

On Clusterings: Good, Bad and Spectral

RAVI KANNAN*

Department of Computer Science,
Yale University, New Haven.
Email: kannan@cs.yale.edu

SANTOSH VEMPALA[†] and ADRIAN VETTA[‡]

Department of Mathematics,
M.I.T., Cambridge.
Email: {vempala, avetta}@math.mit.edu

July 2001

ABSTRACT. We motivate and develop a new bicriteria measure for assessing the quality of a clustering which avoids the drawbacks of existing measures. A simple recursive heuristic is shown to have poly-logarithmic worst-case guarantees under the new measure. The main result of the paper is the analysis of a popular *spectral* algorithm. One variant of spectral clustering turns out to have effective worst-case guarantees; another finds a “good” clustering if it exists.

*Supported in part by NSF grant CCR-9820850.

[†]Supported by NSF CAREER award CCR-9875024.

[‡]Supported in part by NSF CAREER award CCR-9875024 and in part by a Wolfe fellowship.

1 Introduction

Clustering, or partitioning into dissimilar groups of similar items, is a problem with many variants in mathematics and the applied sciences. The availability of vast amounts of data has revitalized research on the problem. Over the years, several clever heuristics have been invented for clustering. While many of these heuristics are problem-specific, the method known as *spectral clustering* has been applied successfully in a variety of different situations. Roughly speaking, spectral clustering is the technique of partitioning the rows of a matrix according to their components in the top few singular vectors of the matrix (see section 1.1 for a detailed description). The main motivation of this paper was to analyze the performance of spectral clustering. However, such an evaluation is inextricably linked to the question of how to measure the quality of a clustering. The justification provided by practitioners is typically case-by-case and experimental (“it works well on my data”). Theoreticians, meanwhile, have been busy studying quality measures that are seductively simple to define (e.g. k -median, minimum sum, minimum diameter, etc.). The measures thus far analyzed by theoreticians are easy to fool, i.e. there are simple examples where the “right” clustering is obvious but optimizing these measures produces undesirable solutions (see section 2). Thus, neither approach has been entirely satisfactory.

In this paper we propose a new bicriteria measure of the quality of a clustering, based on expansion-like properties of the underlying pairwise similarity graph. The quality of a clustering is given by two parameters: α , the minimum conductance of the clusters and ϵ , the ratio of the weight of inter-cluster edges to the total weight of all edges. The objective is to find an (α, ϵ) clustering that maximizes α and minimizes ϵ . Note that the conductance provides a measure of the quality of an individual cluster (and thus of the overall clustering) whilst the weight of the inter-cluster edges provides a measure of the cost of the clustering. Hence, imposing a lower bound, α , on the quality of each individual cluster we strive to minimize the cost, ϵ , of the clustering; or conversely, imposing an upper bound on the cost of the clustering we strive to maximize its quality. In section 2, we motivate the use of this more complex, bicriteria measure by showing that it does not have the obvious drawbacks of the simpler quality measures.

While the new measure might be qualitatively attractive, it would be of little use if optimizing it were computationally intractable. In section 3 we study a recursive heuristic designed to optimize the new measure. Although finding an exact solution is NP-hard, the algorithm is shown to have simultaneous poly-logarithmic approximations guarantees for the two parameters in the bicriteria measure (corollary 2).

In section 4 we turn to spectral algorithms for clustering. These algorithms are attractive in part because of their speed (see section 5). However, their performance had hitherto eluded a rigorous analysis. The use of the new measure turns out to be conducive for such a theoretical analysis. In particular, we show that a recursive version of spectral clustering has effective worst-case approximation guarantees with respect to the bicriteria measure (corollary 4). It is worth noting that both our worst-case guarantees follow from the same general theorem (see theorem 1 in section 3). Another variant of spectral clustering has the following guarantee: if the input data has a rather good clustering (i.e. α is large and ϵ is small), then the spectral algorithm will find a clustering that is “close” to the optimal clustering (theorem 5).

1.1 Spectral Clustering Algorithms

Spectral clustering refers to the general technique of partitioning the rows of a matrix according to their components in the top few singular vectors of the matrix. The underlying problem, that of clustering the rows of a matrix, is ubiquitous. We mention three special cases that are all of independent interest:

- The matrix encodes the pairwise similarities of vertices of a graph.
- The rows of the matrix are points in a d -dimensional Euclidean space. The columns are the coordinates.
- The rows of the matrix are documents of a corpus. The columns are terms. The (i, j) entry encodes information about the occurrence of the j th term in the i th document.

Given a matrix A , the spectral algorithm for clustering the rows of A is given below.

Spectral Algorithm I

Find the top k right singular vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$.

Let C be the matrix whose j th column is given by $A\mathbf{v}_j$.

Place row i in cluster j if C_{ij} is the largest entry in the i th row of C .

The algorithm has the following interpretation². Suppose the rows of A are points in a high-dimensional space. Then the subspace defined by the top k right singular vectors of A is the rank k subspace that best approximates A . The spectral algorithm projects all the points onto this subspace. Each singular vector then defines a cluster; to obtain a clustering we map each projected point to the (cluster defined by the) singular vector that is closest to it in angle.

In section 4, we describe a recursive variant of this algorithm.

2 What is a Good Clustering?

How good is the spectral algorithm? Intuitively, a clustering algorithm performs well if points that are similar are assigned the same cluster and points that are dissimilar are assigned to different clusters. Of course, this may not be possible to do for every pair of points, and so we compare the clustering found by the algorithm to the *optimal* one for the given matrix. This, though, leads to another question: what exactly is an optimal clustering? To provide a quantitative answer, we first need to define a measure of the quality of a clustering. In recent years several combinatorial measures of clustering quality have been investigated in detail. These include *minimum diameter*, *k-center*, *k-median*, and *minimum sum* (for example: [3], [4], [6], [8], [9], etc.).

All these measures, although mathematically attractive due to their simplicity, are easy to fool. That is, one can construct examples with the property that the “best” clustering is obvious and yet an algorithm that optimizes one of these measures finds a clustering that is substantially different (and therefore unsatisfactory). Such examples are presented in Figures 1 and 2, where the goal is to partition the points into two clusters. Observe that all the

²Computationally it is useful to note that the j th column of C is also given by $\lambda_j \mathbf{u}_j$; here λ_j is the j th singular vector of A and \mathbf{u}_j is the j th left singular vector.

measures given above seek to minimize some objective function. In the figures, nearby points (which represent highly similar points) induce low cost edges; points that are farther apart (and represent dissimilar points) induce high cost edges.

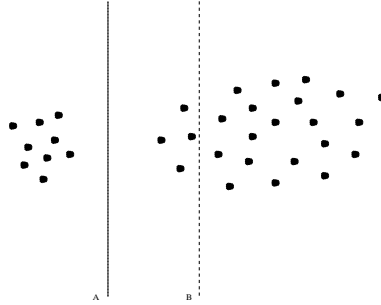


Figure 1: Optimizing the diameter produces B while A is clearly more desirable.

