term is encapsulated in Lemma 3.4.

172 **Upper bound of the deviation term** The following argument is reminiscent of [Vu18]. Say $u \in V_\ell$.
 173 Note that $G\chi_\ell = \sqrt{s_u} \cdot G_u$ where $\chi_\ell = \frac{1}{\sqrt{s_u}} \cdot 1_{V_\ell}$ is the normalized characteristic vector of V_ℓ (i.e.,
 174 $1_{V_\ell}(v) = 1 \iff v \in V_\ell$). It follows that

$$\left\| (P_{\widehat{G}_k} - I)G \right\|_2 \leq \left\| (P_{\widehat{G}_k} - I)\widehat{G} \right\|_2 + \left\| (P_{\widehat{G}_k} - I)E \right\|_2 \leq \left\| G - \widehat{G} \right\|_2 + \left\| (P_{\widehat{G}_k} - I)E \right\|_2 \leq 2 \|E\|_2,$$

175 where the second inequality holds because $P_{\widehat{G}_k} \widehat{G}$ is the best k -rank approximation of \widehat{G} and
 176 $\text{rank}(G) = k$, and in the third inequality, we use $\left\| (P_{\widehat{G}_k} - I) \right\|_2 \leq 1$, as $P_{\widehat{G}_k}$ is a projection
 177 matrix. Therefore,

$$\left\| (P_{\widehat{G}_k} - I)G_u \right\| = \frac{1}{\sqrt{s_u}} \left\| (P_{\widehat{G}_k} - I)G\chi_u \right\| \leq \frac{1}{\sqrt{s_u}} \left\| (P_{\widehat{G}_k} - I)G \right\|_2 \leq \frac{2 \|E\|_2}{\sqrt{s_u}}. \quad (2)$$

178 A typical result in random matrix theory (c.f. Proposition 2.3) states that with high probability,
 179 $\|E\|_2 = O(\sqrt{n})$. Combining Equation (2) and $s_u \approx n/k$, we get $\left\| (P_{\widehat{G}_k} - I)G_u \right\| = O(\sqrt{k})$. And
 180 by our assumption on n , we have $\sqrt{k} = o((p - q)n/k) = o(\Delta)$.

181 **Technical contribution.** The major novelty of our analysis is using the polynomial φ . [MZ22b]
 182 used a centering step to make the mean matrix nicely structured, while in our analysis, we used
 183 polynomial approximation to address this issue. Another difference is that in [MZ22b], the centering
 184 step appears explicitly in the algorithm. By contrast, our polynomial approximation only appears in
 185 the analysis — the algorithm is vanilla.

186 As a byproduct, we developed new techniques for studying *eigenspace perturbation*, a typical topic
 187 in random matrix theory. Our high-level idea is “polynomial approximation + entrywise analysis”.
 188 That is, we reduce the analysis of eigenspace perturbation to the analysis of a simple polynomial (of
 189 matrix) under perturbation. We have more tools to deal with the latter.

190 1.5 Discussion and Future Directions

191 In this paper, we studied the behavior of vanilla-SVD in SSBM, a benchmark signal-plus-noise model
 192 widely studied in random matrix theory. We showed that vanilla spectral algorithms indeed filter
 193 noise in SSBM. In fact, our analysis technique, ‘polynomial approximation + entrywise analysis’, is
 194 not very limited to SSBM. We believe our analysis may provide more applications for some more
 195 realistic models such as the factor model — a model which has been widely used in economics and
 196 model portfolio theory.

197 In the long term, it would be very interesting to understand the behavior of vanilla spectral algorithms
 198 on real data: 1) Why does it succeed in some applications? 2) How could we fix it if it has failed in

199 other cases? A deeper understanding of vanilla spectral algorithms will provide guidelines for using
 200 them in many machine learning tasks.

201 2 Preliminaries

202 **Notations** Let $\mathbf{1}_n$ denote the n -dimensional vector whose entries are all 1's, and let J_n be the $n \times n$
 203 matrix whose entries are all 1's. Let s_u denote the size of the cluster to which u belongs. For a matrix
 204 A , $A[i]$ denotes the row of A indexed by i , and A_i denotes the column indexed by i ; $\lambda_i(A)$ is the i -th
 205 largest eigenvalue of A ; let P_{A_k} denote the orthogonal projection matrix onto the space spanned by
 206 the first k eigenvectors of A . For a vector $x \in \mathbb{R}^n$, $\|x\| \stackrel{\text{def}}{=} \sqrt{x_1^2 + \dots + x_n^2}$ denotes the Euclidean
 207 norm.

208 **Definition 2.1** (Matrix operator norms). Let $A \in \mathbb{R}^{n \times n}$. Define $\|A\|_2 \stackrel{\text{def}}{=} \max_{\|x\|=1} \|Ax\|$ and
 209 $\|A\|_{2 \rightarrow \infty} \stackrel{\text{def}}{=} \max_{x: \|x\|=1} \|Ax\|_\infty$.

210 **Proposition 2.1** (e.g., [CTP19]). For all matrices $A, B \in \mathbb{R}^{n \times n}$, it holds that (1) $\|A\|_{2 \rightarrow \infty} =$
 211 $\max_{i \in [n]} \|A[i]\|$; (2) $\|AB\|_{2 \rightarrow \infty} \leq \|A\|_{2 \rightarrow \infty} \|B\|_2$.

212 **Proposition 2.2** (Weyl's inequality). For all $A, E \in \mathbb{R}^{n \times n}$, we have $|\lambda_i(A) - \lambda_i(A + E)| \leq \|E\|_2$.

213 **Proposition 2.3** (Norm of a random matrix [Vu18]). There is a constant $C_0 > 0$. Let E be
 214 a symmetric matrix whose upper diagonal entries e_{ij} are independent random variables where
 215 $e_{ij} = 1 - p_{ij}$ or $-p_{ij}$ with probabilities p_{ij} and $1 - p_{ij}$ respectively, where $p_{ij} \in [0, 1]$. Let
 216 $\sigma^2 := \max_{i,j} \{p_{ij}(1 - p_{ij})\}$. If $\sigma^2 \geq C_0 \log n/n$, then $\Pr[\|E\|_2 \geq C_0 \sigma n^{1/2}] \leq n^{-3}$.

217 **Proposition 2.4** (Projection of a random vector, lemma 2.1 in [Vu18]). There exists a constant C_1
 218 such that the following holds. Let $X = (\xi_1, \dots, \xi_n)$ be a random vector in \mathbb{R}^n whose coordinates ξ_i
 219 are independent random variables with mean 0 and variance at most $\sigma^2 \leq 1$. Assume furthermore
 220 that the ξ_i are bounded by 1 in absolute value. Let H be a subspace of dimension d and let $\Pi_H \xi$ be
 221 the length of the orthogonal projection of ξ onto H . Then $\Pr[\Pi_H X \geq \sigma \sqrt{d} + C_1 \sqrt{\log n}] \leq n^{-3}$.

222 **Proposition 2.5.** For $a \in [0, 2]$ and $r \in \mathbb{N}$, if $|a - 1| \leq \delta < \frac{1}{2^r}$, then $|a^r - 1| \leq 2r\delta$.

