
Local Spectral Clustering of Density Upper Level Sets

Anonymous Authors¹

Abstract

Spectral clustering methods are a family of popular nonparametric clustering tools. Recent works have proposed and analyzed *local* spectral methods, which extract clusters using locally-biased random walks around a user-specified seed node. Several local methods, such as the personalized PageRank (PPR) algorithm, have been shown to have worst-case guarantees for certain graph-based measures of cluster quality. In contrast to existing works, we analyze PPR in a traditional statistical learning setup, where we obtain samples from an unknown distribution, and aim to identify connected regions of high-density (density clusters). We introduce two natural criteria for cluster quality, and derive bounds for these criteria when evaluated on empirical analogues of density clusters. Moreover, we prove that PPR, run on a neighborhood graph, extracts sufficiently salient density clusters.

1. Introduction

Let $\mathbf{X} = \{x_1, \dots, x_n\}$ be a sample drawn i.i.d. from a distribution \mathbb{P} on \mathbb{R}^d , with density f , and consider the problem of clustering: splitting the data into groups which satisfy some notion of within-group similarity and between-group difference. We focus on spectral clustering methods, a family of powerful nonparametric clustering algorithms. Roughly speaking, a spectral technique first constructs a geometric graph G , where vertices are associated with samples, and edges correspond to proximities between samples. It then learns a feature embedding based on the Laplacian of G , and applies a simple clustering technique (such as k-means clustering) in the embedded feature space.

To be more precise, let $G = (V, E, w)$ denote a weighted, undirected graph constructed from the samples \mathbf{X} , where $V = \{1, \dots, n\}$, and $w_{uv} = K(x_u, x_v) \geq 0$ for $u, v \in V$, and a particular kernel function K . Here $(u, v) \in E$ if and

only if $w_{uv} > 0$. We denote by $\mathbf{A} \in \mathbb{R}^{n \times n}$ the weighted adjacency matrix, which has entries $A_{uv} = w_{uv}$, and by \mathbf{D} the degree matrix, with $D_{uu} = \sum_{v \in V} A_{uv}$. We also denote by \mathbf{W}, \mathbf{L} the random walk transition probability matrix and normalized¹ Laplacian matrix, respectively, which are defined as

$$\mathbf{W} = \mathbf{D}^{-1}\mathbf{A}, \quad \mathbf{L} = \mathbf{I} - \mathbf{W},$$

where $\mathbf{I} \in \mathbb{R}^{n \times n}$ is the identity matrix. Classical global spectral methods take an eigendecomposition $\mathbf{L} = \mathbf{U}\Sigma\mathbf{U}^T$, use some number of eigenvectors (columns in \mathbf{U}) as a feature representation for the samples, and then run (say) k-means in this new feature space.

When applied to geometric graphs constructed from a large number of samples, global spectral clustering methods can be computationally cumbersome and insensitive to the local geometry of the underlying distribution (Leskovec et al., 2010; Mahoney et al., 2012). This has led to recent increased interest in local spectral algorithms, which leverage locally-biased spectra computed using random walks around a user-specified seed node. A popular local clustering algorithm is Personalized PageRank (PPR), first introduced by Haveliwalla (2003), and further developed by Spielman & Teng (2011; 2014); Andersen et al. (2006); Mahoney et al. (2012); Zhu et al. (2013), among others.

Local spectral clustering techniques have been practically very successful (Leskovec et al., 2010; Mahoney et al., 2012; ?), which has led many authors to develop theory that explains their success (Zhu et al., 2013; ?), from the perspective of worst-case guarantees on traditional graph-theoretic notions of cluster quality (like conductance). In this paper, we adopt a more traditional statistical viewpoint, and examine what the output of a local clustering algorithm on \mathbf{X} reveals about the unknown density f . In particular, we examine the ability of the PPR algorithm to recover *density clusters* of f , which are defined as the connected components of the upper level set $\{x : f(x) \geq \lambda\}$ for some threshold $\lambda > 0$ (a central object of central interest in the classical statistical literature on clustering, dating back to Hartigan 1981).

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

¹Other popular choices here include the unnormalized Laplacian, and symmetric normalized Laplacian.

