

A Note On Spectral Clustering^{*†}

Pavel Kolev¹ and Kurt Mehlhorn²

1 Max-Planck-Institut für Informatik, Saarbrücken, Germany
pkolev@mpi-inf.mpg.de

2 Max-Planck-Institut für Informatik, Saarbrücken, Germany
mehlhorn@mpi-inf.mpg.de

Abstract

Spectral clustering is a popular and successful approach for partitioning the nodes of a graph into clusters for which the ratio of outside connections compared to the volume (sum of degrees) is small. In order to partition into k clusters, one first computes an approximation of the bottom k eigenvectors of the (normalized) Laplacian of G , uses it to embed the vertices of G into k -dimensional Euclidean space \mathbb{R}^k , and then partitions the resulting points via a k -means clustering algorithm. It is an important task for theory to explain the success of spectral clustering.

Peng et al. (COLT, 2015) made an important step in this direction. They showed that spectral clustering provably works if the gap between the $(k + 1)$ -th and the k -th eigenvalue of the normalized Laplacian is sufficiently large. They proved a structural and an algorithmic result. The algorithmic result needs a considerably stronger gap assumption and does not analyze the standard spectral clustering paradigm; it replaces spectral embedding by heat kernel embedding and k -means clustering by locality sensitive hashing.

We extend their work in two directions. Structurally, we improve the quality guarantee for spectral clustering by a factor of k and simultaneously weaken the gap assumption. Algorithmically, we show that the standard paradigm for spectral clustering works. Moreover, it even works with the same gap assumption as required for the structural result.

1998 ACM Subject Classification G.2.2 [Discrete Mathematics] Graph Theory, Graph Algorithms, G.3 [Probability and Statistics] Probabilistic Algorithms (including Monte Carlo)

Keywords and phrases spectral embedding, k -means clustering, power method, gap assumption

Digital Object Identifier 10.4230/LIPIcs.ESA.2016.57

1 Introduction

A *cluster* in an undirected graph $G = (V, E)$ is a set S of nodes whose volume is large compared to the number of outside connections. Formally, we define the *conductance* of S by $\phi(S) = |E(S, \bar{S})|/\mu(S)$, where $\mu(S) = \sum_{v \in S} \deg(v)$ is the *volume* of S . The k -way partitioning problem for graphs asks to partition the vertices of a graph such that the conductance of each block of the partition is small (formal definition below). This problem arises in many applications, e.g., image segmentation and exploratory data analysis. We refer to the survey [10] for additional information. A popular and very successful approach to clustering [4, 8, 10] is *spectral clustering*. One first computes an approximation of the bottom k eigenvectors of the (normalized) Laplacian of G , uses it to embed the vertices of G into k -dimensional Euclidean space \mathbb{R}^k , and then partitions the resulting points via a

* The full version of the paper is available at <http://arxiv.org/abs/1509.09188>.

† This work has been funded by the Cluster of Excellence “Multimodal Computing and Interaction” within the Excellence Initiative of the German Federal Government.

k -means clustering algorithm. It is an important task for theory to explain the success of spectral clustering. Recently, Peng et al. [7] made an important step in this direction. They showed that spectral clustering provably works if the $(k + 1)$ -th and the k -th eigenvalue of the normalized Laplacian differ sufficiently. In order to explain their result, we need some notation.

Let $\mathcal{L}_G = I - D^{-1/2}AD^{-1/2}$ be the normalized Laplacian matrix of G , where D is the diagonal degree matrix and A is the adjacency matrix, and let $f_j \in \mathbb{R}^V$ be the eigenvector corresponding to the j -th smallest eigenvalue λ_j of \mathcal{L}_G . The *spectral embedding map* $F : V \rightarrow \mathbb{R}^k$ is defined by

$$F(u) = \frac{1}{\sqrt{d_u}} (f_1(u), \dots, f_k(u))^T, \quad \text{for all vertices } u \in V. \quad (1)$$

Peng et al. [7] construct a k -means instance \mathcal{X}_V by inserting d_u many copies of the vector $F(u)$ into \mathcal{X}_V , for every vertex $u \in V$.

Let \mathcal{X} be a set of vectors of the same dimension. Then

$$\Delta_k(\mathcal{X}) \triangleq \min_{\text{partition } (X_1, \dots, X_k) \text{ of } \mathcal{X}} \sum_{i=1}^k \sum_{x \in X_i} \|x - c_i\|^2, \quad \text{where } c_i = \frac{1}{|X_i|} \sum_{x \in X_i} x,$$

is the optimal cost of clustering \mathcal{X} into k sets. An α -approximate clustering algorithm returns a k -way partition (A_1, \dots, A_k) and centers c_1, \dots, c_k such that

$$\text{Cost}(\{A_i, c_i\}_{i=1}^k) \triangleq \sum_{i=1}^k \sum_{x \in A_i} \|x - c_i\|^2 \leq \alpha \cdot \Delta_k(\mathcal{X}). \quad (2)$$

The *order k conductance constant* $\rho(k)$ is a well studied worst case guarantee for k -way partitioning that is defined by

$$\rho(k) = \min_{\text{disjoint nonempty } Z_1, \dots, Z_k} \Phi(Z_1, \dots, Z_k), \quad \text{where } \Phi(Z_1, \dots, Z_k) = \max_{i \in [1:k]} \phi(Z_i). \quad (3)$$

Lee et al. [3] connected $\rho(k)$ and the k -th smallest eigenvalue of the normalized Laplacian matrix \mathcal{L}_G through the relation, also known as higher order Cheeger inequality,

$$\lambda_k/2 \leq \rho(k) \leq O(k^2)\sqrt{\lambda_k}. \quad (4)$$

In this work, we focus attention on the *order k partition constant* $\widehat{\rho}(k)$ of G , defined by

$$\widehat{\rho}(k) \triangleq \min_{\text{partition } (P_1, \dots, P_k) \text{ of } V} \Phi(P_1, \dots, P_k), \quad \text{where } \Phi(Z_1, \dots, Z_k) = \max_{i \in [1:k]} \phi(Z_i).$$

In a consecutive work, inspired by partitioning graphs into expanders, Oveis Gharan and Trevisan [6] proved the following relation

$$\rho(k) \leq \widehat{\rho}(k) \leq k\rho(k). \quad (5)$$

We are now ready to state the main structural result by Peng et al [7].

► **Theorem 1.1** ([7, Theorem 1.2]). *Let $k \geq 3$ and (P_1, \dots, P_k) be a k -way partition of V with $\Phi(P_1, \dots, P_k) = \widehat{\rho}(k)$. Let G be a graph that satisfies the gap assumption¹*

$$\delta = \frac{2 \cdot 10^5 \cdot k^3}{\Upsilon} \in (0, 1/2], \quad \text{where } \Upsilon \triangleq \frac{\lambda_{k+1}}{\widehat{\rho}(k)}. \quad (6)$$

¹ Note that $\lambda_k/2 \leq \widehat{\rho}(k)$, see (4,5). Thus the assumption implies $\lambda_k/2 \leq \widehat{\rho}(k) = \delta \lambda_{k+1}/(2 \cdot 10^5 \cdot k^3)$, i.e., there is a substantial gap between the k -th and the $(k + 1)$ -th eigenvalue.