Consider a clustering that minimizes the maximum diameter of the clusters; the diameter of a cluster being the largest distance, say, between two points in a cluster. It is NP-hard to find such a clustering, but this is not our main concern. What is worrisome about the example shown in Figure 1 is that the optimal solution (B) produces a cluster which contains points that should have been separated. Clustering with respect to the minimum sum and k -center measures will produce the same result. The reason such a poor cluster is produced is that although we have minimized the maximum dissimilarity between points in a cluster, this was at the expense of creating a cluster with many dissimilar points. The clustering (A) on the other hand, although it leads to a larger maximum diameter, say, is desirable since it better satisfies the goal of “similar points together and dissimilar points apart”. This problem also arises for the k -median measure (see, for example, the case shown in Figure 2); it may produce clusters of poor quality.

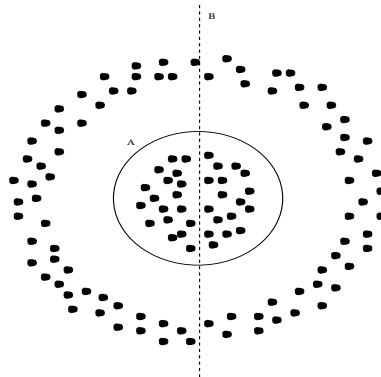


Figure 2: The inferior clustering B is found by optimizing the 2-median measure.

We will find it more convenient to use model the input as a similarity graph rather than as a distance graph. This is indeed often the case in practice. Thus the input is an edge-weighted complete graph whose vertices need to be partitioned. The weight of an edge a_{ij} represents

the similarity of the vertices (points) i and j . Thus, the graph for points in space would have high edge weights for points that are close together and low edge weights for points that are far apart. So the graph is associated with an $n \times n$ symmetric matrix A with entries a_{ij} ; here we assume that the a_{ij} are non-negative.

Let us now return to the question of what a good clustering is. The quality of a cluster should be determined by how similar the points within a cluster are. Note that each cluster is represented by a subgraph. In particular, if there is a cut of small weight that divides the cluster into two pieces of comparable size then the cluster has lots of pairs of vertices that are dissimilar and hence it is of low quality. This might suggest that the quality of a subgraph as a cluster is the minimum cut of the subgraph. However, this is misleading as is illustrated by Figure 3.

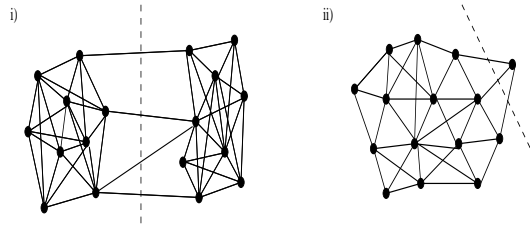


Figure 3: The second subgraph is of higher quality as a cluster even though it has a smaller minimum cut.

In this example edges represent high-similarity pairs and non-edges represent pairs that are highly dissimilar. The minimum cut of the first subgraph is larger than that of the second subgraph. This is because the second subgraph has low degree vertices. However, the second subgraph is a higher quality cluster. This can be attributed to the fact that in the first subgraph there is a cut whose weight is small *relative to the sizes of the pieces it creates*. A quantity that measures the relative cut size is the *expansion*. The expansion of a graph is the minimum ratio over all cuts of the graph of the total weight of edges of the cut to the number or vertices in the smaller part created by the cut. Formally, we denote the expansion of a cut (S, \bar{S}) by:

$$\psi(S) = \frac{\sum_{i \in S, j \notin S} a_{ij}}{\min(|S|, |\bar{S}|)}$$

We say that the expansion of a graph is the minimum expansion over all the cuts of the graph. Our first measure of quality of a cluster is the expansion of the subgraph corresponding to it. The expansion of a clustering is the minimum expansion of one of the clusters.

The measure defined above gives equal importance to all the vertices of the given graph. This, though, may lead to a rather taxing requirement. For example, in order to accommodate a vertex i with very little similarity to all the other vertices combined (i.e., $\sum_j a_{ij}$ is small), then α will have to be very low. Arguably, it is more prudent to give greater importance to vertices that have many similar neighbors and lesser importance to vertices that have few similar neighbors. This can be done by a direct generalization of the expansion, called the *conductance*, in which subsets of vertices are weighted to reflect their importance.

The conductance of a cut (S, \bar{S}) in G is denoted by:

$$\phi(S) = \frac{\sum_{i \in S, j \notin S} a_{ij}}{\min(a(S), a(\bar{S}))}$$

Here $a(S) = a(S, V) = \sum_{i \in S} \sum_{j \in V} a_{ij}$. The conductance of a graph is the minimum conductance over all the cuts of the graph; $\phi(G) = \min_{S \subseteq V} \phi(S)$. In order to quantify the quality of a clustering we generalize the definition of conductance further. Take a cluster $C \subseteq V$ and a cut $(S, C \setminus S)$ within C , where $S \subseteq C$. Then we say that the *conductance* of S in C is:

$$\phi(S, C) = \frac{\sum_{i \in S, j \in C \setminus S} a_{ij}}{\min(a(S), a(C \setminus S))}$$

The conductance of a cluster $\phi(C)$ will then be the smallest conductance of a cut within the cluster. The conductance of a clustering is the minimum conductance of its clusters. This conductance measure seems extremely well suited to achieve our intuitive goal i.e. clustering similar points and separating dissimilar points. We then obtain the following optimization problem: given a graph and an integer k , find a k -clustering with the maximum conductance. Notice that optimizing the expansion/conductance gives the right clustering in the examples of Figures 1 and 2. To see this assume, for example, that the points induce an unweighted graph (i.e. zero-one edge weights). Thus, a pair of vertices induces an edge if and only if the two vertices are close together. Clustering (A) will then be obtained in each example.

There is still a problem with the above clustering measure. The graph might consist mostly of clusters of high quality and maybe a few points that create clusters of very poor quality, so that any clustering necessarily has a poor overall quality (since we have defined the quality of a clustering to be the minimum over all the clusters). In fact, to boost the overall quality, the best clustering might create many clusters of relatively low quality so that the minimum is as large as possible. Such an example is shown in Figure 4.

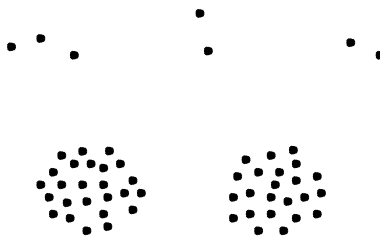


Figure 4: Assigning the outliers leads to poor quality clusters.

One way to handle the problem might be to avoid restricting the number of clusters. But this could lead to a situation where many points are in singleton (or extremely small) clusters. Instead we measure the quality of a clustering using two criteria, the first is the minimum quality of the clusters (called α), and the second is the fraction of the total weight of edges that are not covered by the clusters (called ϵ).