1.1. Graph Connectivity Criteria

Here we define a pair of criteria that reflect the quality of a cluster with respect to $G = (V, E, w)$. There are many graph-based measures of cluster quality that one could consider; see, e.g., Yang & Leskovec (2015); Fortunato (2010) for an overview. The pair of criteria that we focus on are (arguably) quite natural, and moreover, they play a fundamental role in our analysis of the PPR algorithm. Our two criteria capture the *external* and *internal* connectivity of a subset $S \subseteq V$, denoted $\Phi(S; G)$ and $\Psi(S; G)$, respectively, and defined below in turn.

External Connectivity: Normalized Cut. Define the cut between subsets $S, S' \subseteq V$ to be

$$\text{cut}(S, S'; G) = \sum_{u \in S} \sum_{v \in S'} w_{uv},$$

and define $\text{vol}(S; G) = \text{cut}(S, V; G) = \sum_{u \in S} \sum_{v \in V} w_{uv}$. As our notion of external connectivity, we use the *normalized cut* of S , defined as

$$\Phi(S; G) = \frac{\text{cut}(S; G)}{\min\{\text{vol}(S; G), \text{vol}(S^c; G)\}}, \quad (1)$$

where we abbreviate $\text{cut}(S; G) = \text{cut}(S; S^c; G)$.

Internal Connectivity: Inverse Mixing Time. For $S \subseteq V$, denote by $G[S] = (S, E_S, w_S)$ the subgraph induced by S (where the edges are $E_S = E \cap (S \times S)$). Let $\mathbf{A}_S, \mathbf{D}_S$ be the adjacency matrix and degree matrix, respectively, of $G[S]$. Define the random walk matrix as usual, $\mathbf{W} = \mathbf{D}_S^{-1} \mathbf{A}_S$, and for $v \in V$, write

$$q_{vu}^{(t)} = e_v \mathbf{W}_S^t e_u \quad (2)$$

for the t -step transition probability of a random walk over $G[S]$ originating at v .² Also write $\tilde{\pi} = (\tilde{\pi}_u)_{u \in S}$ for the stationary distribution of this random walk. (Given the definition of \mathbf{W}_S , it is well-known that the stationary distribution is given by $\tilde{\pi}_u = (\mathbf{D}_S)_{uu} / \text{vol}(S; G[S])$.)

Our internal connectivity parameter will capture the time it takes for the random walk over $G[S]$ to mix (approach the stationary distribution) uniformly over S . For this, we first define *relative pointwise mixing time* of $G[S]$ as

$$\tau_\infty(\mathbf{q}; G[S]) = \min \left\{ t : \frac{|q_{vu}^{(m)} - \tilde{\pi}_u|}{\tilde{\pi}_u} \leq \frac{1}{4}, \text{ for } u, v \in V \right\},$$

where $\mathbf{q} = (\mathbf{q}_v^{(1)}, \mathbf{q}_v^{(2)}, \dots)_{v \in V}$, and $\mathbf{q}_v^{(m)} = (q_{vu}^{(m)})_{u \in V}$. Now our internal connectivity parameter is simply the inverse mixing time,

$$\Psi(S; G) = \frac{1}{\tau_\infty(\mathbf{q}; G[S])}. \quad (3)$$

²Given a starting node v and a random walk defined by transition probability matrix \mathbf{P} , the rotation $e_v \mathbf{P}^t$ is used to denote the distribution of the random walk after t steps.

If S has normalized cut no greater than Φ , and inverse mixing time no less than Ψ , we call it as a (Φ, Ψ) -cluster. Both local (Zhu et al., 2013) and global (Kannan et al., 2004) spectral algorithms have been shown to output clusters (or partitions) which approximate the optimal (Φ, Ψ) -cluster (or partition) for a given graph G .³

1.2. PPR on a Neighborhood Graph

We now describe the clustering algorithm that will be our focus for the rest of the paper. We start with the geometric graph that we form based on the samples \mathbf{X} : for a radius $r > 0$, we consider the r -neighborhood graph of \mathbf{X} , denoted $G_{n,r} = (V, E)$, an unweighted graph with vertices $V = \{1, \dots, n\}$, and an edge $(u, v) \in E$ if and only if $\|x_u - x_v\| \leq r$, where $\|\cdot\|$ denotes Euclidean norm. Note that this is a special case of the general construction introduced above, with $K(u, v) = 1(\|x_u - x_v\| \leq r)$.