223 3 Analysis of Vanilla SVD Algorithm

224 Write $s_i \stackrel{\text{def}}{=} |V_i|$. We say the partition V_1, \dots, V_k is *balanced* if $\left(1 - \frac{1}{16 \log n}\right) \frac{n}{k} \leq s_i \leq$
 225 $\left(1 + \frac{1}{16 \log n}\right) \frac{n}{k}, \forall i \in [k]$. By Chernoff bound, the partition V_1, \dots, V_k is balanced with proba-
 226 bility at least $1 - n^{-1}$; hence, we assume that the partition is balanced in the following argument.
 227 Since $\sigma^2 \geq C \log n/n$, the event $\|E\| = O(\sqrt{n})$ holds with high probability (see Proposition 2.3).

228 Recall the decomposition into deviation term and noise term in Equation (1). We first state our
 229 upper bound of the deviation term, which readily follows from the argument in Section 1.4, and the
 230 complete proof is in Appendix B.

231 **Lemma 3.1** (Upper bound of deviation term). Let C_0 be the constant in Proposition 2.3. If the
 232 partition is balanced and $n \geq 10^4 \cdot C_0^2 \frac{k^2 \sigma^2}{(p-q)^2}$, then with probability at least $1 - n^{-3}$ we have
 233 $\left\| (P_{\hat{G}_k} - I)G_u \right\| \leq 0.04(p-q)\sqrt{n/k}, \forall u \in V$.

234 Section 3.1 and Section 3.2 lead to an upper bound of the noise term, and Section 3.3 is the proof of
 235 main theorem.

236 3.1 An Approximation of P_{G_k} and $P_{\hat{G}_k}$

237 In order to give some intuition on the choice of φ , we first analyze the spectrum of G , and the result
 238 is summed up in Theorem 3.1.

239 **The eigenvalues of G** Note that $G = H + q\mathbf{1}_n\mathbf{1}_n^\top$, where $H = \begin{pmatrix} (p-q)J_{s_1} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (p-q)J_{s_2} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & (p-q)J_{s_k} \end{pmatrix}$.

240 Without loss of generality, assume that $s_1 \geq s_2 \geq \dots \geq s_k$. It is easy to see that the eigenvalues of
 241 H are $(p-q)s_1, \dots, (p-q)s_k, 0$. Viewing G as a rank-one perturbation of H , we have the following
 242 theorem that characterizes eigenvalues of G . Its proof, in Appendix C, readily follows from a theorem
 243 in [BNS79], which studies eigenvalues under rank-one perturbation.

244 **Theorem 3.1.** Write $s_i \stackrel{\text{def}}{=} |V_i|$ and assume that $s_1 \geq s_2 \geq \dots \geq s_k$. Define $\delta_i \stackrel{\text{def}}{=} \lambda_i(G) - (p-q)s_i$,
 245 then (1) $\delta_i \geq 0$ and $\sum_{i=1}^k \delta_i = nq$; (2) $\lambda_1(G) \geq nq + (p-q)\frac{n}{k}$, and hence $\sum_{i=2}^k \delta_i \leq (p-q)(s_1 - \frac{n}{k})$.

246 **The choice of the polynomial φ** Let $\mu \stackrel{\text{def}}{=} (p-q)\frac{n}{k}$, and let $\psi(t)$ be the quadratic polynomial such
 247 that $\psi(\lambda_1(G)) = \psi(\mu) = 1, \psi(0) = 0$, i.e., $\psi(t) \stackrel{\text{def}}{=} -\frac{1}{\lambda_1(G)\mu}(t - \lambda_1(G))(t - \mu) + 1 \stackrel{\text{def}}{=} At^2 + Bt$,
 248 where $A = -\frac{1}{\lambda_1(G)\mu}, B = \frac{1}{\lambda_1(G)} + \frac{1}{\mu}$. Finally, let $\varphi(t) \stackrel{\text{def}}{=} (\psi(t))^r$ where $r \stackrel{\text{def}}{=} \log n$.

249 Here we give some intuition for the choice of φ . Let $\hat{G} = \sum_{i=1}^n \hat{\lambda}_i v_i v_i^\top$ be the spectral decomposition
 250 of \hat{G} . Then $\varphi(\hat{G}) = \sum_{i=1}^n \varphi(\hat{\lambda}_i) v_i v_i^\top, P_{\hat{G}_k} = \sum_{i=1}^k v_i v_i^\top$. The spectral decomposition of $\varphi(\hat{G}) -$
 251 $P_{\hat{G}_k}$ is $\varphi(\hat{G}) - P_{\hat{G}_k} = \sum_{i=1}^k (\varphi(\hat{\lambda}_i) - 1) v_i v_i^\top + \sum_{i=k+1}^n \varphi(\hat{\lambda}_i) v_i v_i^\top$. Hence,

$$\left\| \varphi(\hat{G}) - P_{\hat{G}_k} \right\|_2 = \max\{|\varphi(\hat{\lambda}_1) - 1|, \dots, |\varphi(\hat{\lambda}_k) - 1|, |\varphi(\hat{\lambda}_{k+1})|, \dots, |\varphi(\hat{\lambda}_n)|\}. \quad (3)$$

252 Recall that $\hat{\lambda}_i - \lambda_i(G)$ is bounded by Weyl's inequality. Plus, when the partition is balanced,
 253 Theorem 3.1 shows that the eigenvalues of G is nicely distributed: except for $\lambda_1(G)$, other eigenvalues
 254 are all close to μ . Hence, our choice of φ makes $\left\| \varphi(\hat{G}) - P_{\hat{G}_k} \right\|_2$ small, and thus $\varphi(\hat{G})$ is a good
 255 approximation of $P_{\hat{G}_k}$. Formally, we have the following lemma.

256 **Lemma 3.2** (Polynomial approximation). Assume that the partition is balanced and $n \geq 10^4 \cdot C_0^2 \cdot$
 257 $\frac{k^2 \cdot p \cdot \log n}{(p-q)^2}$, where C_0 is the constant in Proposition 2.3. Then with probability at least $1 - n^{-3}$, it
 258 holds that for all $x \in \mathbb{R}^n$, $\frac{1}{2} \left\| P_{\hat{G}_k} x \right\| \leq \left\| \varphi(\hat{G}) x \right\| \leq \frac{3}{2} \left\| P_{\hat{G}_k} x \right\| + \|x\| / n^{\log \log n}$, and $\frac{1}{2} \left\| P_{G_k} x \right\| \leq$
 259 $\left\| \varphi(G) x \right\| \leq \frac{3}{2} \left\| P_{G_k} x \right\|$.

260 *Proof.* Let $G = \sum_{i=1}^k \lambda_i u_i u_i^\top$ (resp., $\hat{G} = \sum_{i=1}^n \hat{\lambda}_i v_i v_i^\top$) be the spectral decomposition of G (resp.,
 261 \hat{G}). We shall use the following claim.

262 **Claim 3.1.** The following holds with probability $1 - n^{-3}$ (over the choice of E): for every $i \in [k]$,
 263 $\left| \varphi(\hat{\lambda}_i) - 1 \right| < \frac{1}{2}, |\varphi(\lambda_i) - 1| < \frac{1}{2}$; and for every $i = k+1, \dots, n$, $\left| \varphi(\hat{\lambda}_i) \right| < n^{-\log \log n}$.