Let (A_1, \dots, A_k) be the k -way partition² of V returned by an α -approximate k -means algorithm applied to \mathcal{X}_V . Then the following statements hold (after suitable renumbering of one of the partitions):

1. $\mu(A_i \Delta P_i) \leq \alpha\delta \cdot \mu(P_i)$, and
 2. $\phi(A_i) \leq (1 + 2\alpha\delta) \cdot \phi(P_i) + 2\alpha\delta$,
- where the symmetric difference $A_i \Delta P_i = (A_i \setminus P_i) \cup (P_i \setminus A_i)$.

Under the stronger gap assumption $\delta = 2 \cdot 10^5 \cdot k^5 / \Upsilon \in (0, 1/2]$, they showed how to obtain a partition in time $O(m \cdot \text{poly}(n))$ with essentially the guarantees stated in Theorem 1.1, where $m = |E|$ is the number of edges in G and $n = |V|$ is the number of nodes.

However, their algorithmic result does not analyze the standard spectral clustering paradigm, since it replaces spectral embedding by heat kernel embedding and k -means clustering by locality sensitive hashing. Therefore, their algorithmic result does not explain the success of the standard spectral clustering paradigm.

Our Results

We strengthen the approximation guarantees in Theorem 1.1 by a factor of k and simultaneously weaken the gap assumption. As a consequence, the variant of Lloyd's k -means algorithm analyzed by Ostrovsky et al. [5] applied to³ $\widetilde{\mathcal{X}}_V$ achieves the improved approximation guarantees in time $O(m(k^2 + \frac{\ln n}{\lambda_{k+1}}))$ with constant probability. Table 1 summarizes these results.

Let \mathcal{O} be the set of all k -way partitions (P_1, \dots, P_k) with $\Phi(P_1, \dots, P_k) = \widehat{\rho}(k)$, i.e., the set of all partitions that achieve the order k partition constant. Let

$$\widehat{\rho}_{\text{avr}}(k) \triangleq \min_{(P_1, \dots, P_k) \in \mathcal{O}} \frac{1}{k} \sum_{i=1}^k \phi(P_i)$$

be the *minimal average conductance* over all k -way partitions in \mathcal{O} . The minimal average conductance can be considerably smaller than the order k partition constant. Consider a graph consisting of one clique of size $S_k = f(n) = o(n/k^{3/2})$, $k - 1$ cliques of size $(n - f(n))/(k - 1)$ each, and k additional edges that connect the cliques in the form of a ring. Then $\phi(S_i) \approx k^2/n^2$ for $1 \leq i \leq k - 1$ and $\phi(S_k) \approx 1/f(n)^2$. Thus $\widehat{\rho}(k) = \max_i \phi(S_i) = \phi(S_k) \approx 1/f(n)^2$ and $\widehat{\rho}_{\text{avr}}(k) = (1/k) \sum_{1 \leq i \leq k} \phi(S_i) \approx k^2/n^2 + (1/k) \cdot (1/f(n)^2) \approx \widehat{\rho}(k)/k$.

For the remainder of this paper we denote by (P_1, \dots, P_k) a k -way partition of V that achieves $\widehat{\rho}_{\text{avr}}(k)$. In the full version of the paper, we give an analogous relation to (5) for $\widehat{\rho}_{\text{avr}}(k)$. We state now our main result.

► **Theorem 1.2** (Main Theorem).

- (a) (Existence of a Good Clustering) *Let $k \geq 3$. Let G be a graph satisfying the gap assumption*

$$\delta = \frac{20^4 \cdot k^3}{\Psi} \in (0, 1/2], \quad \text{where} \quad \Psi \triangleq \frac{\lambda_{k+1}}{\widehat{\rho}_{\text{avr}}(k)}. \tag{7}$$

Let (A_1, \dots, A_k) be the k -way partition returned by an α -approximate clustering algorithm applied to the spectral embedding \mathcal{X}_V . Then for every $i \in [1 : k]$ the following two statements hold (after suitable renumbering of one of the partitions):

² The k -means algorithm returns a partition of \mathcal{X}_V . One may assume w.l.o.g. that all copies of $F(u)$ are put into the same cluster of \mathcal{X}_V . Thus the algorithm also partitions V .

³ $\widetilde{\mathcal{X}}_V$ is defined as \mathcal{X}_V but in terms of approximate eigenvectors, see Subsection 2.3.

■ **Table 1** A comparison of the results in Peng et al. [7] and our results. The parameter $\delta \in (0, 1/2]$ relates the approximation guarantees with the gap assumption.

	Gap Assumption	Partition Quality	Running Time
Peng et al. [7]	$\delta = 2 \cdot 10^5 \cdot k^3 / \Upsilon$	$\mu(A_i \Delta P_i) \leq \alpha \delta \cdot \mu(P_i)$ $\phi(A_i) \leq (1 + 2\alpha\delta) \phi(P_i) + 2\alpha\delta$	Existential result
This paper	$\delta = 20^4 \cdot k^3 / \Psi$	$\mu(A_i \Delta P_i) \leq \frac{\alpha\delta}{10^3 k} \cdot \mu(P_i)$ $\phi(A_i) \leq \left(1 + \frac{2\alpha\delta}{10^3 k}\right) \phi(P_i) + \frac{2\alpha\delta}{10^3 k}$	Existential result
Peng et al. [7]	$\delta = 2 \cdot 10^5 \cdot k^5 / \Upsilon$	$\mu(A_i \Delta P_i) \leq \frac{\delta \log^2 k}{k^2} \cdot \mu(P_i)$ $\phi(A_i) \leq \left(1 + \frac{2\delta \log^2 k}{k^2}\right) \phi(P_i) + \frac{2\delta \log^2 k}{k^2}$	$O(m \cdot \text{poly log}(n))$
This paper	$\delta = 20^4 \cdot k^3 / \Psi$ $\delta \leq k/10^9$ $\Delta_k(\mathcal{X}_V) \geq n^{-O(1)}$	$\mu(A_i \Delta P_i) \leq \frac{2\delta}{10^3 k} \cdot \mu(P_i)$ $\phi(A_i) \leq \left(1 + \frac{4\delta}{10^3 k}\right) \phi(P_i) + \frac{4\delta}{10^3 k}$	$O\left(m \left(k^2 + \frac{\ln n}{\lambda_{k+1}}\right)\right)$

1. $\mu(A_i \Delta P_i) \leq \frac{\alpha\delta}{10^3 k} \cdot \mu(P_i)$, and
 2. $\phi(A_i) \leq \left(1 + \frac{2\alpha\delta}{10^3 k}\right) \cdot \phi(P_i) + \frac{2\alpha\delta}{10^3 k}$.
- (b) (An Efficient Algorithm) *If in addition $\delta \leq k/10^9$ and⁴ $\Delta_k(\mathcal{X}_V) \geq n^{-O(1)}$, then the variant of Lloyd's algorithm analyzed by Ostrovsky et al. [5] applied to $\widetilde{\mathcal{X}}_V$ returns in time $O(m(k^2 + \frac{\ln n}{\lambda_{k+1}}))$ with constant probability a partition (A_1, \dots, A_k) such that for every $i \in [1 : k]$ the following two statements hold (after suitable renumbering of one of the partitions):*
3. $\mu(A_i \Delta P_i) \leq \frac{2\delta}{10^3 k} \cdot \mu(P_i)$, and
 4. $\phi(A_i) \leq \left(1 + \frac{4\delta}{10^3 k}\right) \cdot \phi(P_i) + \frac{4\delta}{10^3 k}$.