Next, we define the PPR vector $\mathbf{p} = \mathbf{p}(v, \alpha; G_{n,r})$, with respect to a seed node $v \in V$ and a teleportation parameter $\alpha \in [0, 1]$, to be the solution of the following linear system:

$$\mathbf{p} = \alpha \mathbf{e}_v + (1 - \alpha) \mathbf{p} \mathbf{W}, \quad (4)$$

where \mathbf{W} is the random walk matrix of the underlying graph $G_{n,r}$ and \mathbf{e}_v denotes indicator vector for node v (with a 1 in the v th position and 0 elsewhere). In practice, we can approximately solve the above linear system via a simple, efficient random walk, with appropriate restarts to v .

For a level $\beta > 0$ and a target volume $\pi_0 > 0$, we define a β -sweep cut of \mathbf{p} as

$$S_\beta = \{u \in V : p_u > \beta \pi_0\}. \quad (5)$$

Having computed sweep cuts over a range $\beta \in (\frac{3}{10}, \frac{1}{2})$,⁴ we output a cluster $\hat{C} = S_{\beta^*}$, based on the sweep cut S_{β^*} that minimizes the normalized cut $\Phi(S_{\beta^*}; G_{n,r})$ as defined in (1). For concreteness, we summarize this procedure in Algorithm 1.

1.3. Summary of Results

Let $\mathcal{C}_f(\lambda)$ denote the connected components of the density upper level set $\{x : f(x) > \lambda\}$. For a given density cluster $\mathcal{C} \in \mathcal{C}_f(\lambda)$, we call $\mathcal{C}[\mathbf{X}] = \mathcal{C} \cap \mathbf{X}$ the *empirical density cluster*. Below we define a notion of consistency in density cluster estimation.

³In the case of (Kannan et al., 2004), the internal connectivity parameter ϕ is actually the conductance, i.e., the minimum normalized cut within the subgraph $G[S]$. See Theorem 3.1 in their paper for details; however, note that $\phi^2 / \log(\text{vol}(S)) \leq O(\Psi)$, and so the lower bound on ϕ translates to a lower bound on Ψ .

⁴The choice of a specific range such as $(\frac{3}{10}, \frac{1}{2})$ is standard in the analysis of PPR algorithms, see, e.g., Zhu et al. (2013).

Algorithm 1 PPR on a Neighborhood Graph

Input: data $\mathbf{X} = \{x_1, \dots, x_n\}$, radius $r > 0$, teleportation parameter $\alpha \in [0, 1]$, seed $v \in \mathbf{X}$, target volume $\pi_0 > 0$.

Output: cluster $\hat{\mathcal{C}} \subseteq V$.

- 1: Form the neighborhood graph $G_{n,r}$.
- 2: Compute the PPR vector $\mathbf{p}(v, \alpha; G_{n,r})$ as in (4).
- 3: For $\beta \in (\frac{3}{10}, \frac{1}{2})$ compute sweep cuts S_β as in (5).
- 4: Return $\hat{\mathcal{C}} = S_{\beta^*}$, where

$$\beta^* = \arg \min_{\beta \in (\frac{3}{10}, \frac{1}{2})} \Phi(S_\beta; G_{n,r}).$$

Definition 1 (Consistent density cluster estimation). *For an estimator $\hat{\mathcal{C}} \subseteq \mathbf{X}$ and cluster $\mathcal{C} \in \mathbb{C}_f(\lambda)$, we say $\hat{\mathcal{C}}$ is a consistent estimator of \mathcal{C} if for all $\mathcal{C}' \in \mathbb{C}_f(\lambda)$ with $\mathcal{C} \neq \mathcal{C}'$ the following holds as $n \rightarrow \infty$:*

$$\mathcal{C}[\mathbf{X}] \subseteq \hat{\mathcal{C}} \quad \text{and} \quad \hat{\mathcal{C}} \cap \mathcal{C}'[\mathbf{X}] = \emptyset, \quad (6)$$

with probability tending to 1.