Fix $x \in \mathbb{R}^n$. On the one hand, $\frac{1}{2} \leq \varphi(\hat{\lambda}_i) \leq \frac{3}{2}, \forall i \in [k]$, and hence

$$\left\| \varphi(\hat{G}) x \right\|^2 = \sum_{i=1}^n \varphi(\hat{\lambda}_i)^2 \langle x, v_i \rangle^2 \geq \sum_{i=1}^k \varphi(\hat{\lambda}_i)^2 \langle x, v_i \rangle^2 \geq \sum_{i=1}^k \frac{1}{4} \langle x, v_i \rangle^2 = \frac{1}{4} \left\| P_{\hat{G}_k} x \right\|^2,$$

which means $\left\| \varphi(\hat{G}) x \right\| \geq \frac{1}{2} \left\| P_{\hat{G}_k} x \right\|$. On the other hand,

$$\left\| \varphi(\hat{G}) x \right\|^2 = \sum_{i=1}^n \varphi(\hat{\lambda}_i)^2 \langle x, v_i \rangle^2 \leq \sum_{i=1}^k \left(\frac{3}{2}\right)^2 \langle x, v_i \rangle^2 + \sum_{i=k+1}^n \frac{\langle x, v_i \rangle^2}{n^{2 \log \log n}} \leq \frac{9}{4} \left\| P_{\hat{G}_k} x \right\|^2 + \frac{\|x\|^2}{n^{2 \log \log n}}.$$

264 Since $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$, we have $\left\| \varphi(\hat{G}) x \right\| \leq \frac{3}{2} \left\| P_{\hat{G}_k} x \right\| + \frac{\|x\|}{n^{\log \log n}}$. This establishes the first
 265 part.

266 Note that $\left\| \varphi(G) x \right\| = \sqrt{\sum_{i=1}^k \varphi(\lambda_i)^2 \langle x, u_i \rangle^2}$ and we also have $\frac{1}{2} \leq \varphi(\lambda_i) \leq \frac{3}{2}, \forall i \in [k]$, and thus
 267 similar argument goes for G . This finishes the proof of Lemma 3.2.

268 It remains to prove Claim 3.1. The claim readily follows from the choice of φ and the fact that $\lambda_i, \widehat{\lambda}_i$
 269 are close. A complete proof can be found in Appendix C. \square

270 3.2 The Upper Bound of the Noise Term

271 According to Equation (1), in order to derive an upper bound of $\|P_{G_k} E_u\|$, it remains to bound
 272 $\left\| \left(\varphi(\widehat{G}) - \varphi(G) \right) E_u \right\|$ from above. This is done by the following lemma.

Lemma 3.3. *Let C_0 be the constant in Proposition 2.3. Assume that the partition is balanced and
 $n \geq (100 + C_0)^2 \cdot \frac{k^2 \cdot p \cdot \log^{12} n}{(p-q)^2}$. For every $u \in V$, it holds that*

$$\Pr_E \left[\left\| \left(\varphi(\widehat{G}) - \varphi(G) \right) E_u \right\| \leq C_2 (\sqrt{kp} \log^2 n) + \frac{1}{\log n} \right] \geq 1 - O(n^{-2}),$$

273 where $C_2 \stackrel{\text{def}}{=} 7 \cdot 10^6$ is a constant.

274 Combining Lemma 3.2 and Proposition 2.4, we get an upper bound of the noise term:

275 **Lemma 3.4** (Upper bound of noise term). *Let C_0 be the constant in Proposition 2.3. Assume that
 276 $n \geq (100 + C_0)^2 \cdot \frac{k^2 \cdot p \cdot \log^{12} n}{(p-q)^2}$. Then with probability at least $1 - O(n^{-1})$, we have $\left\| P_{\widehat{G}_k} E_u \right\| \leq$
 277 $C_3 (\sqrt{kp} \log^2 n + \sqrt{\log n})$ for all $u \in V$, where C_3 is a constant.*

278 The proof of Lemma 3.3 is deferred to Section 4. We use it to prove Lemma 3.4 here.

279 *Proof of Lemma 3.4.* It follows from Lemma 3.2 that

$$\begin{aligned} \left\| P_{\widehat{G}_k} E_u \right\| &\leq 2 \left\| \varphi(\widehat{G}) E_u \right\| \leq 2 \left\| \left(\varphi(\widehat{G}) - \varphi(G) \right) E_u \right\| + 2 \left\| \varphi(G) E_u \right\| \\ &\leq 2 \left\| \left(\varphi(\widehat{G}) - \varphi(G) \right) E_u \right\| + 3 \left\| P_{G_k} E_u \right\|. \end{aligned}$$

280 By Proposition 2.4, with high probability at least $1 - n^{-1}$, $\|P_{G_k} E_u\|$ is bounded by $\sigma\sqrt{k} + C_1\sqrt{\log n}$,
 281 where C_1 is a universal constant. Meanwhile, by Lemma 3.3 and union bound over all u , with
 282 probability at least $1 - O(n^{-1})$, $\left\| \left(\varphi(\widehat{G}) - \varphi(G) \right) E_u \right\| \leq 7 \cdot 10^6 (\sqrt{kp} \log^2 n + 1/\log n)$ for every
 283 $u \in V$. Therefore, with probability $1 - O(n^{-1})$, it holds that $\left\| P_{\widehat{G}_k} E_u \right\| \leq 1.4 \times 10^7 \sqrt{kp} \log^2 n +$
 284 $3\sigma\sqrt{k} + 3C_1\sqrt{\log n}$ for all $u \in V$. Setting $C_3 \stackrel{\text{def}}{=} (1.4 \times 10^7 + 3 + 3C_1)$, we have the desired
 285 result. \square

286 3.3 Putting It Together

287 Now we are well-equipped to prove Theorem 1.1.

288 *Proof of Theorem 1.1.* Let $C \stackrel{\text{def}}{=} (100 + 100C_0 + 100C_3)^2$, where C_0, C_3 are the constants in Propo-
 289 sition 2.3 and Lemma 3.4. By our assumption on n , we have $(p - q)\sqrt{n/k} > 100C_3(\sqrt{kp} \log^6 n +$
 290 $\sqrt{\log n})$. It is easy to verify n satisfies the conditions in Lemma 3.4 and Lemma 3.1.

291 Write $\Delta \stackrel{\text{def}}{=} 0.8(p - q)\sqrt{n/k}$. We aim to show that for every $u, v \in V$: if $u, v \in V_\ell$ for some
 292 ℓ , then $\left\| P_{\widehat{G}_k} \widehat{G}_u - P_{\widehat{G}_k} \widehat{G}_v \right\| \leq \Delta/4$; otherwise, $\left\| P_{\widehat{G}_k} \widehat{G}_u - P_{\widehat{G}_k} \widehat{G}_v \right\| \geq \Delta$. Then by calling
 293 `ClusterByDistance`, alg. 1 recovers all large clusters correctly.