Part (b) of Theorem 1.2 gives a theoretical support for the practical success of spectral clustering based on approximate spectral embedding followed by k -means clustering. Moreover, if $k \leq \text{poly}(\log n)$ and $\lambda_{k+1} \geq \text{poly}(\log n)$, our algorithm works in nearly linear time. Previous papers [3, 7, 9] replaced k -means clustering by other techniques for their algorithmic results.

The k -means algorithm in [5] is efficient only for inputs \mathcal{X} for which some partition into k clusters is much better than any partition into $k - 1$ clusters. The authors proved that the algorithm is efficient for inputs \mathcal{X} satisfying $\Delta_k(\mathcal{X}) \leq \varepsilon^2 \Delta_{k-1}(\mathcal{X})$ for some $\varepsilon \in (0, \varepsilon_0]$, where $\varepsilon_0 = 6/10^7$, stated that the result should also hold for a larger ε_0 , and mentioned that they did not attempt to maximize ε_0 . For the proof of part (b) of Theorem 1.2, we show in Section 5 that $\widetilde{\mathcal{X}}_V$ satisfies this assumption. In this proof, we need $\delta \leq k \cdot \varepsilon_0 / 600 = k/10^9$.

One of the reviewers suggested to include a numerical example. Consider a graph consisting of k cliques of size n/k each plus k additional edges that connect the cliques in the form of a ring. Such a graph is about the easiest input for a clustering algorithm. Then $\widehat{\rho}_{\text{avr}}(k) = \widehat{\rho}(k) \approx (k/n)^2$. For the gap assumption to hold we need $\lambda_{k+1} \geq 2 \cdot 20^4 \cdot k^3 \cdot \widehat{\rho}_{\text{avr}}(k)$. Since $\lambda_{k+1} \leq 2$, this implies $n \geq 400 \cdot k^{2.5}$. For small k , this is a modest requirement on the size of the graph.

⁴ The case $\Delta_k(\mathcal{X}_V) \leq n^{-O(1)}$ constitutes a trivial clustering problem. For technical reasons, we have to exclude too easy inputs.

For the algorithmic result, we need in addition $\delta \leq k \cdot \varepsilon_0/600$. For the gap condition to hold, we need $2 \geq \lambda_{k+1} \geq (600/\varepsilon_0 k) \cdot 20^4 \cdot k^3 \cdot (k^2/n^2)$ or $n \geq 4\sqrt{3} \cdot 10^3 \cdot k^2/\sqrt{\varepsilon_0}$. For $\varepsilon_0 = 6/10^7$, this amounts to $n \geq 4\sqrt{5} \cdot 10^6 \cdot k^2$, a quite large lower bound on n .

Our statement above that Part (b) of Theorem 1.2 gives a theoretical support for the practical success of spectral clustering based on approximate spectral embedding followed by k -means clustering therefore has to be taken with a grain of salt. It is only an asymptotic statement and does not explain the good behavior on small graphs.

2 Highlights of Our Technical Contribution

2.1 Exact Spectral Embedding – Notation

We use the notation adopted by Peng et al. [7]. Let $f_j \in \mathbb{R}^V$ be the eigenvector corresponding to the j -th smallest eigenvalue λ_j of \mathcal{L}_G , and let $\bar{g}_i = \frac{D^{1/2}\chi_{P_i}}{\|D^{1/2}\chi_{P_i}\|}$ be the normalized indicator vector associated with the i -th optimal cluster $P_i \subset V$.

Since the eigenvectors $\{f_i\}_{i=1}^n$ form an orthonormal basis of \mathbb{R}^n , each normalized indicator vector \bar{g}_i can be expressed as $\bar{g}_i = \sum_{j=1}^n \alpha_j^{(i)} f_j$, for all $i \in [1 : k]$. Its *projection* into the subspace spanned by the bottom k eigenvectors is given by $\hat{f}_i = \sum_{j=1}^k \alpha_j^{(i)} f_j$. Peng et al. [7] proved that if the gap parameter Υ is large enough then $\text{span}(\{\hat{f}_i\}_{i=1}^k) = \text{span}(\{f_i\}_{i=1}^k)$ and hence the bottom k eigenvectors can be expressed by $f_i = \sum_{j=1}^k \beta_j^{(i)} \hat{f}_j$, for all $i \in [1 : k]$. We show that similar statements hold with substituted gap parameter Ψ .

A corner stone in the analysis of spectral clustering is to prove the existence of exactly k directions near which all spectrally embedded vectors are closely concentrated. These estimation centers are defined by

$$p^{(i)} = \frac{1}{\sqrt{\mu(P_i)}} \left(\beta_i^{(1)}, \dots, \beta_i^{(k)} \right)^T. \tag{8}$$

Our analysis crucially relies on the isometric properties of the following square matrix. Let $\mathbf{B} \in \mathbb{R}^{k \times k}$ be a matrix defined by $\mathbf{B}_{j,i} = \beta_j^{(i)}$, for every $i, j \in [1 : k]$.

2.2 Exact Spectral Embedding – Structural Results

The proof of Theorem 1.2 (a) follows the proof-structure of [7, Theorem 1.2] in Peng et al., but improves upon it in essential ways.

Our key technical insight is that the matrices $\mathbf{B}\mathbf{B}^T$ and $\mathbf{B}^T\mathbf{B}$ are close to the identity matrix. The proof of Theorem 2.1 appears in the full version of the paper.

► **Theorem 2.1** (Matrix $\mathbf{B}\mathbf{B}^T$ is Close to Identity Matrix). *If $\Psi \geq 10^4 \cdot k^3/\varepsilon^2$ and $\varepsilon \in (0, 1)$ then for all distinct $i, j \in [1 : k]$ it holds*

$$1 - \varepsilon \leq \langle \mathbf{B}_{i,:}, \mathbf{B}_{i,:} \rangle \leq 1 + \varepsilon \quad \text{and} \quad |\langle \mathbf{B}_{i,:}, \mathbf{B}_{j,:} \rangle| \leq \sqrt{\varepsilon}.$$

Informally, Theorem 2.1 implies that each normalized indicator vector \bar{g}_i is close to the corresponding eigenvector f_i . This gives a simple and intuitive explanation for the success of spectral clustering.

To see this, let $\mathbf{F}_k, \hat{\mathbf{F}}_k \in \mathbb{R}^{k \times k}$ be matrices whose i -th column is f_i and \hat{f}_i , respectively. The projection matrix \mathbf{P}_k into the k -th principle subspace of \mathcal{L}_G is given by $\mathbf{P}_k = \mathbf{F}_k \mathbf{F}_k^T$ and since $\hat{\mathbf{F}}_k \mathbf{B} = \mathbf{F}_k$, by Theorem 2.1 it follows that $\hat{\mathbf{F}}_k \hat{\mathbf{F}}_k^T \approx \mathbf{F}_k \mathbf{F}_k^T$. Therefore, each projected indicator vector satisfies $\hat{f}_i \approx f_i$. This implies $\alpha^{(i)} \approx \chi_i$ and hence we have $\bar{g}_i \approx f_i$.

Formally, Theorem 2.1 allows us to improve the separation guarantee between any pair of estimation centers by a factor of k over [7, Lemma 4.3], measured in terms of the Euclidean distance.