A summary of our main results (and outline for the rest of this paper) is as follows.

1. In Section 2, we derive in Theorem 1 an upper bound on the normalized cut of an empirical density cluster $\mathcal{C}[\mathbf{X}]$, under a natural set of geometric conditions (precluding clusters which are too thin and long).
2. Under largely the same set of geometric conditions, we derive in Theorem 2 a lower bound on the inverse mixing time of a random walk over $\mathcal{C}[\mathbf{X}]$. We also provide in Theorem 3 a tighter lower bound, but under more restrictive assumptions (convexity of $\mathcal{C}[\mathbf{X}]$).
3. In Section 3, we show in Theorem 4 that these bounds in bounds in Theorems 1 and 2, on the cluster quality criteria, have algorithmic consequences for PPR: properly initialized, Algorithm 1 performs consistent density cluster estimation in the sense of (6).
4. We show in Corollary 1 that Theorems 1 and 2, along with the results in Zhu et al. (2013), lead to alternative, graph-theoretic guarantees on cluster quality: an upper bound on the normalized cut of the estimated cluster $\hat{\mathcal{C}}$, and an upper bound on volume of the symmetric set difference between $\hat{\mathcal{C}}$ and $\mathcal{C}[\mathbf{X}]$.
5. In Section 4, we discuss the implications of our results for some example density functions of interest, and empirically demonstrate that violations of the geometric conditions we require manifestly impact density cluster recovery.

On the topic of conditions, it is worth mentioning that, as density clusters are inherently local, focusing on the PPR al-

gorithm actually eases our analysis and allows us to require fewer global regularity conditions relative to those needed for more classical global spectral algorithms.

1.4. Related Work

In addition to the background given above, a few related lines of work are worth highlighting. Global spectral clustering methods were first developed in the context of graph partitioning (Fiedler, 1973; Donath & Hoffman, 1973) and their performance is well-understood in this context (see, e.g., Tolliver & Miller 2006; von Luxburg 2007). In a similar vein, several recent works (McSherry, 2001; Rohe et al., 2011; Chaudhuri et al., 2012; Balakrishnan et al., 2011; Lei & Rinaldo, 2015; Abbe, 2018) have studied the efficacy of spectral methods in successfully recovering the community structure in the stochastic block model and variants.

von Luxburg et al. (2008); Hein et al. (2005), building on earlier work of Koltchinskii & Gine (2000), studied the limiting behaviour of spectral clustering algorithms. These authors show that when samples are obtained from a distribution, and we appropriately construct a geometric graph, the spectrum of the Laplacian converges to that of the Laplace-Beltrami operator on the data-manifold. However, relating the partition obtained using the Laplace-Beltrami operator to the more intuitively defined high-density clusters can be challenging in general.

Perhaps most similar to our results are the works Vempala & Wang (2004); Shi et al. (2009); Schiebinger et al. (2015), who study the consistency of spectral algorithms in recovering the latent labels in certain parametric and nonparametric mixture models. These results focus on global rather than local algorithms, and as such impose global rather than local conditions on the nature of the density. Moreover, they do not in general ensure recovery of density clusters, which is the focus in our work.

2. Cluster Quality Criteria on Empirical Density Clusters

In order to provide meaningful bounds on the normalized cut and inverse mixing time of an empirical density cluster $\mathcal{C}[\mathbf{X}]$, we must introduce some assumptions on the density f .