294 Let $\varepsilon(u) \stackrel{\text{def}}{=} \left\| P_{\widehat{G}_k} \widehat{G}_u - G_u \right\|$. According to the argument in Section 1.4, it suffices to show that
 295 $\varepsilon(u) \leq 0.1(p - q)\sqrt{n/k}$ for all $u \in V$. We further decompose $\varepsilon(u)$ into noise term and deviation
 296 term, i.e., $\varepsilon(u) \leq \text{noise}(u) + \text{dev}(u)$, where $\text{noise}(u) \stackrel{\text{def}}{=} \left\| P_{\widehat{G}_k} E_u \right\|$ and $\text{dev}(u) \stackrel{\text{def}}{=} \left\| (P_{\widehat{G}} - I)G_u \right\|$.
 297 By Lemma 3.4 and Lemma 3.1, with probability at least $1 - O(n^{-1})$, the following hold for all $u \in V$:
 298 (1) $\text{noise}(u) \leq C_3(\sqrt{kp} \log^2 n + \sqrt{\log n}) \leq 0.01(p - q)\sqrt{n/k}$; (2) $\text{dev}(u) \leq 0.04(p - q)\sqrt{n/k}$.
 299 Therefore, with probability at least $1 - O(n^{-1})$, we indeed have $\varepsilon(u) \leq 0.1(p - q)\sqrt{n/k}, \forall u \in V$.
 300 This completes the proof. \square

301 4 Proof of Lemma 3.3: Entrywise Analysis

302 This section is dedicated to proving Lemma 3.3.

303 Since both $(\varphi(\widehat{G}) - \varphi(G))$ and E are symmetric, we have $\left\|(\varphi(\widehat{G}) - \varphi(G))E_u\right\| \leq$
 304 $\left\|E(\varphi(\widehat{G}) - \varphi(G))\right\|_{2 \rightarrow \infty}$. The high-level idea is to write $E(\varphi(\widehat{G}) - \varphi(G))$ as a sum of ma-
 305 trices, where each matrix is of the form $E^t S Q$ such that $\|Q\|_2 = O(1)$. This way, we have
 306 $\|E^t S Q\|_{2 \rightarrow \infty} \leq \|E^t S\|_{2 \rightarrow \infty} \cdot O(1)$, and $\|E^t S\|_{2 \rightarrow \infty}$ is bounded by a lemma from [MZ22b].

307 Let $D \stackrel{\text{def}}{=} \psi(\widehat{G}) - \psi(G) = A(EG + GE + E^2) + BE$ and write $F \stackrel{\text{def}}{=} \psi(G)$, $\widehat{F} \stackrel{\text{def}}{=} \psi(\widehat{G})$. Then

$$\begin{aligned} \varphi(\widehat{G}) - \varphi(G) &= \psi(\widehat{G})^r - \psi(G)^r = (F + D)^r - F^r \\ &= \underbrace{F^{r-1}D + F^{r-2}D\widehat{F} + \dots + FD\widehat{F}^{r-2}}_{\stackrel{\text{def}}{=} M} + D\widehat{F}^{r-1}, \end{aligned}$$

308 where the last step is a decomposition based on the first location of D in the product terms. And

$$D\widehat{F}^{r-1} = D(D + F)^{r-1} = D^r + \underbrace{DF\widehat{F}^{r-2} + D^2F\widehat{F}^{r-3} + \dots + D^{r-1}F}_{\stackrel{\text{def}}{=} M'}.$$

309 That is, $E(\varphi(\widehat{G}) - \varphi(G)) = EM + ED^r + EM'$. We bound the three terms respectively.

310 Here we first list some definitions and estimations of the quantities involved.

- 311 • According to Proposition 2.3, with probability at least $1 - n^{-3}$, we have $\|E\|_2 \leq C_0 \sigma \sqrt{n}$,
 312 where C_0 is a constant. In the following argument, we always assume this holds.
- 313 • $\mu \stackrel{\text{def}}{=} (p - q)n/k$. By our assumption on n , we have $\mu \geq (100 + C_0)\sqrt{np} \log^6 n$.
- 314 • $A = -\frac{1}{\lambda_1(G)\mu}$, $B = (\frac{1}{\lambda_1(G)} + \frac{1}{\mu})$, $r = \log n$; $\lambda_1(G) > \mu$, and thus $B \leq \frac{2}{\mu}$, $|A| \leq \frac{1}{\mu^2}$.
- 315 • By Claim 3.1, $\|F\|_2 \leq 1 + \frac{1}{4 \log n}$, $\|\widehat{F}\|_2 \leq 1 + \frac{1}{4 \log n}$. By Proposition 2.5, $\|F\|_2^t, \|\widehat{F}\|_2^t \leq$
 316 $2, \forall t \leq \log n$.

317 **Upper bound of $\|EM\|_{2 \rightarrow \infty}$** Note that $\|EF^t D\|_{2 \rightarrow \infty} \leq \|EF\|_{2 \rightarrow \infty} \|F\|_2^{t-1} \|D\|_2$, and
 318 $\|F\|_2^{t-1} \leq 2$ for all $t \leq r$. Moreover, $\|D\|_2 \leq |A|(2\|E\|_2 \|G\|_2 + \|E\|_2^2) + B\|E\|_2 \leq$
 319 $3 \frac{\|E\|_2}{\mu} + \frac{\|E\|_2^2}{\mu^2} + \leq 4 \frac{\|E\|_2}{\mu} \leq 4(\log^6 n)^{-1} < \frac{1}{\log^3 n}$. And the following lemma gives an upper
 320 bound of $\|EF\|_{2 \rightarrow \infty}$.

321 **Lemma 4.1.** $\Pr_E [\|EF\|_{2 \rightarrow \infty} \leq 10(\sqrt{kp \log n} + \sqrt{\log n})] \geq 1 - 2n^{-2}$.

322 Therefore, by union bound, we have the following holds with probability at least $1 - n^{-1}$:

$$\|EM\|_{2 \rightarrow \infty} \leq r \cdot 10(\sqrt{kp \log n} + \sqrt{\log n}) \cdot 2 \cdot \frac{1}{\log^3 n} \leq \frac{40(\sqrt{kp} + 1)}{\log n}. \quad (4)$$

323 **Upper bound of $\|ED^r\|_{2 \rightarrow \infty}$** Since $\|D\|_2 < \frac{1}{\log^3 n}$, we have

$$\|ED^r\|_{2 \rightarrow \infty} \leq \|E\|_{2 \rightarrow \infty} \|D\|_2^r \leq \sqrt{n} \cdot (\log^3 n)^{-\log n} < \frac{1}{n}. \quad (5)$$

324 **Lemma 4.2** (Upper bound of $\|EM'\|_{2 \rightarrow \infty}$). *With probability $1 - O(n^{-2})$ (over the choice of E),*
 325 *we have $\|EM'\|_{2 \rightarrow \infty} \leq 6C_2 \sqrt{kp} \log^2 n$, where $C_2 = 10^6$ is a constant.*

Finally, combining Equation (4), Equation (5), and the above lemma, we conclude that with probability at least $1 - O(n^{-2})$,

$$\left\|E(\varphi(\widehat{G}) - \varphi(G))\right\|_{2 \rightarrow \infty} \leq \frac{40(\sqrt{kp} + 1)}{\log n} + \frac{1}{n} + 6C_2 \sqrt{kp} \log^2 n \leq 7C_2(\sqrt{kp} \log^2 n + \frac{1}{\log n}).$$