Definition 1. We call a partition $\{C_1, C_2, \dots, C_l\}$ of V an (α, ϵ) -partition if :

1. The conductance of each C_i is at least α .
2. The total weight of inter-cluster edges is at most an ϵ fraction of the total edge weight.

Thus we obtain a bicriteria measure of the quality of a clustering. Associated with this bicriteria measure is the following optimization problem (note that the number of clusters is not restricted).

Problem 1. *Given α find an (α, ϵ) -partition that minimizes ϵ (alternatively, given ϵ find an (α, ϵ) -partition that maximizes α).*

For this measure of cluster quality, we are not aware of any “bad” examples for the (α, ϵ) -clustering problem, i.e. examples where solving the bicriteria optimization problem gives a clustering that is clearly inferior to the best clustering. An important question is whether this observation holds up under a systematic experimental study. The focus of the rest of this paper, however, is to consider the measure from a theoretical standpoint and to examine in detail the performance of spectral clustering algorithms.

It may be noted that there is a monotonic function f that represents the optimal (α, ϵ) pairings. For example, for each α there is a minimum value of ϵ , equal to $f(\alpha)$, such that an (α, ϵ) -partition exists. In the following sections we present two approximation algorithms for the clustering property. One nice characteristic of these algorithms is that in a single application they can be used to obtain an approximation f' for the entire function f , not just for f evaluated at a single point. Thus the user need not specify a desired value of α or ϵ a priori. Rather, the desired conductance/cost trade-off may be determined after consideration of the approximation function f' .

3 Approximation Algorithms

Problem 1 is NP-hard. To see this, consider maximizing α whilst setting ϵ to zero. This problem is equivalent to finding the conductance of a given graph. Here we present a simple heuristic and provide worst-case approximation guarantees for it.

Approximate-Cluster Algorithm

Find a cut that approximates the minimum conductance cut in G .
Recurse on the pieces induced by the cut.

The idea behind our algorithm is simple. Given G , find a cut (S, \bar{S}) of minimum conductance. Then recurse on the subgraphs induced by S and \bar{S} . Finding a cut of minimum conductance is hard, and hence we need to use an approximately minimum cut. There are two well-known approximations for the minimum conductance cut, one is based on a linear programming relaxation and the other is derived from the second eigenvector of the graph. Before we discuss these approximations, we prove a general theorem for general approximation heuristics.

Let A be an approximation algorithm that produces a cut of conductance at most Kx^ν if the minimum conductance is x , where K is independent of x (K could be a function of n , for example) and ν is a fixed constant between 0 and 1. The following theorem provides a guarantee for the approximate-cluster algorithm using A as a subroutine.

Theorem 1. *If G has an (α, ϵ) -partition, then the approximate-cluster algorithm will find a partition of quality*

$$\left(\frac{\alpha}{6K \log \frac{n}{\epsilon}} \right)^{1/\nu}, (12K + 2)\epsilon^\nu \log \frac{n}{\epsilon}.$$

Proof. Let the cuts produced by the algorithm be $(S_1, T_1), (S_2, T_2), \dots$, where we adopt the convention that S_j is the “smaller” side (i.e., $a(S_j) \leq a(T_j)$). Let C_1, C_2, \dots, C_l be the optimal (α, ϵ) -clustering. We use the termination condition of $\alpha^* = \frac{\alpha}{6 \log \frac{n}{\epsilon}}$. We will assume that we apply the recursive step in the algorithm only if the conductance of a given piece as detected by the heuristic for the minimum conductance cut is less than α^* . In addition, purely for the sake of analysis we consider a slightly modified algorithm. If at any point we have a cluster C_t with the property that $a(C_t) < \frac{\epsilon}{n} a(V)$ then we split C_t into singletons. Clearly, upon termination, each cluster has conductance at least

$$\left(\frac{\alpha^*}{K}\right)^{1/\nu} = \left(\frac{\alpha}{6K \log \frac{n}{\epsilon}}\right)^{1/\nu}.$$

Thus it remains to bound the weight of the inter-cluster edges. Observe that $a(V)$ is twice the total edge weight in the graph, and so $W = \frac{\epsilon}{2} a(V)$ is the weight of the inter-cluster edges in this optimal solution.

Now we divide the cuts into two groups; the first group H are the ones with “high” conductance within clusters. The second group consists of the remaining cuts. We will use the notation $w(S_j, T_j) = \sum_{u \in S_j, v \in T_j} a_{uv}$. In addition, we denote by $w_I(S_j, T_j)$ the sum of the weights of the intra-cluster edges of the cut (S_j, T_j) , i.e., $w_I(S_j, T_j) = \sum_{i=1}^l w(S_j \cap C_i, T_j \cap C_i)$.

$$H = \{j : w_I(S_j, T_j) \geq 2\alpha^* \sum_{i=1}^l \min(a(S_j \cap C_i), a(T_j \cap C_i))\}$$

We now bound the cost of the high conductance group. For all $j \in H$, we have,

$$\alpha^* a(S_j) \geq w(S_j, T_j) \geq w_I(S_j, T_j) \geq 2\alpha^* \sum_i \min(a(S_j \cap C_i), a(T_j \cap C_i))$$

Consequently we observe that

$$\sum_i \min(a(S_j \cap C_i), a(T_j \cap C_i)) \leq \frac{1}{2} a(S_j)$$

For each $j \in H$, we define a cluster-avoiding cut (S'_j, T'_j) in $S_j \cup T_j$ in the following manner. For each $i, 1 \leq i \leq l$, if $a(S_j \cap C_i) \geq a(T_j \cap C_i)$, then place all of $(S_j \cup T_j) \cap C_i$ into S'_j . If $a(S_j \cap C_i) < a(T_j \cap C_i)$, then place all of $(S_j \cup T_j) \cap C_i$ into T'_j . An example is given in Figure 5, where the original cut is shown by the solid line and the cluster-avoiding cut by the dashed line. Notice that, since $|a(S_j) - a(S'_j)| \leq \frac{1}{2} a(S_j)$, we have that $\min(a(S'_j), a(T'_j)) \geq \frac{1}{2} a(S_j)$. Now we will use the approximation guarantee for the cut procedure to get an upper bound on $w(S_j, T_j)$ in terms of $w(S'_j, T'_j)$.

$$\begin{aligned} \frac{w(S_j, T_j)}{a(S_j)} &\leq K \left(\frac{w(S'_j, T'_j)}{\min\{a(S'_j), a(T'_j)\}} \right)^\nu \\ &\leq K \left(\frac{2w(S'_j, T'_j)}{a(S_j)} \right)^\nu \end{aligned}$$

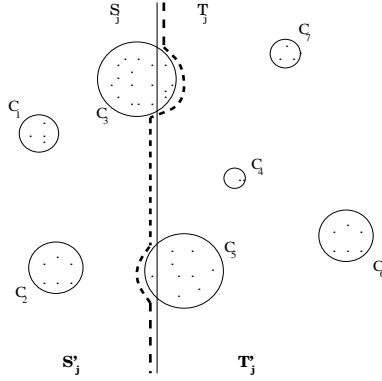


Figure 5:

Hence we have bounded the overall cost of the high conductance cuts with respect to the cost of the cluster-avoiding cuts. We now bound the cost of these cluster-avoiding cuts. Let $P(S)$ denote the set of inter-cluster edges incident at a vertex in S , for any subset S of V . Then $w(S'_j, T'_j) \leq w(P(S'_j))$, since every edge in (S'_j, T'_j) is an inter-cluster edge. So we have,

$$w(S_j, T_j) \leq K(2w(P(S'_j)))^\nu a(S_j)^{1-\nu} \quad (1)$$

Next we prove the following claim.