Let $B(x, r) = \{y \in \mathbb{R}^d : \|y - x\| \leq r\}$ be a closed ball of radius r around the point x . Given a set $\mathcal{A} \subseteq \mathbb{R}^d$, and a number $\sigma > 0$, define the σ -expansion of \mathcal{A} to be $\mathcal{A}_\sigma = \mathcal{A} + B(0, \sigma) = \{y \in \mathbb{R}^d : \inf_{x \in \mathcal{A}} \|y - x\| \leq \sigma\}$. We are now ready to give the assumptions, which we state with respect to a density cluster $\mathcal{C} \in \mathbb{C}_f(\lambda)$ for some $\lambda > 0$, and expansion parameter $\sigma > 0$:

(A1) *Cluster separation*: For all $C' \in \mathbb{C}_f(\lambda)$,

$$\text{dist}(C_\sigma, C'_\sigma) > \sigma,$$

where $\text{dist}(\mathcal{A}, \mathcal{A}') = \min_{x \in \mathcal{A}} \text{dist}(x, \mathcal{A}')$ for $\mathcal{A}' \subseteq \mathbb{R}^d$.

(A2) *Cluster diameter*: There exists $D < \infty$ such that for all $x, x' \in C_\sigma$:

$$\|x - x'\| \leq D.$$

(A3) *Bounded density within cluster*: There exist numbers $0 < \lambda_\sigma < \Lambda_\sigma < \infty$ such that:

$$\lambda_\sigma = \inf_{x \in C_\sigma} f(x) \leq \sup_{x \in C_\sigma} f(x) \leq \Lambda_\sigma \quad (7)$$

(A4) *Low noise density*: For some $\gamma > 0$, there exists a constant $c_1 > 0$ such that for all $x \in \mathbb{R}^d$ with $0 < \text{dist}(x, C_\sigma) \leq \sigma$,

$$\inf_{x' \in C_\sigma} f(x') - f(x) \geq c_1 \text{dist}(x, C_\sigma)^\gamma,$$

where $\text{dist}(x, \mathcal{A}) = \min_{x_0 \in \mathcal{A}} \|x - x_0\|$ for $\mathcal{A} \subseteq \mathbb{R}^d$.

We note that σ plays several roles here, precluding arbitrarily narrow clusters and long clusters in (A2) and (A3), flat densities around the level set in (A4), and poorly separated clusters in (A1).

Assumptions (A1), (A3) and (A4) are used to upper bound $\Phi(\mathcal{C}[\mathbf{X}]; G_{n,r})$, whereas (A2) and (A3) are necessary to lower bound $\Psi(\mathcal{C}[\mathbf{X}]; G_{n,r})$. We note that the lower bound on minimum density in (7) and (A1) combined are similar to the (σ, ϵ) -saliency of (Chaudhuri & Dasgupta, 2010), a standard density clustering assumption, while (A4) is seen in, for instance, (Singh et al., 2009), (as well as many other works on density clustering and level set estimation.) It is worth highlighting that these assumptions are all local in nature, a benefit of studying a local algorithm such as PPR.

We are ready to provide bounds on the graph quality bi-criteria. For notational simplicity, hereafter for $S \subseteq \mathbf{X}$ we will refer to $\Phi(S; G_{n,r})$ as $\Phi_{n,r}(S)$, and likewise with $\Psi(S; G_{n,r})$ and $\Psi_{n,r}(S)$. We will also use $\nu(\cdot)$ to denote the uniform measure over \mathbb{R}^d , and $\nu_d = \nu(B(0, d))$ as the measure of the unit ball.

We begin with an upper bound on the normalized cut in Theorem 1. We will require Assumptions (A1), (A3) and (A4) to hold; however, the upper bound on density in (7) will not be needed and so we omit the parameter Λ_σ from the statement of the theorem.

Theorem 1. For some $\lambda > 0$, let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ satisfy Assumptions (A1), (A3) and (A4) for some $\sigma, \lambda_\sigma, c_1, \gamma > 0$. Then, for any $r < \sigma/4d$ and $\delta \in (0, 1]$, the following

statements hold with probability at least $1 - \delta$: Fix $\epsilon > 0$. Then, for

$$n \geq \frac{9 \log(2/\delta)}{\epsilon^2} \left(\frac{1}{\lambda_\sigma^2 \nu(C_\sigma) \nu_d r^d} \right)^2 \quad (8)$$

we have

$$\frac{\Phi_{n,r}(C_\sigma[\mathbf{X}])}{r} \leq 4c_\sigma d \frac{\lambda}{\lambda_\sigma} \frac{(\lambda_\sigma - \frac{r^\gamma}{\gamma+1})}{\lambda_\sigma} + \epsilon \quad (9)$$

where $c_\sigma = 1/\sigma$.