► **Lemma 2.2.** *If $\delta = 20^4 \cdot k^3 / \Psi \in (0, 1]$ then for every $i \in [1 : k]$ it holds that $\|p^{(i)}\|^2 \in [1 \pm \sqrt{\delta}/4] \frac{1}{\mu(P_i)}$.*

Proof. By definition $p^{(i)} = \frac{1}{\sqrt{\mu(P_i)}} \cdot \mathbf{B}_{i,:}$ and Theorem 2.1 yields $\|\mathbf{B}_{i,:}\|^2 \in [1 \pm \sqrt{\delta}/4]$. ◀

► **Lemma 2.3 (Larger Distance Between Estimation Centers).** *If $\delta = 20^4 \cdot k^3 / \Psi \in (0, 1/2]$ then for any distinct $i, j \in [1 : k]$ it holds that $\|p^{(i)} - p^{(j)}\|^2 \geq [2 \cdot \min\{\mu(P_i), \mu(P_j)\}]^{-1}$.*

Proof. Since $p^{(i)}$ is a row of matrix B , Theorem 2.1 with $\varepsilon = \sqrt{\delta}/4$ yields

$$\left\langle \frac{p^{(i)}}{\|p^{(i)}\|}, \frac{p^{(j)}}{\|p^{(j)}\|} \right\rangle = \frac{\langle \mathbf{B}_{i,:}, \mathbf{B}_{j,:} \rangle}{\|\mathbf{B}_{i,:}\| \|\mathbf{B}_{j,:}\|} \leq \frac{\sqrt{\varepsilon}}{1 - \varepsilon} = \frac{2\delta^{1/4}}{3}.$$

W.l.o.g. assume that $\|p^{(i)}\|^2 \geq \|p^{(j)}\|^2$, say $\|p^{(j)}\| = \alpha \cdot \|p^{(i)}\|$ for some $\alpha \in (0, 1]$. Then by Lemma 2.2 we have $\|p^{(i)}\|^2 \geq (1 - \sqrt{\delta}/4) \cdot [\min\{\mu(P_i), \mu(P_j)\}]^{-1}$, and hence

$$\begin{aligned} \|p^{(i)} - p^{(j)}\|^2 &= \|p^{(i)}\|^2 + \|p^{(j)}\|^2 - 2 \left\langle \frac{p^{(i)}}{\|p^{(i)}\|}, \frac{p^{(j)}}{\|p^{(j)}\|} \right\rangle \|p^{(i)}\| \|p^{(j)}\| \\ &\geq \left(\alpha^2 - \frac{4\delta^{1/4}}{3} \cdot \alpha + 1 \right) \|p^{(i)}\|^2 \geq [2 \cdot \min\{\mu(P_i), \mu(P_j)\}]^{-1}. \end{aligned} \quad \blacktriangleleft$$

The observation that Υ can be replaced by Ψ in all statements in [7] is technically easy. However, this is crucial for Theorem 1.2 (b), since it yields an improved version of [7, Lemma 4.5] showing that a weaker by a factor of k assumption is sufficient. Due to space limitation, we defer the proof of Lemma 2.4 to the full version of the paper.

► **Lemma 2.4.** *Let (P_1, \dots, P_k) and (A_1, \dots, A_k) are partitions of the vector set. Suppose for every permutation $\pi : [1 : k] \rightarrow [1 : k]$ there is an index $i \in [1 : k]$ such that*

$$\mu(A_i \triangle P_{\pi(i)}) \geq \frac{2\varepsilon}{k} \cdot \mu(P_{\pi(i)}), \quad (9)$$

where $\varepsilon \in (0, 1)$ is a parameter. If $\delta = 20^4 \cdot k^3 / \Psi \in (0, 1/2]$ and $\varepsilon \geq 64\alpha \cdot k^3 / \Psi$ then

$$\text{Cost}(\{A_i, c_i\}_{i=1}^k) > \frac{2k^2}{\Psi} \alpha.$$

With the above Lemmas in place, the proof of Theorem 1.2 (a) is then completed as in [7]. We give more details in Section 3.

Before we turn to Theorem 1.2 (b), we consider the variant of Lloyd's algorithm analyzed by Ostrovsky et al. [5] applied to \mathcal{X}_V . This algorithm is efficient for inputs \mathcal{X} satisfying: some partition into k clusters is much better than any partition into $k - 1$ clusters.

► **Theorem 2.5.** [5, Theorem 4.15] *Assuming that $\Delta_k(\mathcal{X}) \leq \varepsilon^2 \Delta_{k-1}(\mathcal{X})$ for $\varepsilon \in (0, 6/10^7]$, there is an algorithm that returns a solution of cost at most $[(1 - \varepsilon^2)/(1 - 37\varepsilon^2)] \Delta_k(\mathcal{X})$ with probability at least $1 - O(\sqrt{\varepsilon})$ in time $O(nkd + k^3d)$.*

In Section 4, we establish the assumption of Ostrovsky et al. [5] for \mathcal{X}_V .

► **Theorem 2.6** (Normalized Spectral Embedding is ε -separated). *Let G be a graph that satisfies the gap assumption $\delta = 20^4 \cdot k^3 / \Psi \in (0, 1/2]$ and $\delta \leq k \cdot \varepsilon / 600$, where $\varepsilon = 6/10^7$ is the Ostrovsky et al.'s constant. Then it holds*

$$\Delta_k(\mathcal{X}_V) \leq \varepsilon^2 \Delta_{k-1}(\mathcal{X}_V). \quad (10)$$

However, Theorem 2.6 is insufficient for Theorem 1.2 (b), since we need a similar result for the set $\widetilde{\mathcal{X}}_V$ formed by approximate eigenvectors. To overcome this issue we build upon the recent work by Boutsidis et al. [1] which shows that running an approximate k -means clustering algorithm on approximate eigenvectors obtained via the power method, yields an additive approximation to solving the k -means clustering problem on exact eigenvectors.

In order to state the connection, we need to introduce some of their notation.

2.3 Approximate Spectral Embedding – Notation

Let $Z \in \mathbb{R}^{n \times k}$ be a matrix whose rows represent n vectors that are to be partitioned into k clusters. For every k -way partition we associate an indicator matrix $X \in \mathbb{R}^{n \times k}$ that satisfies $X_{ij} = 1/\sqrt{|C_j|}$ if the i -th row $Z_{i,:}$ belongs to the j -th cluster C_j , and $X_{ij} = 0$ otherwise. We denote the optimal indicator matrix X_{opt} by

$$X_{\text{opt}} = \arg \min_{X \in \mathbb{R}^{n \times k}} \|Z - XX^T Z\|_F^2 = \arg \min_{X \in \mathbb{R}^{n \times k}} \sum_{j=1}^k \sum_{u \in X_j} \|Z_{u,:} - c_j\|_2^2, \quad (11)$$

where $c_j = (1/|X_j|) \sum_{u \in X_j} Z_{u,:}$ is the center point of cluster C_j .

The normalized Laplacian matrix $\mathcal{L}_G \in \mathbb{R}^{n \times n}$ of a graph G is define by $\mathcal{L}_G = I - \mathcal{A}$, where $\mathcal{A} = D^{-1/2}AD^{-1/2}$ is the normalized adjacency matrix. Let $U_k \in \mathbb{R}^{n \times k}$ be a matrix composed of the bottom k orthonormal eigenvectors of \mathcal{L}_G corresponding to the smallest eigenvalues $\lambda_1, \dots, \lambda_k$. We define by $Y \triangleq U_k$ the canonical spectral embedding.