Claim 2. For each vertex $u \in V$, there are at most $\log \frac{n}{\epsilon}$ values of j such that u belongs to S_j . Further, there are at most $2 \log \frac{n}{\epsilon}$ values of j such that u belongs to S'_j .

To prove the claim, fix a vertex $u \in V$. Let

$$I = \{j : u \in S_j\} \quad J = \{j : u \in S'_j \setminus S_j\}$$

Clearly if $u \in S_j, S_k$ (with $k > j$), then (S_k, T_k) must be a partition of S_j or a subset of S_j . Now we have, $a(S_k) \leq \frac{1}{2}a(S_k \cup T_k) \leq \frac{1}{2}a(S_j)$. So $a(S_j)$ reduces by a factor of 2 or greater between two successive times u belongs to S_j . The maximum value of $a(S_j)$ is at most $a(V)$ and the minimum value is at least $\frac{\epsilon}{n}a(V)$, so the first statement of the claim follows.

Now suppose $j, k \in J; j < k$. Suppose also $u \in C_i$. Then $u \in T_j \cap C_i$. Also later T_j (or a subset of T_j) is being partitioned into (S_k, T_k) and since $u \in S'_k \setminus S_k$, we have $a(T_k \cap C_i) \leq a(S_k \cap C_i)$. Thus $a(T_k \cap C_i) \leq \frac{1}{2}a(S_k \cup T_k) \leq \frac{1}{2}a(T_j \cap C_i)$. Thus $a(T_j \cap C_i)$ halves between two successive times that $j \in J$. So, $|J| \leq \log \frac{n}{\epsilon}$. This proves the second statement in the claim (since $u \in S'_j$ implies that $u \in S_j$ or $u \in S'_j \setminus S_j$). These concepts are shown pictorially in Figure 6, where the cuts (S_j, T_j) and (S_k, T_k) are represented by solid lines and the cuts (S'_j, T'_j) and (S'_k, T'_k) by dashed lines.

Using this claim we can bound the overall cost of the group of cuts with high conductance within clusters with respect to the cost of the optimal clustering as follows:

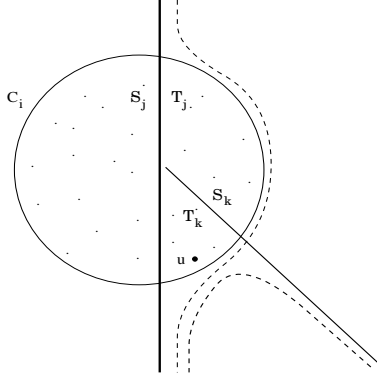


Figure 6:

$$\begin{aligned}
\sum_{j \in H} w(S_j, T_j) &\leq \sum_{\text{all } j} K(2w(P(S'_j)))^\nu a(S_j)^{1-\nu} \\
&\leq K \left(2 \sum_{\text{all } j} w(P(S'_j)) \right)^\nu \left(\sum_j a(S_j) \right)^{1-\nu} \\
&\leq K \left(2\epsilon \log \frac{n}{\epsilon} a(V) \right)^\nu \left(2 \log \frac{n}{\epsilon} a(V) \right)^{1-\nu} \\
&\leq 2K\epsilon^\nu \log \frac{n}{\epsilon} a(V)
\end{aligned} \tag{2}$$

Here we used Holder's inequality. Next we deal with the group of cuts with low conductance within clusters i.e. those j not in H . First, suppose that all the cuts together induce a partition of C_i into $P_1^i, P_2^i, \dots, P_{r_i}^i$. Every edge between two vertices in C_i which belong to different sets of the partition must be cut by some cut (S_j, T_j) and conversely, every edge of every cut $(S_j \cap C_i, T_j \cap C_i)$ must have its two end points in different sets of the partition. So, given that C_i has α conductance, we obtain

$$\sum_{\text{all } j} w_I(S_j \cap C_i, T_j \cap C_i) = \frac{1}{2} \sum_{s=1}^{r_i} w(P_s^i, C_i \setminus P_s^i) \geq \frac{1}{2} \alpha \sum_s \min(a(P_s^i), a(C_i \setminus P_s^i))$$

For each vertex $u \in C_i$ there can be at most $\log \frac{n}{\epsilon}$ values of j such that u belongs to the smaller (according to $a(\cdot)$) of the two sets $S_j \cap C_i, T_j \cap C_i$. So, we have that

$$\sum_{s=1}^{r_i} \min(a(P_s^i), a(C_i \setminus P_s^i)) \geq \frac{1}{\log \frac{n}{\epsilon}} \sum_j \min(a(S_j \cap C_i), a(T_j \cap C_i))$$

Thus,

$$\sum_{\text{all } j} w_I(S_j, T_j) \geq \frac{\alpha}{2 \log \frac{n}{\epsilon}} \sum_{i=1}^l \sum_j \min(a(S_j \cap C_i), a(T_j \cap C_i))$$

Therefore, from the definition of H , we have

$$\sum_{j \notin H} w_I(S_j, T_j) \leq 2\alpha^* \sum_{\text{all } j} \sum_{i=1}^l \min(a(S_j \cap C_i), a(T_j \cap C_i)) \leq \frac{2}{3} \sum_{\text{all } j} w_I(S_j, T_j)$$

Thus, we are able to bound the intra-cluster cost of the low conductance group of cuts in terms of the intra-cluster cost of the high conductance group. Applying (2) then gives

$$\sum_{j \notin H} w_I(S_j, T_j) \leq 2 \sum_{j \in H} w_I(S_j, T_j) \leq 4K\epsilon^\nu \log \frac{n}{\epsilon} a(V) \quad (3)$$

In addition, since each inter-cluster edge belongs to at most one cut S_j, T_j , we have that

$$\sum_{j \notin H} (w(S_j, T_j) - w_I(S_j, T_j)) \leq \frac{\epsilon}{2} a(V) \quad (4)$$

We then sum up (2), (3) and (4). To get the total cost we note that splitting up all the V_t with $a(V_t) \leq \frac{\epsilon}{n} a(V)$ into singletons costs us at most $\frac{\epsilon}{2} a(V)$ on the whole. Substituting $a(V)$ as twice the total sum of edge weights gives the bound on the cost of inter-cluster edge weights follows. \square

The Leighton-Rao algorithm for approximating the conductance finds a cut of conductance at most $2 \log n$ times the minimum [10]. In our terminology, it is an approximation algorithm with $K = 2 \log n$ and $\nu = 1$. Applying theorem 1 leads to the following guarantee.