Remark 1. The proof of Theorem 1, along with all other proofs, can be found in the supplementary document. The key point is to note that for any $x \in \mathcal{C}$, the simple, (possibly loose) $B(x, \sigma) \subseteq C_\sigma$ translates to the upper bound $\nu(C_{\sigma+r}) \leq (1 + 2dr/\sigma)\nu(C_\sigma)$. We leverage (A4) to find a corresponding bound on the weighted volume, before applying standard concentration inequalities to convert from population to sample based results.

Remark 2. (9) is almost tight. Specifically, choosing

$$\mathcal{A}_\sigma = B(0, \sigma),$$

$$f(x) = \begin{cases} \lambda & \text{for } x \in \mathcal{A}_\sigma, \\ \lambda - \text{dist}(x, \mathcal{A}_\sigma)^\gamma & \text{for } 0 < \text{dist}(x, \mathcal{A}_\sigma) < r \end{cases}$$

we have that for n within constant order of the lower bound in (8), with probability at least $1 - \delta$

$$\frac{\Phi_{n,r}(\mathcal{A}_\sigma[\mathbf{X}])}{r} \geq c \frac{(\lambda - \frac{r^{\epsilon+1}}{\epsilon+1})}{\lambda} - \epsilon$$

for some constant c . (Note that a factor of $1/\sigma$ in c_σ is not replicated in this lower bound.)

We now provide a lower bound on $\Psi_{n,r}(\mathcal{C}[\mathbf{X}])$.

Theorem 2. Fix $\lambda > 0$, and let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ satisfy Assumptions (A2) and (A3) for some $\sigma, \lambda_\sigma, \Lambda_\sigma, D > 0$. Then, for any $r < \sigma/4d$, the following statement holds with probability at least $1 - \delta$: Fix $0 < \epsilon > 1$. Then, for n satisfying:

$$\sqrt{3^{d+1} \frac{(\log n + d \log \mu + \log(4/\delta))}{n \nu_d r^d \lambda_\sigma}} \leq \epsilon$$

we have

$$\frac{1}{\Psi_{n,r}(C_\sigma[\mathbf{X}])} \leq (d \log \mu + c_\lambda) \cdot \left(c_1 + c_2 \frac{\Lambda_\sigma^4 D^{2d}}{\lambda_\sigma^4 r^{2d}} (d \log \mu + c_\lambda) \right) \quad (10)$$

where $\mu = \log(\frac{2D}{r})$; c_1, c_2 , and c_3 are constants which depend only on dimension; and $c_\lambda = \log(\Lambda_\sigma^2/\lambda_\sigma^2)$.

Remark 3. The proof of Theorem 2 relies on upper bounding the mixing time using the *conductance* of $G_{n,r}[\mathcal{C}_\sigma[X]]$,

$$\tilde{\Phi} = \min_{S \subseteq \mathcal{C}_\sigma[X]} \Phi(S; G_{n,r}[\mathcal{C}_\sigma[X]])$$

The factor of $\frac{1}{r^{2d}}$ present in the bound of (10) is suboptimal. This exponential dependence on d stems from a loose bound on the aforementioned conductance in the proof of Theorem 2. In particular, we assert only that any set $S \subseteq \mathcal{C}_\sigma[X]$ must have $\text{cut}(S; \mathcal{C}_\sigma[X])$ on the order of $n^2 r^{2d}$, while upper bounding $\text{vol}(S; \mathcal{C}_\sigma[X])$ by roughly $n^2 r^d$, for a bound on the conductance of order r^d . The presence of r^{2d} comes from upper bounding the mixing time by about $1/\tilde{\Phi}^2$, this latter bound being a variant of classic results on rapid mixing (Jerrum & Sinclair, 1989).

It is possible to sharpen the dependency on d , but at the cost of an additional assumption:

(A5) \mathcal{C} is convex.

Additionally, the resulting bound holds only asymptotically.