Our approximate spectral embedding is computed by the so called “**Power method**”. Let $S \in \mathbb{R}^{n \times k}$ be a matrix whose entries are i.i.d. samples from the standard Gaussian distribution $N(0, 1)$ and p be a positive integer. Then the approximate spectral embedding \widetilde{Y} is defined by the following process:

$$1) \mathcal{B} \triangleq I + \mathcal{A}; \quad 2) \text{ Let } \widetilde{U}\widetilde{\Sigma}\widetilde{V}^T \text{ be the SVD of } \mathcal{B}^p S; \quad \text{and} \quad 3) \widetilde{Y} \triangleq \widetilde{U} \in \mathbb{R}^{n \times k}. \quad (12)$$

We proceed by defining the normalized (approximate) spectral embedding. We construct a matrix $Y' \in \mathbb{R}^{m \times k}$ such that for every vertex $u \in V$ we add $\deg(u)$ many copies of the normalized row $U_k(u, :)/\sqrt{\deg(u)}$ to Y' . Formally, the normalized (approximate) spectral embedding Y' (\widetilde{Y}') is defined by

$$Y' = \begin{pmatrix} \mathbf{1}_{\deg(1)} \frac{U_k(1,:)}{\sqrt{\deg(1)}} \\ \dots \\ \mathbf{1}_{\deg(n)} \frac{U_k(n,:)}{\sqrt{\deg(n)}} \end{pmatrix}_{m \times k} \quad \text{and} \quad \widetilde{Y}' = \begin{pmatrix} \mathbf{1}_{\deg(1)} \frac{\widetilde{U}(1,:)}{\sqrt{\deg(1)}} \\ \dots \\ \mathbf{1}_{\deg(n)} \frac{\widetilde{U}(n,:)}{\sqrt{\deg(n)}} \end{pmatrix}_{m \times k}, \quad (13)$$

where $\mathbf{1}_{\deg(i)}$ is all-one column vector with dimension $\deg(i)$.

Similarly to (11) we associate to Y' (\widetilde{Y}') an indicator matrix X' (\widetilde{X}') that satisfies $X'_{ij} = 1/\sqrt{\mu(C_j)}$ if the i -th row $Y'_{i,:}$ belongs to the j -th cluster C_j , and $X'_{ij} = 0$ otherwise. We may assume w.l.o.g. that a k -means algorithm outputs an indicator matrix X' such that all copies of row $U_k(v, :)/\sqrt{\deg(v)}$ belong to the same cluster, for every vertex $v \in V$.

We associate to matrices Y' and \widetilde{Y}' the sets of points \mathcal{X}_V and $\widetilde{\mathcal{X}}_V$ respectively. We present now a key connection between the spectral embedding map $F(\cdot)$, the optimal k -means cost $\Delta_k(\mathcal{X}_V)$ and matrices Y', X'_{opt} :

$$\left\| Y' - X'_{\text{opt}} (X'_{\text{opt}})^{\text{T}} Y' \right\|_F^2 = \sum_{j=1}^k \sum_{v \in C_j^*} \deg(v) \|F(v) - c_j^*\|_F^2 = \Delta_k(\mathcal{X}_V), \quad (14)$$

where each center satisfies $c_j^* = \mu(C_j^*)^{-1} \cdot \sum_{v \in C_j^*} \deg(v) F(v)$ and $F(v) = Y_{v,:} / \sqrt{\deg(v)}$.

2.4 Approximate Spectral Embedding – Algorithmic Results

Our analysis relies on the proof techniques developed in [1, 2]. By adjusting these techniques (c.f. [1, Lemma 5] and [2, Lemma 7]) to our setting, we prove (in the full version of the paper) the following result for the symmetric positive semi-definite matrix \mathcal{B} whose largest k singular values (eigenvalues) correspond to the eigenvectors u_1, \dots, u_k of \mathcal{L}_G .

► **Lemma 2.7.** *Let $\widetilde{U}\widetilde{\Sigma}\widetilde{V}^{\text{T}}$ be the SVD of $\mathcal{B}^p S \in \mathbb{R}^{n \times k}$, where $p \geq 1$ and S is an $n \times k$ matrix of i.i.d. standard Gaussians. Let $\gamma_k = \frac{2-\lambda_{k+1}}{2-\lambda_k} < 1$ and fix $\delta, \epsilon \in (0, 1)$. Then for any $p \geq \ln(8nk/\epsilon\delta) / \ln(1/\gamma_k)$ with probability at least $1 - 2e^{-2n} - 3\delta$ it holds*

$$\left\| U_k U_k^{\text{T}} - \widetilde{U}\widetilde{U}^{\text{T}} \right\|_F \leq \epsilon.$$

We establish several technical Lemmas that combined with Lemma 2.7 allow us to apply the proof techniques in [1, Theorem 6]. More precisely, we prove in Subsection 5.1 that running an approximate k -means algorithm on a normalized approximate spectral embedding \widetilde{Y}' computed by the power method, yields an approximate clustering of the normalized spectral embedding Y' .

► **Theorem 2.8.** *Compute matrix \widetilde{Y}' via the power method with $p \geq \ln(8nk/\epsilon\delta) / \ln(1/\gamma_k)$, where $\gamma_k = (2 - \lambda_{k+1}) / (2 - \lambda_k) < 1$. Run on the rows of \widetilde{Y}' an α -approximate k -means algorithm with failure probability δ_α . Let the outcome be a clustering indicator matrix $\widetilde{X}'_\alpha \in \mathbb{R}^{n \times k}$. Then with probability at least $1 - 2e^{-2n} - 3\delta_p - \delta_\alpha$ it holds*

$$\left\| Y' - \widetilde{X}'_\alpha (\widetilde{X}'_\alpha)^{\text{T}} Y' \right\|_F^2 \leq (1 + 4\epsilon) \cdot \alpha \cdot \left\| Y' - X'_{\text{opt}} (X'_{\text{opt}})^{\text{T}} Y' \right\|_F^2 + 4\epsilon^2.$$

Our main technical contribution is to prove, in Subsection 5.2, that $\widetilde{\mathcal{X}}_V$ satisfies the assumption of Ostrovsky et al. [5]. Our analysis builds upon Theorem 2.6 and Theorem 2.8.

► **Theorem 2.9** (Approximate Normalized Spectral Embedding is ϵ -separated). *Let G be a graph that satisfies the gap assumption $\delta = 20^4 \cdot k^3 / \Psi \in (0, 1/2]$ and $\delta \leq k \cdot \epsilon / 600$, where $\epsilon = 6/10^7$ is the Ostrovsky et al.'s constant. If the optimum cost⁵ $\|Y' - X'_{\text{opt}} (X'_{\text{opt}})^{\text{T}} Y'\|_F \geq n^{-O(1)}$ and the matrix \widetilde{Y}' is constructed via the power method with $p \geq \Omega(\frac{\ln n}{\lambda_{k+1}})$, then w.h.p it holds $\Delta_k(\widetilde{\mathcal{X}}_V) < 5\epsilon^2 \cdot \Delta_{k-1}(\widetilde{\mathcal{X}}_V)$.*

Based on the preceding results, we prove Theorem 1.2 (b) in Subsection 5.3.

⁵ $\|Y' - X'_{\text{opt}} (X'_{\text{opt}})^{\text{T}} Y'\|_F \geq n^{-O(1)}$ asserts a multiplicative approximation guarantee in Theorem 2.8.