Corollary 2. *Using the Leighton-Rao heuristic, the approximate-cluster algorithm finds an*

$$\left(\frac{\alpha}{12 \log n \log \frac{n}{\epsilon}}, 26\epsilon \log n \log \frac{n}{\epsilon} \right) \text{ bicriteria approximation.}$$

We now assess the running time of the algorithm using this heuristic. The fastest implementation for this heuristic, due to Benczur and Karger [2], runs in $\tilde{O}(n^2)$ time. Since the algorithm makes less than n cuts, the total running time is $\tilde{O}(n^3)$. This might be too slow for some real-world applications. We discuss a potentially more practical algorithm in the next section.

4 Performance Guarantees for Spectral Clustering

In this section we describe and analyse a recursive variant of the spectral algorithm. This algorithm, outlined below, has been used in the field of computer vision [12] and also in the field of web search engines [16]. Note that the algorithm is a special case of the approximate-cluster algorithm described in the previous section; here we use a spectral heuristic to approximate the minimum conductance cut.

Spectral Algorithm II

Normalize A and find its 2nd right eigenvector \mathbf{v} .
 Find the best ratio cut wrt \mathbf{v} .
 Recurse on the pieces induced by the cut.

Observe, also, that upon normalization of the matrix our conductance measure corresponds to the familiar Markov Chain conductance measure i.e.

$$\phi(S) = \frac{\sum_{i \in S, j \notin S} a_{ij}}{\min(a(S), a(\bar{S}))} = \frac{\sum_{i \in S, j \notin S} \pi_i b_{ij}}{\min(\pi(S), \pi(\bar{S}))}$$

where π is the stationary distribution of the Markov Chain.

We now elaborate upon the basic description of this variant of the spectral algorithm. Initially we normalize our matrix A so that the rows sums are equal to one. Now at any stage in the algorithm we have a clustering $\{C_1, C_2, \dots, C_s\}$. Then for each C_t we consider the $|C_t| \times |C_t|$ submatrix B of A restricted to C_t . We normalise B by setting b_{ii} to $1 - \sum_{j \in C_t, j \neq i} b_{ij}$. Hence B is also non-negative with row sums equal to one.

We then find the second eigenvector of B . This is the right eigenvector \mathbf{v} corresponding to the second largest eigenvalue λ_2 , i.e. $B\mathbf{v} = \lambda_2\mathbf{v}$. Then order the elements (rows) of C_t decreasingly with respect to their component in the direction of \mathbf{v} . Given this ordering, say $\{u_1, u_2, \dots, u_r\}$, find the minimum *ratio cut* in C_t . This is the cut that minimises $\phi(\{u_1, u_2, \dots, u_j\}, C_t)$ for some j , $1 \leq j \leq r - 1$. We then recurse on the pieces $\{u_1, \dots, u_j\}$ and $C_t \setminus \{u_1, \dots, u_j\}$.

4.1 Worst-case guarantees

We will use the following theorem to prove a worst-case guarantee for the algorithm. This result was essentially proved by Sinclair and Jerrum (in their proof of Lemma 3.3 in [13], although not mentioned in the statement of the lemma). For completeness, and due to the fact that theorem 3 is usually not explicitly stated in the Markov Chain literature (or usually includes some other conditions which are not relevant here), we include a proof of this result. Observe that, via the use of the second eigenvalue, the theorem bounds the conductance of the cut found by the heuristic with respect to that of the optimal cut.

Theorem 3. *Suppose B is a $N \times N$ matrix with non-negative entries with each row sum equal to 1 and suppose there are positive real numbers $\pi_1, \pi_2, \dots, \pi_N$ summing to 1 such that $\pi_i b_{ij} = \pi_j b_{ji}$ for all i, j . If \mathbf{v} the right eigenvector of B corresponding to the second largest eigenvalue λ_2 , and i_1, i_2, \dots, i_N is an ordering of $1, 2, \dots, N$ so that $v_{i_1} \geq v_{i_2} \dots \geq v_{i_N}$, then*

$$\min_{S \subseteq \{1, 2, \dots, N\}} \frac{\sum_{i \in S, j \notin S} \pi_i b_{ij}}{\min(\sum_{i \in S} \pi_i, \sum_{j \notin S} \pi_j)} \geq 1 - \lambda_2 \geq \frac{1}{2} \left(\min_{l, 1 \leq l \leq N} \frac{\sum_{1 \leq u \leq l, l+1 \leq v \leq N} \pi_{i_u} b_{i_u i_v}}{\min(\sum_{1 \leq u \leq l} \pi_{i_u}, \sum_{l+1 \leq v \leq N} \pi_{i_v})} \right)^2$$

Before proving this theorem, let us use it along with theorem 1 to get a worst-case guarantee for spectral algorithm II. In our terminology, the above theorem says that the spectral heuristic for minimum conductance is an approximation algorithm with $C = \sqrt{2}$ and $\nu = 1/2$.

Corollary 4. *Using the spectral heuristic, the approximate-cluster algorithm finds an*

$$\left(\frac{\alpha^2}{72 \log^2 \frac{n}{\epsilon}}, 20\sqrt{\epsilon} \log \frac{n}{\epsilon} \right) \text{ bicriteria approximation.}$$

Proof (of theorem 3). We first evaluate the second eigenvalue. Towards this end, let $D^2 = \text{diag}(\pi)$. Then, from the time-reversibility property of B , we have $D^2B = B^TD^2$. Hence $Q = DBD^{-1}$ is symmetric. The eigenvalues of B and Q are the same, with their largest eigenvalue equal to 1. In addition, $\pi^TD^{-1}Q = \pi^TD^{-1}$ and therefore π^TD^{-1} is the left eigenvector of Q corresponding to the eigenvalue 1. So we have,

$$\lambda_2 = \max_{\pi^TD^{-1}\mathbf{x}=0} \frac{\mathbf{x}^TDBD^{-1}\mathbf{x}}{\mathbf{x}^T\mathbf{x}}$$

Thus, substituting $\mathbf{y} = D^{-1}\mathbf{x}$, we obtain

$$1 - \lambda_2 = \min_{\pi^TD^{-1}\mathbf{x}=0} \frac{\mathbf{x}^TD(I-B)D^{-1}\mathbf{x}}{\mathbf{x}^T\mathbf{x}} = \min_{\pi^T\mathbf{y}=0} \frac{\mathbf{y}^TD^2(I-B)\mathbf{y}}{\mathbf{y}^TD^2\mathbf{y}}$$