326 This establishes Lemma 3.3.

327 Proofs of Lemma 4.1 and Lemma 4.2 are deferred to Appendix D.

References

- 328
- 329 [AFWZ20] Emmanuel Abbe, Jianqing Fan, Kaizheng Wang, and Yiqiao Zhong. Entrywise eigenvector
330 analysis of random matrices with low expected rank. *Ann. Statist.*, 48(3):1452–1474,
331 2020.
- 332 [ARS⁺04] Paolo Antonelli, HE Revercomb, LA Sromovsky, WL Smith, RO Knuteson, DC Tobin,
333 RK Garcia, HB Howell, H-L Huang, and FA Best. A principal component noise filter
334 for high spectral resolution infrared measurements. *Journal of Geophysical Research:
335 Atmospheres*, 109(D23), 2004.
- 336 [BNS79] James R. Bunch, Christopher P. Nielsen, and Danny C. Sorensen. Rank-one modification
337 of the symmetric eigenproblem. *Numer. Math.*, 31(1):31–48, 1978/79.
- 338 [Col19] Sam Cole. Recovering nonuniform planted partitions via iterated projection. *Linear
339 Algebra and its Applications*, 576:79–107, 2019.
- 340 [CTP19] Joshua Cape, Minh Tang, and Carey E Priebe. The two-to-infinity norm and singular
341 subspace geometry with applications to high-dimensional statistics. *The Annals of
342 Statistics*, 47(5):2405–2439, 2019.
- 343 [EBW18] Justin Eldridge, Mikhail Belkin, and Yusu Wang. Unperturbed: spectral analysis beyond
344 Davis-Kahan. In *Algorithmic learning theory 2018*, volume 83 of *Proc. Mach. Learn.
345 Res. (PMLR)*, page 38. Proceedings of Machine Learning Research PMLR, [place of
346 publication not identified], 2018.
- 347 [GM05] Joachim Giesen and Dieter Mitsche. Reconstructing many partitions using spectral
348 techniques. In *International Symposium on Fundamentals of Computation Theory*, pages
349 433–444. Springer, 2005.
- 350 [HLL83] Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic block-
351 models: First steps. *Social networks*, 5(2):109–137, 1983.
- 352 [KAH19] Vladimir Yu Kiselev, Tallulah S Andrews, and Martin Hemberg. Challenges in unsuper-
353 vised clustering of single-cell rna-seq data. *Nature Reviews Genetics*, 20(5):273–282,
354 2019.
- 355 [LR15] Jing Lei and Alessandro Rinaldo. Consistency of spectral clustering in stochastic block
356 models. *The Annals of Statistics*, 43(1):215–237, 2015.
- 357 [LZK22] Boris Landa, Thomas TCK Zhang, and Yuval Kluger. Biwhitening reveals the rank of a
358 count matrix. *SIAM Journal on Mathematics of Data Science*, 4(4):1420–1446, 2022.
- 359 [McS01] Frank McSherry. Spectral partitioning of random graphs. In *42nd IEEE Symposium
360 on Foundations of Computer Science (Las Vegas, NV, 2001)*, pages 529–537. IEEE
361 Computer Soc., Los Alamitos, CA, 2001.
- 362 [MZ22a] Chandra Sekhar Mukherjee and Jiapeng Zhang. Compressibility: Power of pca in
363 clustering problems beyond dimensionality reduction. *arXiv preprint arXiv:2204.10888*,
364 2022.
- 365 [MZ22b] Chandra Sekhar Mukherjee and Jiapeng Zhang. Detecting hidden communities by power
366 iterations with connections to vanilla spectral algorithms, 2022.
- 367 [PPV⁺19] Carey E Priebe, Youngser Park, Joshua T Vogelstein, John M Conroy, Vince Lyzinski,
368 Minh Tang, Avanti Athreya, Joshua Cape, and Eric Bridgeford. On a two-truths phe-
369 nomenon in spectral graph clustering. *Proceedings of the National Academy of Sciences*,
370 116(13):5995–6000, 2019.
- 371 [SEK04] Michael Steinbach, Levent Ertöz, and Vipin Kumar. The challenges of clustering high
372 dimensional data. In *New directions in statistical physics*, pages 273–309. Springer,
373 2004.
- 374 [Vu18] Van Vu. A simple SVD algorithm for finding hidden partitions. *Combin. Probab.
375 Comput.*, 27(1):124–140, 2018.

376 [ZLWZ09] Lei Zhang, Rastislav Lukac, Xiaolin Wu, and David Zhang. Pca-based spatially adaptive
377 denoising of cfa images for single-sensor digital cameras. *IEEE transactions on image*
378 *processing*, 18(4):797–812, 2009.

379 **A Useful inequalities**

Proposition A.1 (Chernoff bound). *Let X_1, \dots, X_m be i.i.d random variables that can take values in $\{0, 1\}$, with $\mathbf{E}[X_i] \leq p$ for $1 \leq i \leq m$. Then it holds that*

$$\Pr \left[\left| \sum_{i=1}^n X_i - mp \right| \geq t \right] \leq \exp \left(-\frac{3t^2}{mp} \right).$$

Proposition A.2 (Hoeffding bound). *Let X_1, \dots, X_m be independent random variables such that $a_i \leq X_1 \leq b_i$, and write $S \stackrel{\text{def}}{=} \sum_{i=1}^m X_i$. Then it holds that*

$$\Pr [|S - \mathbf{E}[S]| > t] \leq 2 \exp \left(-\frac{2t^2}{\sum_{i=1}^m (b_i - a_i)^2} \right).$$

380 **Definition A.1.** Let X be a Bernoulli random variable with parameter p , i.e., $\Pr[X = 1] =$
 381 $p, \Pr[X = 0] = 1 - p$. The random variable $Y \stackrel{\text{def}}{=} X - p = X - \mathbf{E}[X]$ is called *centered*
 382 *Bernoulli random variable with parameter p .*

Proposition A.3 (Adapted from [MZ22b]). *Let $S \in \mathbb{R}^{n \times n}$, and let $E = (\xi_{ij})$ be an $n \times n$ symmetric random matrix, where*

$$\{\xi_{ij} : 1 \leq i \leq j \leq n\}$$

are independent, centered Bernoulli random variables with parameter at most α for all i, j . Suppose that every entry of S takes value in $[-\beta, \beta]$, and each column of S has at most γ non-zero entries. Then for every $t \in [\log n]$, it holds that

$$\Pr [|(E^t S)_{ij}| > (\log n)^{5t} C_t] = O(n^{-4}), \forall i, j \in [n],$$

where

$$C_t \stackrel{\text{def}}{=} 500\beta\sqrt{\alpha}\sqrt{\gamma} \cdot (100\sqrt{n\alpha})^{t-1}.$$

By union bound,

$$\Pr [\|E^t S\|_{2 \rightarrow \infty} > \sqrt{n}(\log n)^{5t} C_t] = O(n^{-2}).$$

383 **Remark A.1.** The parameter α is determined by E , which equals to p in our case. The above bound
 384 is particularly useful when β, γ are small, that is, we want the matrix S to have small entries and
 385 sparse columns.

386 **Proposition A.4** (Proposition 2.5 restated). *For $a \in [0, 2]$ and $r \in \mathbb{N}$, if $|a - 1| \leq \delta < \frac{1}{2r}$, then*
 387 $|a^r - 1| \leq 2r\delta$.