3 The Proof of Part (a) of Theorem 1.2

The proof of part (a.1) builds upon the following Lemmas. Recall that \mathcal{X}_V contains d_u copies of $F(u)$ for each $u \in V$. W.l.o.g. we may restrict attention to clusterings of \mathcal{X}_V that put all copies of $F(u)$ into the same cluster and hence induce a clustering of V . Let (A_1, \dots, A_k) with cluster centers c_1 to c_k be a clustering of V . Its k -means cost is

$$\text{Cost}(\{A_i, c_i\}_{i=1}^k) = \sum_{i=1}^k \sum_{u \in A_i} d_u \|F(u) - c_i\|^2.$$

The proofs of Lemma 3.1 and Lemma 3.2 appear in the full version of the paper.

► **Lemma 3.1** ((P_1, \dots, P_k) is a good k -means partition). *If $\Psi > 4 \cdot k^{3/2}$ then there are vectors $\{p^{(i)}\}_{i=1}^k$ such that $\text{Cost}(\{P_i, p^{(i)}\}_{i=1}^k) \leq (1 + \frac{3k}{\Psi}) \cdot \frac{k^2}{\Psi}$.*

► **Lemma 3.2** (Only partitions close to (P_1, \dots, P_k) are good). *Under the hypothesis of Theorem 1.2, the following holds. If for every permutation $\sigma : [1 : k] \rightarrow [1 : k]$ there exists an index $i \in [1 : k]$ such that*

$$\mu(A_i \Delta P_{\sigma(i)}) \geq \frac{8\alpha\delta}{10^4 k} \cdot \mu(P_{\sigma(i)}), \quad \text{then it holds } \text{Cost}(\{A_i, c_i\}_{i=1}^k) > \frac{2\alpha k^2}{\Psi}.$$

We note that Lemma 3.2 follows directly by applying Lemma 2.4 with $\varepsilon = 64\alpha \cdot k^3/\Psi$. Substituting these bounds into (2) yields a contradiction, since

$$\frac{2\alpha k^2}{\Psi} < \text{Cost}(\{A_i, c_i\}_{i=1}^k) \leq \alpha \cdot \Delta_k(\mathcal{X}_V) \leq \alpha \cdot \text{Cost}(\{P_i, p^{(i)}\}_{i=1}^k) \leq \left(1 + \frac{3k}{\Psi}\right) \cdot \frac{\alpha k^2}{\Psi}.$$

Therefore, there exists a permutation π (the identity after suitable renumbering of one of the partitions) such that $\mu(A_i \Delta P_i) < \frac{8\alpha\delta}{10^4 k} \cdot \mu(P_i)$ for all $i \in [1 : k]$.

Part (a.2) follows from part (a.1). Indeed, for $\delta' = 8\delta/10^4$ we have

$$\mu(A_i) \geq \mu(P_i \cap A_i) = \mu(P_i) - \mu(P_i \setminus A_i) \geq \mu(P_i) - \mu(A_i \Delta P_i) \geq \left(1 - \frac{\alpha\delta'}{k}\right) \cdot \mu(P_i)$$

and $|E(A_i, \overline{A_i})| \leq |E(P_i, \overline{P_i})| + \mu(A_i \Delta P_i)$ since every edge that is counted in $|E(A_i, \overline{A_i})|$ but not in $|E(P_i, \overline{P_i})|$ must have an endpoint in $A_i \Delta P_i$. Thus

$$\Phi(A_i) = \frac{|E(A_i, \overline{A_i})|}{\mu(A_i)} \leq \frac{|E(P_i, \overline{P_i})| + \frac{\alpha\delta'}{k} \cdot \mu(P_i)}{\left(1 - \frac{\alpha\delta'}{k}\right) \cdot \mu(P_i)} \leq \left(1 + \frac{2\alpha\delta'}{k}\right) \cdot \phi(P_i) + \frac{2\alpha\delta'}{k}.$$

This completes the proof of Theorem 1.2 (a).

4 The Normalized Spectral Embedding is ε -separated

In this section, we prove that the normalized spectral embedding \mathcal{X}_V is ε -separated.

Proof of Theorem 2.6

We establish first a lower bound on $\Delta_{k-1}(\mathcal{X}_V)$.

► **Lemma 4.1.** *Let G be a graph that satisfies the gap assumption $\delta = 20^4 \cdot k^3/\Psi \in (0, 1/2]$. Then for $\delta' = 2\delta/20^4$ it holds $\Delta_{k-1}(\mathcal{X}_V) \geq 1/12 - \delta'/k$.*

57:10 A Note On Spectral Clustering

Before we prove Lemma 4.1 we show that it implies (10). By Lemma 3.1 we have $\Delta_k(\mathcal{X}_V) \leq 2k^2/\Psi = \delta'/k$. Then, we apply Lemma 4.1 with $\delta \leq k \cdot \varepsilon/600$, where $\varepsilon = 6/10^7$ is the Ostrovsky et al.'s constant, yielding

$$\Delta_{k-1}(\mathcal{X}_V) \geq \frac{1}{12} - \frac{\delta'}{k} = \frac{1}{12} - \frac{2}{20^4} \cdot \frac{\delta}{k} \geq \frac{10^{10}}{9 \cdot 2^5} \cdot \frac{\delta}{k} = \frac{1}{\varepsilon^2} \cdot \frac{\delta'}{k} \geq \frac{1}{\varepsilon^2} \cdot \Delta_k(\mathcal{X}_V).$$

Proof of Lemma 4.1

Let (P_1, \dots, P_k) and (Z_1, \dots, Z_{k-1}) be partitions of V . We define a mapping $\sigma : [1 : k-1] \mapsto [1 : k]$ by

$$\sigma(i) = \arg \max_{j \in [1:k]} \frac{\mu(Z_i \cap P_j)}{\mu(P_j)}, \quad \text{for every } i \in [1 : k-1].$$

We lower bound now the clusters overlapping in terms of the volume between any k -way and $(k-1)$ -way partitions of V .

► **Lemma 4.2.** *Suppose (P_1, \dots, P_k) and (Z_1, \dots, Z_{k-1}) are partitions of V . Then for any index $\ell \in [1 : k] \setminus \{\sigma(1), \dots, \sigma(k-1)\}$ (there is at least one such ℓ) and for every $i \in [1 : k-1]$ it holds*

$$\{\mu(Z_i \cap P_{\sigma(i)}), \mu(Z_i \cap P_\ell)\} \geq \tau_i \cdot \min \{\mu(P_\ell), \mu(P_{\sigma(i)})\},$$

where $\sum_{i=1}^{k-1} \tau_i = 1$ and $\tau_i \geq 0$.