The numerator can be rewritten:

$$\begin{aligned} \mathbf{y}^TD^2(I-B)\mathbf{y} &= -\sum_{i \neq j} y_i y_j \pi_i b_{ij} + \sum_i \pi_i (1 - b_{ii}) y_i^2 \\ &= -\sum_{i \neq j} y_i y_j \pi_i b_{ij} + \sum_{i \neq j} \pi_i b_{ij} \frac{y_i^2 + y_j^2}{2} = \sum_{i < j} \pi_i b_{ij} (y_i - y_j)^2 \end{aligned}$$

Denote this final term by $\mathcal{E}(\mathbf{y}, \mathbf{y})$. Then

$$1 - \lambda_2 = \min_{\pi^T\mathbf{y}=0} \frac{\mathcal{E}(\mathbf{y}, \mathbf{y})}{\sum_i \pi_i y_i^2}$$

To prove the first inequality of the theorem, let (S, \bar{S}) be the cut with the minimum conductance. Define a vector \mathbf{w} as follows

$$w_i = \begin{cases} \sqrt{\frac{1}{\sum_u a(u)} \frac{\pi(\bar{S})}{\pi(S)}} & \text{if } i \in S \\ -\sqrt{\frac{1}{\sum_u a(u)} \frac{\pi(S)}{\pi(\bar{S})}} & \text{if } i \in \bar{S} \end{cases}$$

It is then easy to check that $\sum_i \pi_i w_i = 0$ and that

$$\phi(S) \geq \frac{\mathcal{E}(\mathbf{w}, \mathbf{w})}{\sum_i \pi_i w_i^2} \geq 1 - \lambda_2$$

Hence we obtain the desired lower bound on the conductance.

We will now prove the second inequality. Suppose that the minimum above is attained when \mathbf{y} is equal to \mathbf{v} . Then $D\mathbf{v}$ is the eigenvector of Q corresponding to the eigenvalue λ_2 and, \mathbf{v} is the right eigenvector of B corresponding to λ_2 . Our ordering is then with respect to \mathbf{v} in accordance with the statement of the theorem. Assume that, for simplicity of notation, the indices are reordered (i.e. the rows and corresponding columns of B and D are reordered) so that $v_1 \geq v_2 \geq \dots \geq v_N$. Now define r to satisfy $\pi_1 + \pi_2 + \dots + \pi_{r-1} \leq \frac{1}{2} < \pi_1 + \pi_2 + \dots + \pi_r$,

and let $z_i = v_i - v_r$ for $i = 1, \dots, n$. Then $z_1 \geq z_2 \geq \dots \geq z_r = 0 \geq z_{r+1} \geq \dots \geq z_n$, and

$$\begin{aligned} \frac{\mathcal{E}(\mathbf{v}, \mathbf{v})}{\sum_i \pi_i v_i^2} &= \frac{\mathcal{E}(\mathbf{z}, \mathbf{z})}{-v_r^2 + \sum_i \pi_i z_i^2} \geq \frac{\mathcal{E}(\mathbf{z}, \mathbf{z})}{\sum_i \pi_i z_i^2} \\ &= \frac{\left(\sum_{i < j} \pi_i b_{ij} (z_i - z_j)^2 \right) \left(\sum_{i < j} \pi_i b_{ij} (|z_i| + |z_j|)^2 \right)}{\left(\sum_i \pi_i z_i^2 \right) \left(\sum_{i < j} \pi_i b_{ij} (|z_i| + |z_j|)^2 \right)} \end{aligned}$$

Consider the numerator of this final term. By Cauchy-Schwartz

$$\begin{aligned} \left(\sum_{i < j} \pi_i b_{ij} (z_i - z_j)^2 \right) \left(\sum_{i < j} \pi_i b_{ij} (|z_i| + |z_j|)^2 \right) &\geq \left(\sum_{i < j} \pi_i b_{ij} |z_i - z_j| (|z_i| + |z_j|) \right)^2 \\ &\geq \left(\sum_{i < j} \pi_i b_{ij} \sum_{k=i}^{j-1} |z_{k+1}^2 - z_k^2| \right)^2 \end{aligned} \quad (5)$$

Here the second inequality follows from the fact that if $i < j$ then $|z_i - z_j| (|z_i| + |z_j|) \geq \sum_{k=i}^{j-1} |z_{k+1}^2 - z_k^2|$. This follows from observations that

- i) If z_i and z_j have the same sign (i.e. $r \notin \{i, i+1, \dots, j\}$) then $|z_i - z_j| (|z_i| + |z_j|) = |z_i^2 - z_j^2|$.
 - ii) Otherwise, if z_i and z_j have different signs then $|z_i - z_j| (|z_i| + |z_j|) = (|z_i| + |z_j|)^2 > z_i^2 + z_j^2$.
- Also,

$$\sum_{i < j} \pi_i b_{ij} (|z_i| + |z_j|)^2 \leq 2 \sum_{i < j} \pi_i b_{ij} (z_i^2 + z_j^2) = 2 \sum_i \pi_i z_i^2$$

As a result we have,

$$\begin{aligned} \frac{\mathcal{E}(\mathbf{v}, \mathbf{v})}{\sum_i \pi_i v_i^2} &\geq \frac{\left(\sum_{i < j} \pi_i b_{ij} (z_i - z_j)^2 \right) \left(\sum_{i < j} \pi_i b_{ij} (|z_i| + |z_j|)^2 \right)}{\left(\sum_i \pi_i z_i^2 \right) \left(\sum_{i < j} \pi_i b_{ij} (|z_i| + |z_j|)^2 \right)} \\ &\geq \frac{\left(\sum_{i < j} \pi_i b_{ij} \sum_{k=i}^{j-1} |z_{k+1}^2 - z_k^2| \right)^2}{2 \left(\sum_i \pi_i z_i^2 \right)^2} \end{aligned}$$

Set $S_k = \{1, 2, \dots, k\}$, $C_k = \{(i, j) : i \leq k < j\}$ and

$$\hat{\alpha} = \min_{k, 1 \leq k \leq N} \frac{\sum_{(i,j) \in C_k} \pi_i b_{ij}}{\min \left(\sum_{i:i \leq k} \pi_i, \sum_{i:i > k} \pi_i \right)}$$

Since $z_r = 0$, we obtain

$$\begin{aligned}
\sum_{i < j} \pi_i b_{ij} \sum_{k=i}^{j-1} |z_{k+1}^2 - z_k^2| &= \sum_{k=1}^{N-1} |z_{k+1}^2 - z_k^2| \sum_{(i,j) \in C_k} \pi_i b_{ij} \\
&\geq \hat{\alpha} \left(\sum_{k=1}^{r-1} (z_k^2 - z_{k+1}^2) \pi(S_k) + \sum_{k=r}^{N-1} (z_{k+1}^2 - z_k^2) (1 - \pi(S_k)) \right) \\
&= \hat{\alpha} \left(\sum_{k=1}^{N-1} (z_k^2 - z_{k+1}^2) \pi(S_k) + (z_N^2 - z_r^2) \right) = \hat{\alpha} \sum_{k=1}^N \pi_k z_k^2
\end{aligned}$$

Consequently, if $\pi^T \mathbf{y} = 0$ then

$$1 - \lambda_2 = \frac{\mathcal{E}(\mathbf{v}, \mathbf{v})}{\sum_i \pi_i v_i^2} \geq \frac{\hat{\alpha}^2}{2}$$

□

4.2 In the presence of a good clustering

In this section we consider the situation in which a given input matrix A has a particularly good clustering. Here the matrix can be partitioned into blocks such that the conductance of each block as a cluster is high and the total weight of inter-cluster edges is small. We present a result which shows that, in such a circumstance, the spectral algorithm will find a clustering that is close to the optimal clustering i.e. only a small number of rows will be placed in the incorrect cluster.