Proof. Let $x = a - 1 \in [-\delta, \delta]$. If $0 \leq a \leq 1$, we have $1 \geq a^r = (1 + x)^r \geq 1 + rx \geq 1 - r\delta$. If $1 < a < 1 + 1/r$, then $0 < x < 1/r$ and hence

$$1 \leq a^r = (1 + x)^r = \sum_{i=0}^r \binom{r}{i} x^i \leq \sum_{i=0}^r r^i x^i < \sum_{i=0}^{\infty} (rx)^i = \frac{1}{1 - rx} = 1 + \frac{rx}{1 - rx} \leq 1 + 2r\delta.$$

388 □

389 **B Bounding the Deviation Term**

Proof of Lemma 3.1. Our assumption on n implies that $(p - q)n/k > 100C_0\sigma\sqrt{n}$. By Proposition 2.3, with probability at least $1 - n^{-3}$, we have

$$\|E\|_2 \leq C_0\sigma\sqrt{n} \leq 0.01(p - q)n/k.$$

390 According to Equation (2) and $s_u \geq \frac{n}{2k}$, we have

$$\|(P_{\hat{G}} - I)G_u\| \leq \frac{2\|E\|_2}{\sqrt{s_u}} \leq 0.04(p - q)n/k.$$

391 □

392 **C Polynomial Approximation**

393 The proof of Theorem 3.1 rely on the following result on rank-one perturbation.

Proposition C.1 (Eigenvalues under rank-one perturbation, Theorem 1 in [BNS79]). *Let $C = D + \rho zz^T$, where D is diagonal, $\|z\|_2 = 1$. Let $d_1 \geq d_2 \geq \dots \geq d_n$ be the eigenvalues of D , and let $\tilde{d}_1 \geq \tilde{d}_2 \geq \dots \geq \tilde{d}_n$ be the eigenvalues of C . Then*

$$\tilde{d}_i = d_i + \rho \mu_i, \quad 1 \leq i \leq n,$$

394 where $\sum_{i=1}^n \mu_i = 1$ and $0 \leq \mu_i \leq 1$.

Proof of Theorem 3.1. Let $\chi_i \in \{0, 1\}^V$ be the indicator vector for V_i , i.e., $\chi_i(u) = 1$ iff $\phi(u) = i$. It is easy to see that the eigenvectors of H are $\frac{1}{\sqrt{s_1}}\chi_1, \dots, \frac{1}{\sqrt{s_k}}\chi_k$. Write $V = \left(\frac{1}{\sqrt{s_1}}\chi_1, \dots, \frac{1}{\sqrt{s_k}}\chi_k \right) \in \mathbb{R}^{X \times X}$, $D = \text{diag}((p-q)s_1, \dots, (p-q)s_k, 0, \dots, 0)$, then we have $H = VDV^T$. Note that $\mathbf{1}_n = V(\sqrt{s_1}, \dots, \sqrt{s_n})^T$, and hence

$$G = H + q\mathbf{1}_n\mathbf{1}_n^T = V(D + \rho zz^T)V^T,$$

395 where $\rho = nq$, $z = \frac{1}{\sqrt{n}}(\sqrt{s_1}, \dots, \sqrt{s_n})^T$. This means the eigenvalues of G are the same as those of
 396 $D + \rho zz^T$. Since $\|z\| = 1$, Item 1 follow directly from Proposition C.1. To see Item 2, we use the
 397 Rayleigh quotient characterization of the largest eigenvalue:

$$\begin{aligned} \lambda_1(G) &= \max_v \frac{v^T G v}{\|v\|^2} \geq \frac{\mathbf{1}_n^T G \mathbf{1}_n}{n} = \frac{\sum_{u,v \in X} G_{uv}}{n} = \frac{n^2 q + (p-q) \cdot (s_1^2 + \dots + s_k^2)}{n} \\ &\geq nq + (p-q) \frac{n}{k}. \end{aligned}$$

398 where the last inequality follows from $\frac{n}{k} = \frac{1}{k} \sum_{i=1}^k s_i \leq \sqrt{\sum_{i=1}^k s_i^2 / k}$, □

399 *Proof of Claim 3.1.* The assumption on n in Lemma 3.2 implies that $\mu = (p-q)n/k \geq 100C_0\sigma\sqrt{n} \cdot$
 400 $\log n$. By Weyl's inequality, we have

$$\left| \hat{\lambda}_i - \lambda_i \right| \leq \|E\|_2 \leq C_0\sigma\sqrt{n} \leq \frac{\mu}{100 \log n}, \forall i \in [n].$$

401 Meanwhile, by Theorem 3.1,

$$|\lambda_i - \mu| \leq |\lambda_i(G) - (p-q)s_i| + |(p-q)s_i - \mu| \leq \frac{(p-q)n/k}{16 \log n} + \frac{(p-q)n/k}{16 \log n} \leq \frac{\mu}{8 \log n}.$$

402 for $i = 2, 3, \dots, k$. Hence, write $\varepsilon \stackrel{\text{def}}{=} \frac{\mu}{6 \log n}$, we have

- 403 1. $|\hat{\lambda}_1 - \lambda_1| \leq \varepsilon$;
- 404 2. $\lambda_2, \dots, \lambda_k, \hat{\lambda}_2, \dots, \hat{\lambda}_k \in [\mu - \varepsilon, \mu + \varepsilon]$;
- 405 3. for every $i \geq k+1$, $|\hat{\lambda}_i| \leq \varepsilon$.

406 First, $\psi(\lambda_1) = 1$ according to the definition of ψ , and hence $\varphi(\lambda_1) = 1$. As for $\hat{\lambda}_1$,

$$\begin{aligned} \left| \psi(\hat{\lambda}_1) - 1 \right| &= \left| \psi(\hat{\lambda}_1) - \psi(\lambda_1) \right| \leq |A|\varepsilon^2 + |2A\lambda_1 + B|\varepsilon && \text{(by definition of } \psi) \\ &\leq \frac{\varepsilon^2}{\lambda_1 \mu} + \frac{\varepsilon}{\mu} && \text{(since } 2A\lambda_1 + B = \frac{1}{\lambda_1} - \frac{1}{\mu} \geq -\frac{1}{\mu}) \\ &\leq \frac{1}{36 \log^2 n} + \frac{1}{6 \log n} && \text{(as } \frac{\varepsilon}{\mu} = \frac{1}{6 \log n}) \\ &< \frac{1}{4 \log n}. \end{aligned}$$

407 Consequently, $|\varphi(\widehat{\lambda}_1) - 1| < \frac{2r}{4 \log n} \leq 1/2$ by Proposition 2.5.

408 Next, for $a \in \{\lambda_2, \dots, \lambda_k, \widehat{\lambda}_2, \dots, \widehat{\lambda}_k\}$, the argument is similar:

$$|\psi(a) - 1| = |\psi(a) - \psi(\mu)| \leq |A|\varepsilon^2 + |2A\mu + B|\varepsilon \leq \frac{\varepsilon^2}{\lambda_1\mu} + \frac{\varepsilon}{\mu} < \frac{1}{4 \log n},$$

409 where the second inequality follows from $2A\mu + B = \frac{1}{\mu} - \frac{1}{\lambda_1} \leq \frac{1}{\mu}$. This yields $|\varphi(a) - 1| \leq 1/2$
410 by Proposition 2.5.

Finally, for $i \geq k + 1$, it holds that

$$\left| \psi(\widehat{\lambda}_i) \right| \leq |A|\varepsilon^2 + B\varepsilon = \frac{\varepsilon^2}{\lambda_1\mu} + \frac{\varepsilon}{\mu} < \frac{1}{4 \log n},$$

411 which means $\left| \varphi(\widehat{\lambda}_i) \right|^r = \left| \psi(\widehat{\lambda}_i) \right|^r < \left(\frac{1}{4 \log n} \right)^{\log n} < n^{-\log \log n}$. \square

412 D Bounding the Noise Term

413 *Proof of Lemma 4.1.* The lemma readily follows from the following entrywise bound and Chernoff
414 bound.