Proof. By pigeonhole principle there is an index $\ell \in [1 : k]$ such that $\ell \notin \{\sigma(1), \dots, \sigma(k-1)\}$. Thus, for every $i \in [1 : k-1]$ we have $\sigma(i) \neq \ell$ and

$$\frac{\mu(Z_i \cap P_{\sigma(i)})}{\mu(P_{\sigma(i)})} \geq \frac{\mu(Z_i \cap P_\ell)}{\mu(P_\ell)} \triangleq \tau_i,$$

where $\sum_{i=1}^{k-1} \tau_i = 1$ and $\tau_i \geq 0$ for all i . Hence, the statement follows. ◀

Proof of Lemma 4.1. Let (Z_1, \dots, Z_{k-1}) be a $(k-1)$ -way partition of V with centers c'_1, \dots, c'_{k-1} that achieves $\Delta_{k-1}(\mathcal{X}_V)$, and (P_1, \dots, P_k) be a k -way partition of V achieving $\widehat{\rho}_{\text{avr}}(k)$. Our goal now is to lower bound the optimal $(k-1)$ -means cost

$$\Delta_{k-1}(\mathcal{X}_V) = \sum_{i=1}^{k-1} \sum_{j=1}^k \sum_{u \in Z_i \cap P_j} d_u \|F(u) - c'_i\|^2. \quad (15)$$

By Lemma 4.2 there is an index $\ell \in [1 : k] \setminus \{\sigma(1), \dots, \sigma(k-1)\}$. For $i \in [1 : k-1]$ let

$$p^{\gamma(i)} = \begin{cases} p^\ell & , \text{ if } \|p^\ell - c'_i\| \geq \|p^{\sigma(i)} - c'_i\|; \\ p^{\sigma(i)} & , \text{ otherwise.} \end{cases}$$

Then by combining Lemma 2.3 and Lemma 4.2, we have

$$\|p^{\gamma(i)} - c'_i\|^2 \geq [8 \cdot \min \{\mu(P_\ell), \mu(P_{\sigma(i)})\}]^{-1} \quad \text{and} \quad \mu(Z_i \cap P_{\gamma(i)}) \geq \tau_i \cdot \min \{\mu(P_\ell), \mu(P_{\sigma(i)})\}, \quad (16)$$

where $\sum_{i=1}^{k-1} \tau_i = 1$. We now lower bound the expression in (15). Since

$$\|F(u) - c'_i\|^2 \geq \frac{1}{2} \left(\|p^{\gamma(i)} - c'_i\|^2 - \|F(u) - p^{\gamma(i)}\|^2 \right),$$

it follows for $\delta' = 2\delta/20^4$ that

$$\begin{aligned}
 \Delta_{k-1}(\mathcal{X}_V) &= \sum_{i=1}^{k-1} \sum_{j=1}^k \sum_{u \in Z_i \cap P_j} d_u \|F(u) - c'_i\|^2 \geq \sum_{i=1}^{k-1} \sum_{u \in Z_i \cap P_{\gamma(i)}} d_u \|F(u) - c'_i\|^2 \\
 &\geq \frac{1}{2} \sum_{i=1}^{k-1} \sum_{u \in Z_i \cap P_{\gamma(i)}} d_u \|p^{\gamma(i)} - c'_i\|^2 - \sum_{i=1}^{k-1} \sum_{u \in Z_i \cap P_{\gamma(i)}} d_u \|F(u) - p^{\gamma(i)}\|^2 \\
 &\geq \frac{1}{2} \sum_{i=1}^{k-1} \frac{\mu(Z_i \cap P_{\gamma(i)})}{8 \cdot \min\{\mu(P_{\gamma(i)}), \mu(P_{\sigma(i)})\}} - \sum_{i=1}^k \sum_{u \in P_i} d_u \|F(u) - p^i\|^2 \\
 &\geq \frac{1}{16} - \frac{\delta'}{k},
 \end{aligned}$$

where the last inequality holds due to (16) and Lemma 3.1. \blacktriangleleft

5 An Efficient Spectral Clustering Algorithm

Here, we apply the proof techniques developed by Boutsidis et al. [1, 2] to our setting. More precisely, we prove that any α -approximate k -means algorithm that runs on an approximate normalized spectral embedding \widetilde{Y}' computed by the power method, yields an approximate clustering \widetilde{X}'_α of the normalized spectral embedding Y' .

Furthermore, we prove under our gap assumption that \widetilde{Y}' is ε -separated. This allows us to apply the variant of Lloyd's k -means algorithm analyzed by Ostrovsky et al. [5] to efficiently compute \widetilde{X}'_α . Then we use Theorem 1.2 (a) to establish the desired statement.

This section is organized as follows. In Subsection 5.1, we prove Theorem 2.8. Then in Subsection 5.2, we present the proof of Theorem 2.9. Based on the results from the preceding two subsections, we prove Theorem 1.2 (b) in Subsection 5.3.

5.1 Proof of Theorem 2.8

Due to space limits, we defer the proofs of the next Lemmas to the full version of the paper.

► **Lemma 5.1.** $X'X'^T$ is a projection matrix.

► **Lemma 5.2.** It holds that $Y'^TY' = I_{k \times k} = \widetilde{Y}'^T \widetilde{Y}'$.

► **Lemma 5.3.** It holds that $\|Y'Y'^T - \widetilde{Y}'\widetilde{Y}'^T\|_F = \|YY^T - \widetilde{Y}\widetilde{Y}^T\|_F$.

► **Lemma 5.4.** For any matrix U with orthonormal columns and every matrix A it holds

$$\|UU^T - AA^TUU^T\|_F = \|U - AA^TU\|_F. \tag{17}$$

Proof Sketch of Theorem 2.8. By combining Lemma 2.7 and Lemma 5.3, with probability at least $1 - 2e^{-2n} - 3\delta_p$ we have $\|Y'Y'^T - \widetilde{Y}'\widetilde{Y}'^T\|_F = \|YY^T - \widetilde{Y}\widetilde{Y}^T\|_F \leq \varepsilon$.

Let $Y'Y'^T = \widetilde{Y}'\widetilde{Y}'^T + E$ such that $\|E\|_F \leq \varepsilon$. Based on Lemma 5.2 and Lemma 5.4 we have that (17) holds for the matrices Y' and \widetilde{Y}' . Hence, by Lemma 5.1 we can apply the proof in [1, Theorem 6] to obtain $\|Y' - \widetilde{X}'_\alpha (\widetilde{X}'_\alpha)^T Y'\|_F \leq \sqrt{\alpha} \cdot (\|Y' - X'_{\text{opt}} (X'_{\text{opt}})^T Y'\|_F + 2\varepsilon)$. The desired statement follows by simple algebraic manipulations. \blacktriangleleft

5.2 Proof of Theorem 2.9

In this subsection, we show under our gap assumption that the approximate normalized spectral embedding \widetilde{Y}' is ε -separated, i.e. $\Delta_k(\widetilde{\mathcal{X}}_V) < 5\varepsilon^2 \cdot \Delta_{k-1}(\widetilde{\mathcal{X}}_V)$. Our analysis builds upon Theorem 2.6, Theorem 2.8 and the proof techniques in [1, Theorem 6].

We use interchangeably X'_{opt} and $X'^{(k)}_{\text{opt}}$ to denote the optimal indicator matrix for the k -means problem on \mathcal{X}_V that is induced by the rows of matrix Y' . Similarly, we denote by $X'^{(k-1)}_{\text{opt}}$ the optimal indicator matrix for the $(k-1)$ -means problem on \mathcal{X}_V .

Proof Sketch of Theorem 2.9. By Theorem 2.6 we have

$$\left\| Y' - X'_{\text{opt}} \left(X'_{\text{opt}} \right)^{\text{T}} Y' \right\|_F \leq \varepsilon \left\| Y' - X'^{(k-1)}_{\text{opt}} \left(X'^{(k-1)}_{\text{opt}} \right)^{\text{T}} Y' \right\|_F. \quad (18)$$

We set the approximation parameter in Theorem 2.8 to

$$\varepsilon' \triangleq \frac{1}{4} \sqrt{\Delta_k(\mathcal{X}_V)} = \frac{1}{4} \left\| Y' - X'_{\text{opt}} \left(X'_{\text{opt}} \right)^{\text{T}} Y' \right\|_F \geq n^{-O(1)}, \quad (19)$$

and we note that by Theorem 2.6 it holds $\varepsilon' \leq \frac{\varepsilon}{4} \sqrt{\Delta_{k-1}(\mathcal{X}_V)}$.