First we show how to model this situation. We will use the following terminology. Denote by $|\mathbf{x}|$ the 1-norm, and by $\|\mathbf{x}\|$ the 2-norm, of a vector \mathbf{x} . The 2-norm of an $n \times m$ matrix A is

$$\max_{\mathbf{x} \in \mathbb{R}^m, \|\mathbf{x}\|=1} \|A\mathbf{x}\|$$

We assume that A is normalized so that 2-norm of each row is one. In addition, assume A can be written as $B + E$, where B is a block diagonal matrix and E corresponds to the set of edges that run between clusters. Let B be consist of blocks, B_1, B_2, \dots, B_k , which induce the clusters of the optimal clustering.

Rather than conductance, it will be easier to state the result in terms of the minimum *eigenvalue gap* of the blocks B_1, B_2, \dots, B_k . The eigenvalue gap of a matrix is $\beta = 1 - \frac{\lambda_2}{\lambda_1}$. and is closely related to the conductance ($\frac{\phi^2}{2} \leq \beta \leq 2\phi$). Ideally we would like to show that if $A = B + E$ where β is large for each block in B and the total weight of edges in E is small then the spectral algorithm works. While this might be expected for a typical input, it is possible to construct examples, where applying spectral algorithm I does not work.

To handle the problem, we consider a modified version of spectral algorithm I. We make two alterations. Firstly, we will not assign clusters with regards to the matrix C^A (whose r th column is given by $A\mathbf{v}_r$). Instead, the assignment will be with regards to $C = C^A U$, where U is a randomly chosen $k \times k$ orthonormal matrix. Secondly, we adjust the assignment process. For each column r let λ_r be chosen uniformly from the range $[-\frac{q}{2\sqrt{k}}, \frac{q}{2\sqrt{k}}]$, where q will be

a suitably chosen constant. The cluster assignment is then as follows. A pair of rows i and j are placed in separate clusters if $c_{ir} < \lambda_r < c_{jr}$ or $c_{ir} > \lambda_r > c_{jr}$ for any r , $1 \leq r \leq k$. Otherwise rows i and j are placed in the same cluster. Theorem 5 shows that if the β 's are large and the 2-norm of E is small then this spectral algorithm finds a clustering close to the optimal one. Note that this algorithm may create more than k clusters. However, the spectral algorithm finds k clusters “close” to the k optimal ones. For analyses in a similar spirit, see Papadimitriou et al. [11] and Fiat et al. [1].

Theorem 5. *Let the eigenvalue gap of each cluster B_i , $i = 1, \dots, k$, of the optimal clustering be at least β ($0 < \beta \leq 1$). In addition, let the difference between the k th and $k+1$ th eigenvalues of B also be at least β . If $E = A - B$, where $3\|E\| < \beta$ then the spectral algorithm applied to A finds a clustering that differs from the optimal clustering in $o(n)$ rows.*

Proof. The matrix A can be viewed as a perturbation of B by the matrix E . Let the first k eigenvectors of A be denoted by X_1 , and the last $n - k$ by X_2 . Similarly, let the first k eigenvectors of B be denoted by Y_1 , and the last $n - k$ by Y_2 .

Let $B_{ij} = Y_i^T B Y_j$ for $i, j = 1, 2$. Similarly define $E_{ij} = Y_i^T E Y_j$. Then $B_{12} = B_{21} = 0$ and B_{11} is the diagonal matrix of the top k eigenvalues of B . Now define $A_{11} = X_1^T A X_1$. A theorem of Stewart [15] (theorem 4.11, page 745) states that, for some matrix P with 2-norm at most $\frac{2\|E\|}{\beta - 2\|E\|}$, the columns of

$$\bar{X}_1 = (Y_1 + Y_2 P)(I + P^T P)^{-\frac{1}{2}}$$

form an orthonormal basis for a subspace that is invariant for A . Hence there is a $k \times k$ orthonormal matrix U' such that $X_1 = \bar{X}_1 U'$. It follows that $X_1 = Y_1 U' + F$, where $\|F\|$ is less than $\frac{3\|E\|}{\beta - 2\|E\|}$. By a suitable orthonormal transformation, we can take Y_1 to be the matrix whose j th column has as support the rows of B_j .

Recall that we partition our rows, after the transformation by U , with respect to the $n \times k$ matrix $C^A = A X_1 = X_1 A_{11}$. To analyse the performance of our algorithm we first compare the matrices C^A and $C^B = B Y_1 U' = Y_1 B_{11} U'$, then we examine the effect of applying to C^A the transformation U . Denote by \mathbf{a}_i , \mathbf{b}_i and \mathbf{c}_i the i th row of C^A , C^B and $C = C^A U$, respectively. Let $D = C^B - C^A$, then

$$\|D\| = \|C^B - C^A\| = \|B Y_1 U' - (B + E)(Y_1 U' + F)\| = O(\|E\|) + O(\|F\|) = O(1)$$

In particular, if \mathbf{d}_i is the i th row of D then $3\|\mathbf{d}_i\|^2 \leq d^*$ for all but $O(\frac{k^2}{d^*})$ the \mathbf{d}_i . To see this consider the summation $\sum_j \alpha_j d_{ij}$, where α_j is a random variable taking the value $\frac{1}{\sqrt{k}}$ with probability one half, and $-\frac{1}{\sqrt{k}}$ with probability one half. Therefore $E(\alpha_j) = 0$ and $\text{Var}(\alpha_j) = \frac{1}{k}$. Consequently $E(\sum_j \alpha_j d_{ij}) = 0$ and $\text{Var}(\sum_j \alpha_j d_{ij}) = \frac{1}{k} \|\mathbf{d}_i\|^2$. Since the maximum value of $\sum_j \alpha_j d_{ij}$ is $\frac{1}{\sqrt{k}} \|\mathbf{d}_i\|$, it follows that $\sum_j \alpha_j d_{ij}$ is at least $\frac{1}{\sqrt{2k}} \|\mathbf{d}_i\|$ with probability at least $\frac{1}{2k}$. Then, by our bound on $\|D\|$, we see that $3\|\mathbf{d}_i\|^2 \leq d^*$ for the claimed number of rows. Thus most of the \mathbf{a}_i and \mathbf{b}_i are very close together (we will assume the other rows are misplaced, and consider only these “good” rows from now on). Thus if we knew the matrix U' we would be done. The fact that we do not know U' is the reason we multiply C^A by the orthonormal matrix U before clustering. We now have two cases to consider.