Theorem 3. Fix $\lambda > 0$, and let the conditions of Theorem 2 hold with respect to $\mathcal{C} \in \mathbb{C}_f(\lambda)$. Additionally, let \mathcal{C} satisfy (A5). Then, for any $r < \sigma/4d$, the following statement holds: with probability one

$$\liminf_{n \rightarrow \infty} \frac{1}{\Psi_{n,r}(\mathcal{C}_\sigma[X])} \geq c_1 \frac{\Lambda_\sigma^8}{\lambda_\sigma^8} (\log \mu + c_\lambda) \left(c_\lambda (d^3 \log \mu + c_2) + \frac{d^2 D^2}{r^2} (\log d + c_2 \frac{\Lambda_\sigma^2}{\lambda_\sigma^2} c_\lambda + c_3 + \log \log \mu) \right) + c_d \frac{o_r(1)}{r^2}, \quad (11)$$

where $\mu = \log(\frac{2D}{r})$, c_1, c_2 and c_3 are all global constants, c_d is constant in r but may depend on other quantities, and $c_\lambda = \log(\Lambda_\sigma^2/\lambda_\sigma^2)$. Additionally, $\lim_{r \rightarrow 0} o_r(1) = 0$.

We achieve superior rates in Theorem 3 to those of Theorem 2 in part by working with a generalization of the conductance, the *conductance function*

$$\tilde{\Phi}_n(t) = \min_{\substack{S \subseteq \mathcal{C}_\sigma[X] \\ \pi(S) \leq t}} \Phi(S; G_{n,r}[\mathcal{C}_\sigma[X]])$$

where $\pi(S)$ is the stationary distribution over $G_{n,r}[\mathcal{C}_\sigma[X]]$. The utility of the conductance function comes from the known upper bound of mixing time by

$$\int \frac{1}{t \tilde{\Phi}_n(t)^2} dt \quad (12)$$

which results in a tighter bound than merely using the conductance when $\tilde{\Phi}_n(t)$ is large for small values of t .

Remark 4. The only potentially exponential dependence on dimension comes from the factor of c_d . However, for sufficiently small values of r this will be dominated by the preceding factors, due to the presence of the $o_r(1)$ term.

Remark 5. The proof of Theorem 3 relies on a to the best of our knowledge novel uniform lower bound on this conductance function over $G_{n,r}$ in terms of a population-level analogue, which we term $\tilde{\Phi}_{\mathbb{P},r}$. Plugging $\tilde{\Phi}_{\mathbb{P},r}$ into (12) yields a bound on mixing time (with respect to total variation distance) of order $dD^2/r^2 + d^2 \log(D/r)$ over convex sets. By contrast, (11) is of order $d^2 D^2/r^2 + d^3 \log(D/r)$ (ignoring the $o_r(1)$ term) but handles mixing time with respect to relative pointwise distance (which is known to be a stricter metric.)

Remark 6. The convexity requirement is necessary only to lower bound the $\tilde{\Phi}_{\mathbb{P},r}$; any bound on $\tilde{\Phi}_{\mathbb{P},r}$ which does not require convexity could immediately be plugged into the machinery of the proof of Theorem 3 to achieve a corresponding result. Recent work (Abbasi-Yadkori et al., 2017) has developed such population-level bounds on the conductance function, however the Markov chain dealt with there is somewhat different than the one we consider.

3. Consistent cluster estimation with PPR.

For PPR to successfully recover density clusters, the ratio $\Phi_{n,r}(\mathcal{C}_\sigma[X])/\Psi_{n,r}(\mathcal{C}_\sigma[X])$ must be small.

We introduce

$$\begin{aligned} \Phi(\sigma, \lambda, \lambda_\sigma, \gamma) &:= \frac{\lambda}{\lambda_\sigma} \frac{(\lambda_\sigma - \frac{\sigma^\gamma}{\gamma+1})}{\lambda_\sigma} \\ 1/\Psi(\sigma, \lambda_\sigma, \Lambda_\sigma, D) &:= (d \log \mu + c_\lambda) \cdot \\ &\quad \left(c_1 + c_2 \frac{\Lambda_\sigma^4 D^{2d}}{\lambda_\sigma^4 r^{2d}} (d \log \mu + c_\lambda) \right) \end{aligned}$$

Well-conditioned density clusters satisfy all of the given assumptions, for parameters which results in ‘good’ values of Φ and Ψ .