415 **Claim D.1** (Entries of $\psi(G)$). *For every $u, v \in X$, if $u, v \in V_\ell$ for some ℓ , then $0 \leq F_{uv} \leq \frac{5k}{n}$;*
416 *otherwise, $|F_{uv}| \leq \frac{10}{n}$.*

417 We decompose $F = F' + F''$, where F' is the intra-cluster part, i.e., $F'_{uv} = F_{uv}$ if $u, v \in V_\ell$ for
418 some ℓ , and $F'_{uv} = 0$ otherwise. Since for every column of F'_v , its non-zero entries are identical
419 and at most $5k/n$ by the above claim. Hence, every entry of EF' equals to the sum of at most
420 $2n/k$ independent, centered Bernoulli variables with parameter p , scaled by some factor at most
421 $\frac{5k}{n}$. By Chernoff bound, $\Pr_E \left[|(EF')_{uv}| > 10\sqrt{kp \log n/n} \right] \leq n^{-4}, \forall u, v \in V$, and we have
422 $\Pr_E \left[\|EF'\|_{2 \rightarrow \infty} \leq 10\sqrt{kp \log n} \right] \geq 1 - n^{-2}$ by union bound. Analogously, by Hoeffding bound,
423 $\Pr_E \left[\|EF''\|_{2 \rightarrow \infty} \leq 10\sqrt{\log n} \right] \geq 1 - n^{-2}$. Since $\|EF\|_{2 \rightarrow \infty} \leq \|EF'\|_{2 \rightarrow \infty} + \|EF''\|_{2 \rightarrow \infty}$, the
424 lemma follows from the above two inequalities and union bound. \square

425 *Proof of Claim D.1.* Write $\lambda = \lambda_1(G)$ and recall that (i) $(p - q)s_u \leq 2\mu$ for all u (ii) $nq + \mu \leq \lambda \leq$
426 $nq + (p - q)s_1 < nq + 2\mu$, (iii) $\lambda > p \cdot \mu$, and for all $u \in V$. Assume that $u, v \in V_\ell$ for some ℓ .
427 Then

$$\begin{aligned} F_{uv} &= AG_u^\top G_v + BG_{uv} = -\frac{nq^2 + (p^2 - q^2)s_u}{\lambda\mu} + \left(\frac{1}{\lambda} + \frac{1}{\mu} \right) p \\ &= \frac{-nq^2 - (p^2 - q^2)s_u + (p - q)(\lambda + \mu) + q(\lambda + \mu)}{\lambda\mu} \\ &= \frac{q(\mu + \lambda - nq) + (p - q)(\lambda + \mu - (p + q)s_u)}{\lambda\mu}. \end{aligned}$$

Since $\lambda - nq \geq \mu$, the numerator is at least

$$2q\mu + (p - q)(\lambda + \mu - (p + q)s_u) = (p - q)(2qn/k + \lambda + \mu - (p + q)s_u).$$

Because $s_u \leq 2n/k, \lambda \geq nq + (p - q)n/k$, we have

$$2qn/k + \lambda + \mu - (p + q)s_u > 2qn/k + nq + (p - q)2n/k - (p + q)2n/k = (n - 2n/k)q \geq 0,$$

which means $F_{uv} \geq 0$. Meanwhile,

$$F_{uv} \leq \frac{q(\mu + \lambda - nq)}{\lambda\mu} + \frac{(p - q)(\lambda + \mu)}{\lambda\mu},$$

428 where the first term is at most $\frac{3q}{\lambda} \leq \frac{3}{n}$ by (ii); second term is at most $\frac{2(p - q)}{\mu} \leq \frac{2k}{n}$. Therefore,
429 $|F_{uv}| \leq \frac{5k}{n}$.

430 For the second part, assume that u, v are not in the same cluster. Then

$$F_{uv} = AG_u^\top G_v + BG_{uv} = -\frac{nq^2 + (pq - q^2)(s_u + s_v)}{\lambda\mu} + \left(\frac{1}{\lambda} + \frac{1}{\mu}\right)q.$$

Hence,

$$|F_{uv}| \leq \left| \frac{q(\lambda + \mu - nq)}{\lambda\mu} \right| + \left| \frac{q(p - q)(s_u + s_v)}{\lambda\mu} \right|.$$

431 By (ii), the first term is at most $\frac{3q}{\lambda} < \frac{3}{n}$; by (i), the second term is at most $\frac{4q}{\lambda} \leq \frac{4}{n}$; hence,
 432 $|F_{uv}| \leq \frac{10}{n}$. \square

433 **Upper bound of $\|EM'\|_{2 \rightarrow \infty}$ (Proof of Lemma 4.2)** Write $L \stackrel{\text{def}}{=} A(EG + GE)$, $R \stackrel{\text{def}}{=} AE^2 + BE$.
 434 Then $D^t F = (L + R)^t F = R^t F + R^{t-1} L F + R^{t-2} L D F + \dots + R L D^{t-2} F + L D^{t-1} F$. It suffices
 435 to derive a good upper bound of $\|E^\eta L\|_{2 \rightarrow \infty}$ and $\|E^\eta F\|_{2 \rightarrow \infty}$, as R^w can be further expressed as
 436 sum of powers of E . This is done by the following lemma:

437 **Lemma D.1.** *The following holds with probability $1 - O(n^{-2})$ over the choice of E : for all $\eta \leq \log n$,*
 438 *it holds that*

439 $\bullet \|E^\eta L\|_{2 \rightarrow \infty} \leq C_2 \sqrt{kp} (100\sqrt{np})^{\eta-1} \log^{5\eta} n,$

440 $\bullet \|E^\eta F\|_{2 \rightarrow \infty} \leq C_2 \sqrt{kp} (100\sqrt{np})^{\eta-1} \log^{5\eta} n,$

441 where $C_2 \stackrel{\text{def}}{=} 10^6$ is an absolute constant.

442 Specifically,

$$\begin{aligned} ED^t F &= ER^t F + \sum_{i=0}^{t-1} ER^i L D^{t-1-i} F \\ &= \underbrace{E \sum_{j=0}^t \binom{t}{j} A^j B^{t-j} E^{j+t} F}_{\stackrel{\text{def}}{=} M_t} + \underbrace{\sum_{i=0}^{t-1} E \sum_{j=0}^i \binom{i}{j} A^j B^{i-j} E^{i+j} L D^{t-1-i} F}_{\stackrel{\text{def}}{=} N_t}. \end{aligned}$$

443 Note that $|A^j B^{w-j}| \leq 2^w \cdot \mu^{-(w+j)} \leq 2^w \cdot (100\sqrt{np} \log^6 n)^{-(w+j)}$. It follows that for every
 444 $t \in [r-1]$,

$$\begin{aligned} \|M_t\|_{2 \rightarrow \infty} &\leq \sum_{j=0}^t \binom{t}{j} |A^j B^{t-j}| \|E^{t+j+1} F\|_{2 \rightarrow \infty} \\ &\leq \sum_{j=0}^t \binom{t}{j} 2^t \cdot (100\sqrt{np} \log^6 n)^{-(t+j)} \cdot C_2 \sqrt{kp} \cdot (100\sqrt{np})^{t+j} \log^{5(t+j)} n \\ &\leq C_2 \sqrt{kp} \cdot 2^t \sum_{j=0}^t \binom{t}{j} (\log n)^{-(t+j)} \\ &= C_2 \sqrt{kp} \cdot 2^t (\log n)^{-t} \cdot \left(1 + \frac{1}{\log n}\right)^t \leq C_2 \sqrt{kp}, \end{aligned}$$