(i) Take two rows \mathbf{c}_i and \mathbf{c}_j of C that are that are in the same optimal cluster with respect to

B. Since these two rows are good, it follows that $\|\mathbf{a}_i - \mathbf{a}_j\|^2 \leq d^*$. We now show that \mathbf{c}_i and \mathbf{c}_j are close.

$$\begin{aligned} \sum_r |c_{ir} - c_{jr}| &= \sum_r |\mathbf{a}_i \cdot \mathbf{u}_r - \mathbf{a}_j \cdot \mathbf{u}_r| \\ &= \sum_r |\mathbf{u}_r \cdot (\mathbf{a}_i - \mathbf{a}_j)| \\ &\leq \sum_r \|\mathbf{u}_r\| \|(\mathbf{a}_i - \mathbf{a}_j)\| \\ &\leq \sqrt{d^*} \sum_r \|\mathbf{u}_r\| \leq \sqrt{d^*} k \end{aligned}$$

Then the probability that they are separated by the spectral algorithm is less than

$$\sum_r \frac{|c_{ir} - c_{jr}|}{\frac{q}{\sqrt{k}}} \leq \frac{\sum_r |c_{ir} - c_{jr}|}{\frac{q}{\sqrt{k}}} \leq O\left(k^{\frac{3}{2}} \sqrt{d^*}\right)$$

(ii) Take two rows \mathbf{c}_i and \mathbf{c}_j of C that are that are in different optimal clusters with respect to B . Now, since \mathbf{a}_i and \mathbf{a}_j are good rows, we have $\mathbf{a}_i \cdot \mathbf{a}_j \leq 3\sqrt{d^*}$. Thus we may bound the probability that \mathbf{a}_i and \mathbf{a}_j are correctly separated as follows. Now with constant probability p_0 , both the following events occur: (1) The quantity $|\mathbf{u}_r \cdot \mathbf{a}_i - \mathbf{u}_r \cdot \mathbf{a}_j|$ is at least $\frac{p}{\sqrt{k}}$ for a constant p . (2) The magnitudes of $\mathbf{u}_r \cdot \mathbf{a}_i$ and $\mathbf{u}_r \cdot \mathbf{a}_j$ are each at most $\frac{q}{2\sqrt{k}}$. Therefore, the probability of \mathbf{c}_i and \mathbf{c}_j not being separated is at most $(1 - p_0 \frac{2p}{q})^k \leq 2^{-c_0 k}$ for an absolute constant c_0 .

So the total number of misclassified rows is $O(\frac{k^2}{d^*} + nk^{\frac{3}{2}}\sqrt{d^*} + n2^{-c_0 k})$. Setting $d^* = \theta(k^{\frac{1}{3}}n^{-\frac{2}{3}})$ gives the result, provided that $\Omega(1) \leq k \leq o(n^{\frac{1}{5}})$. \square

5 Conclusion

There are two basic aspects to analyzing a clustering algorithm

- Quality: how good is the clustering produced?
- Speed: how fast can it be found?

In this paper we have mostly dealt with the former issue while taking care that the algorithms are polynomial time. The spectral algorithms depend on the time it takes to find the top (or top k) singular vector(s). While this can be done exactly in polynomial time, it can still too expensive for applications such as information retrieval. The recent work of [7] and [5] on randomized algorithms for low-rank approximation addresses this problem. The running time of their first algorithm depends only on the quality of the desired approximation and not on the size of the matrix, but it assumes that the entries of the matrix can be sampled in a specific manner. Their second algorithm needs no assumptions and has a running time that is linear in the number of non-zero entries.

Acknowledgement. We thank Anna Karlin for pointing out an error in an earlier version.

References

- [1] Y. Azar, A. Fiat, A. Karlin, F. McSherry and J. Saia. "Spectral analysis of data." In *Proc. of 33rd STOC*.
- [2] A. Benczur and D. Karger. "Approximate $s - t$ min-cuts in $O(n^2)$ time." In *Proc. of 28th STOC*, pp47-55, 1996.
- [3] M. Charikar, C. Chekuri, T. Feder and R. Motwani. "Incremental clustering and dynamic information retrieval." In *Proc. of 29th STOC*, 1997.
- [4] M. Charikar, S. Guha, D. Shmoys and E. Tardos. "A constant-factor approximation for the k -median problem." In *Proc. of 31st STOC*, pp1-10, 1999.
- [5] P. Drineas, A. Frieze, R. Kannan, S. Vempala and V. Vinay. "Clustering in large graphs and matrices." In *Proc. of 10th SODA*, 1999.
- [6] M. Dyer and A. Frieze. "A simple heuristic for the p -center problem." In *Oper. Res. Let.*, **3**, pp285-288, 1985.
- [7] A. Frieze, R. Kannan and S. Vempala. "Fast Monte-Carlo algorithms for finding low-rank approximations." In *Proc. of 39th FOCS*, 1998.
- [8] P. Indyk. "A sublinear time approximation scheme for clustering in metric spaces." In *Proc of 40th FOCS*, pp154-59, 1999.
- [9] K. Jain and V. Vazirani. "Primal-dual approximation algorithms for metric facility location and k -median problems." In *Proc. of 40th FOCS*, pp2-13, 1999.
- [10] T. Leighton and S. Rao. "An approximate max-flow min-cut theorem for uniform multicommodity flow problems with applications to approximation algorithms." In *Proc. of 28th FOCS*, pp256-69, 1988.
- [11] C. Papadimitriou, P. Raghaven, H. Tamaki and S. Vempala. "Latent semantic indexing: a probabilistic analysis." In *Proc. of 17th Symposium on the Principles of Database Systems*, 1998.
- [12] J. Shi and J. Malik. "Normalized cuts and image segmentation." In *IEEE Conf. on Computer Vision and Pattern Recognition*, 1997. See <http://www.cs.berkeley.edu/jshi/Grouping>
- [13] A. Sinclair and M. Jerrum. "Approximate counting, uniform generation and rapidly mixing Markov chains." In *Information and Computation*, **82**, pp93-133, 1989.
- [14] D. Spielman and S. Teng. "Spectral partitioning works: planar graphs and finite element meshes." In *Proc. of 37th FOCS*, 1996.
- [15] G. Stewart, "Error and perturbation bounds for subspaces associated with certain eigenvalue problems." In *SIAM review*, **15(4)**, pp727-64, 1973.
- [16] <http://cluster.cs.yale.edu/>