We construct now matrix \widetilde{Y} via the power method with $p \geq \Omega\left(\frac{\ln n}{\lambda_{k+1}}\right)$. By Lemma 5.3 we have $\left\| Y' Y'^{\text{T}} - \widetilde{Y} \widetilde{Y}'^{\text{T}} \right\|_F = \left\| Y Y^{\text{T}} - \widetilde{Y} \widetilde{Y}^{\text{T}} \right\|_F$ and thus by Lemma 2.7 with high probability it holds that $\left\| Y' (Y')^{\text{T}} - \widetilde{Y}' \widetilde{Y}'^{\text{T}} \right\|_F \leq \varepsilon'$.

Let $Y' (Y')^{\text{T}} = \widetilde{Y}' \widetilde{Y}'^{\text{T}} + E$ such that $\|E\|_F \leq \varepsilon'$. By combining Lemma 5.1, Lemma 5.2, Lemma 5.3 and by applying the proof techniques in [1, Theorem 6] we obtain

$$\sqrt{\Delta_k(\widetilde{\mathcal{X}}_V)} \leq \|E\|_F + \left\| Y' - \widetilde{X}'_{\text{opt}} \left(\widetilde{X}'_{\text{opt}} \right)^{\text{T}} Y' \right\|_F.$$

Furthermore, we can show that $\ln\left(\frac{2-\lambda_k}{2-\lambda_{k+1}}\right) \geq \frac{1}{2} \left(1 - \frac{4\delta}{20^4 k^2}\right) \lambda_{k+1}$. Then Theorem 2.8 yields

$$\left\| Y' - \widetilde{X}'_{\text{opt}} \left(\widetilde{X}'_{\text{opt}} \right)^{\text{T}} Y' \right\|_F^2 \leq (1 + 4\varepsilon') \cdot \left\| Y' - X'^{(k)}_{\text{opt}} \left(X'^{(k)}_{\text{opt}} \right)^{\text{T}} Y' \right\|_F^2 + 4\varepsilon'^2.$$

Also, we can show $\left\| Y' - X'^{(k)}_{\text{opt}} \left(X'^{(k)}_{\text{opt}} \right)^{\text{T}} Y' \right\|_F^2 \leq \frac{1}{8 \cdot 10^{13}}$ and by the definition of ε' it follows

$$\sqrt{\Delta_k(\widetilde{\mathcal{X}}_V)} \leq 2\sqrt{\Delta_k(\mathcal{X}_V)} \leq 2\varepsilon \cdot \sqrt{\Delta_{k-1}(\mathcal{X}_V)}. \quad (20)$$

Moreover, by applying similar arguments as in the proof of [1, Theorem 6] we can prove that

$$\sqrt{\Delta_{k-1}(\mathcal{X}_V)} \leq \left(1 + \frac{\varepsilon}{2}\right) \sqrt{\Delta_{k-1}(\widetilde{\mathcal{X}}_V)}. \quad (21)$$

The statement follows by combining (20) and (21). \blacktriangleleft

5.3 Proof of Part (b) of Theorem 1.2

Let $p = \Theta(\frac{\ln n}{\lambda_{k+1}})$. We can compute the matrix $\mathcal{B}^p S$ in time $O(mkp)$ and its singular value decomposition $\widetilde{U}\widetilde{\Sigma}\widetilde{V}^T$ in time $O(nk^2)$. Based on it, we construct in time $O(mk)$ matrix \widetilde{Y}' (c.f. (13)).

By Theorem 2.9, $\widetilde{\mathcal{X}}_V$ is ε -separated for $\varepsilon = 6/10^7$, i.e. $\Delta_k(\widetilde{\mathcal{X}}_V) < 5\varepsilon^2 \cdot \Delta_{k-1}(\widetilde{\mathcal{X}}_V)$. Hence, by Theorem 2.5 there is an algorithm that outputs in time $O(mk^2 + k^4)$ a clustering with indicator matrix \widetilde{X}'_α satisfying

$$\left\| \widetilde{Y}' - \widetilde{X}'_\alpha (\widetilde{X}'_\alpha)^T \widetilde{Y}' \right\|_F^2 \leq \left(1 + \frac{1}{10^{10}} \right) \cdot \left\| \widetilde{Y}' - \widetilde{X}'_{\text{opt}} (\widetilde{X}'_{\text{opt}})^T \widetilde{Y}' \right\|_F^2$$

with constant probability (close to 1), where $\alpha = 1 + 1/10^{10}$.

Moreover, by applying Theorem 2.8 with $\varepsilon' = \frac{1}{4 \cdot 10^8} \left\| Y' - X'_{\text{opt}} (X'_{\text{opt}})^T Y' \right\|_F$ we can prove that the indicator matrix \widetilde{X}'_α yields a multiplicative approximation of \mathcal{X}_V , i.e.

$$\left\| Y' - \widetilde{X}'_\alpha (\widetilde{X}'_\alpha)^T Y' \right\|_F^2 \leq \left(1 + \frac{1}{10^6} \right) \left\| Y' - X'_{\text{opt}} (X'_{\text{opt}})^T Y' \right\|_F^2. \quad (22)$$

The statement follows by Theorem 1.2 (a) applied to the partition (A_1, \dots, A_k) of V that is induced by the indicator matrix \widetilde{X}'_α .

Acknowledgement. We would like to thank the anonymous reviewers for their constructive comments.

References

- 1 Christos Boutsidis, Prabhajan Kambadur, and Alex Gittens. Spectral clustering via the power method - provably. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 40–48, 2015.
- 2 Christos Boutsidis and Malik Magdon-Ismael. Faster svd-truncated regularized least-squares. In *2014 IEEE International Symposium on Information Theory, Honolulu, HI, USA, June 29 - July 4, 2014*, pages 1321–1325, 2014.
- 3 James R. Lee, Shayan Oveis Gharan, and Luca Trevisan. Multi-way spectral partitioning and higher-order Cheeger inequalities. In *Proceedings of the Forty-fourth Annual ACM Symposium on Theory of Computing, STOC'12*, pages 1117–1130, New York, NY, USA, 2012. ACM.
- 4 Andrew Y. Ng, Michael I. Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in Neural Information Processing Systems 14*, pages 849–856. MIT Press, 2002.
- 5 Rafail Ostrovsky, Yuval Rabani, Leonard J. Schulman, and Chaitanya Swamy. The effectiveness of Lloyd-type methods for the k-means problem. *J. ACM*, 59(6):28:1–28:22, January 2013.
- 6 Shayan Oveis Gharan and Luca Trevisan. Partitioning into expanders. In *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2014, Portland, Oregon, USA, January 5-7, 2014*, pages 1256–1266, 2014.
- 7 Richard Peng, He Sun, and Luca Zannetti. Partitioning well-clustered graphs: Spectral clustering works! In *Proceedings of The 28th Conference on Learning Theory, COLT 2015, Paris, France, July 3-6, 2015*, pages 1423–1455, 2015.
- 8 Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.

57:14 A Note On Spectral Clustering

- 9 Ali Kemal Sinop. How to round subspaces: A new spectral clustering algorithm. In *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2016, Arlington, VA, USA, January 10-12, 2016*, pages 1832–1847, 2016.
- 10 Ulrike von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, pages 395–416, 2007.