445 where the second inequality is by Lemma D.1, and the last step follows from Proposition 2.5.
 446 Similarly, for every $t \in [r-1]$,

$$\begin{aligned}
 \|N_t\|_{2 \rightarrow \infty} &\leq \sum_{i=0}^{t-1} \sum_{j=0}^i \binom{i}{j} |A^j B^{i-j}| \|E^{i+j+1} L\|_{2 \rightarrow \infty} \|D^{t-1-i} F\|_2 \\
 &\leq 2 \sum_{i=0}^{t-1} \sum_{j=0}^i \binom{i}{j} |A^j B^{i-j}| \|E^{i+j+1} L\|_{2 \rightarrow \infty} \\
 &\leq 2 \sum_{i=0}^{t-1} \sum_{j=0}^i \binom{i}{j} 2^i \cdot (100\sqrt{np} \log^6 n)^{-(i+j)} \cdot C_2 \sqrt{kp} \cdot (100\sqrt{np})^{i+j} \log^5(i+j) n \\
 &\leq 2C_2 \sqrt{kp} \sum_{i=0}^{t-1} 2^i \sum_{j=0}^i \binom{t}{j} (\log n)^{-(i+j)} \\
 &\leq 2C_2 \sqrt{kp} \cdot t \leq 2C_2 \sqrt{kp} \cdot \log n,
 \end{aligned}$$

447 where the second inequality follows from $\|D^{t-1-i} F\|_2 \leq 2$, and the third inequality is by
 448 Lemma D.1. In sum,

$$\|EM'\|_{2 \rightarrow \infty} \leq \sum_{t=1}^{r-1} (\|M_t\|_{2 \rightarrow \infty} + \|N_t\|_{2 \rightarrow \infty}) \|\widehat{F}\|_2^{r-1-t} \leq 6C_2 \sqrt{kp} \log^2 n, \quad (6)$$

449 where in the last inequality we also use $\|\widehat{F}\|_2^t \leq 2$.

450 The proof of Lemma D.1 draw on the entrywise bound in Proposition A.3.

451 *Proof of Lemma D.1.* Fix an $\eta \leq \log n$. Write $s^* \stackrel{\text{def}}{=} n/k$ for the ease of notation. Observe that
 452 $E^\eta L = A(E^{\eta+1}G + E^\eta GE) = A(E^{\eta+1}H + E^\eta HE) + Aq(E^{\eta+1}J_n + E^\eta J_n E)$. We can apply
 453 Proposition A.3 to $E^{\eta+1}H$ and $E^\eta H$, with $\alpha = p, \beta = p - q, \gamma = 2s^*$. That is, with probability at
 454 least $1 - O(n^{-2})$,

$$\|E^j H\|_{2 \rightarrow \infty} \leq 500(\log n)^{5j} \sqrt{n}(p-q) \sqrt{p} \sqrt{2s^*} \cdot (100\sqrt{np})^{j-1} \text{ for } j = \eta, \eta + 1.$$

455 Our assumption on n yields $\mu \geq C\sqrt{np} \log^6 n$; moreover, $|A|(p-q) \leq \frac{p-q}{\mu^2} \leq 1/s^* \cdot \frac{1}{\mu} \leq$
 456 $1/s^* \cdot (C\sqrt{np} \log^6 n)^{-1}$. Therefore,

$$\begin{aligned}
 &\|A(E^{\eta+1}H + E^\eta HE)\|_{2 \rightarrow \infty} \\
 &\leq A \left(\|E^{\eta+1}H\|_{2 \rightarrow \infty} + \|E^\eta H\|_{2 \rightarrow \infty} \|E\|_2 \right) \\
 &\leq \frac{1}{s^*} \cdot (C\sqrt{np} \log^6 n)^{-1} \cdot 500(\log n)^{5\eta+5} \cdot \sqrt{n} \cdot \sqrt{p} \cdot \sqrt{2s^*} \left((100\sqrt{np})^\eta + (100\sqrt{np})^{\eta-1} C_0 \sigma \sqrt{n} \right) \\
 &\leq 500000 \sqrt{np/s^*} (\log n)^{5\eta} (100\sqrt{np})^{\eta-1}. \quad (7)
 \end{aligned}$$

Similarly, by applying Proposition A.3 to $E^j J_n$ with $\alpha = p, \beta = 1, \gamma = n$, we have, with probability at least $1 - O(n^{-2})$,

$$\|E^j J_n\|_{2 \rightarrow \infty} \leq 500(\log n)^{5j} \cdot \sqrt{p} \cdot n \cdot (100\sqrt{np})^{j-1}, j = \eta, \eta + 1.$$

457 Since $|Aq| = \frac{q}{\lambda} \cdot \frac{1}{\mu} \leq \frac{1}{n} \cdot (C\sqrt{n} \log^6 n)^{-1}$, we have

$$\begin{aligned}
 &\|Aq(E^{\eta+1}J_n + E^\eta J_n E)\|_{2 \rightarrow \infty} \\
 &\leq |Aq| \left(\|E^{\eta+1}J_n\|_{2 \rightarrow \infty} + \|E^\eta J_n\|_{2 \rightarrow \infty} \|E\|_2 \right) \\
 &\leq \frac{1}{n} \cdot (C\sqrt{np} \log^6 n)^{-1} \cdot 500 \cdot (\log n)^{5\eta+5} \cdot n \cdot \left((100\sqrt{np})^\eta + (100\sqrt{np})^{\eta-1} C_0 \sigma \sqrt{n} \right) \\
 &\leq 500000 \sqrt{p} (\log n)^{5\eta} (100\sqrt{np})^{\eta-1}. \quad (8)
 \end{aligned}$$

458 Combining Equation (7) and Equation (8), we have, with probability at least $1 - O(n^{-2})$,
 459 $\|E^\eta L\|_{2 \rightarrow \infty} \leq C_2 \sqrt{np/s^*} (\log n)^{5\eta} (100\sqrt{np})^{\eta-1}$ where $C_2 \stackrel{\text{def}}{=} 10^6$.

460 For the second part, we decompose $F = F' + F''$, where F' is the intra-cluster part, i.e., $F'_{uv} = F_{uv}$ if
 461 $u, v \in V_\ell$ for some ℓ ; and $F'_{uv} = 0$ otherwise. Equipped with Claim D.1, we can apply Proposition A.3
 462 on $E^\eta F'$ (with $\alpha = p, \beta = 10/s^*, \gamma = 2s^*$), and $E^\eta F''$ (with $\alpha = p, \beta = \frac{5}{n}, \gamma = n$):

$$\begin{aligned} \|E^\eta F\|_{2 \rightarrow \infty} &\leq \|E^\eta F'\|_{2 \rightarrow \infty} + \|E^\eta F''\|_{2 \rightarrow \infty} \\ &\leq 20000 \log^{5\eta} \sqrt{n} \sqrt{p} \left(\frac{10}{s^*} \cdot \sqrt{2s^*} + \frac{10}{n} \cdot \sqrt{n} \right) (100\sqrt{n})^{\eta-1} \\ &\leq C_2 (\log n)^{5\eta} \sqrt{np/s^*} (100\sqrt{n})^{\eta-1}, \end{aligned}$$

463 where the second inequality holds with probability at least $1 - O(n^{-2})$. The lemma follows from a
 464 union bound over all $\eta \leq \log n$. □