Definition 2 (Well-conditioned density clusters). For $\lambda > 0$ and $\mathcal{C} \in \mathbb{C}_f(\lambda)$, let \mathcal{C} satisfy (A1) - (A4) with respect to parameters $\sigma, \lambda_\sigma, \gamma > 0$ and $\Lambda_\sigma, D < \infty$. Letting $\kappa_1(\mathcal{C})$ and $\kappa_2(\mathcal{C})$ be given by

$$\begin{aligned} \kappa_1(\mathcal{C}) &:= \frac{\Phi(\sigma, \lambda, \lambda_\sigma, \gamma)}{\Psi(\sigma, \lambda_\sigma, \Lambda_\sigma, D)} \\ \kappa_2(\mathcal{C}) &:= \kappa_1(\mathcal{C}) \cdot \sqrt{\Psi(\sigma, \lambda_\sigma, \Lambda_\sigma, D)}, \end{aligned}$$

we call \mathcal{C} a (κ_1, κ_2) -well-conditioned density cluster (with respect to $\sigma, \lambda_\sigma, \gamma, \Lambda_\sigma$ and D).

Φ and Ψ are familiar; they are exactly the upper and lower bounds on $\Phi_{n,r}(\mathcal{C}_\sigma[\mathbf{X}])$ and $\Psi_{n,r}(\mathcal{C}_\sigma[\mathbf{X}])$ derived in Theorems 1 and 2, respectively.

Remark 7. For convenience and maximum generality, we define $\Psi(\sigma, \lambda_\sigma, \Lambda_\sigma, D)$ to correspond with the bound given by (10), and assume only (A1) - (A4). However, if we additionally have (A5), then we could sharpen $\Psi(\sigma, \lambda_\sigma, \Lambda_\sigma, D)$ to the tighter rate of (11), with nothing changing hereafter.

As is typical in the local clustering literature, our results will be stated with respect to specific choices or ranges of each of the user-specified parameters, which in this case may depend on the underlying (unknown) density.

In particular, for a well conditioned density cluster \mathcal{C} (with respect to some $\sigma, \lambda_\sigma, \gamma, \Lambda_\sigma$ and D), we require

$$\alpha \in [1/10, 1/9] \cdot \Psi(\sigma, \lambda_\sigma, \Lambda_\sigma, D), r \leq \sigma/4d$$

$$\pi_0 \in [2/3, 6/5] \frac{\lambda_\sigma}{\nu(\mathcal{C}_\sigma)\Lambda_\sigma^2}, v \in \mathcal{C}_\sigma[\mathbf{X}]^g \quad (13)$$

where $\mathcal{C}_\sigma[\mathbf{X}]^g \subseteq \mathcal{C}_\sigma[\mathbf{X}]$ is some 'good' subset of $\mathcal{C}_\sigma[\mathbf{X}]$ which, as we will see, satisfies $\text{vol}(\mathcal{C}_\sigma[\mathbf{X}]^g) \geq \text{vol}(\mathcal{C}_\sigma[\mathbf{X}])/2$. (Intuitively one can think of $\mathcal{C}_\sigma[\mathbf{X}]^g$ as being the nodes sufficiently close to the center of $\mathcal{C}_\sigma[\mathbf{X}]$, although we provide no formal justification to this effect.)

Definition 3. If the input parameters to Algorithm 1 satisfy 13 with respect to some $\mathcal{C}_\sigma[\mathbf{X}]$, we say the algorithm is well-initialized.

Theorem 4. Fix $\lambda > 0$, and let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ be a (κ_1, κ_2) -well conditioned cluster (with respect to some $\sigma, \lambda_\sigma, \gamma, \Lambda_\sigma$ and D). If

$$\kappa_2 \leq \frac{1}{40 \cdot 36} \frac{\lambda_\sigma^2}{\Lambda_\sigma^2} \frac{r^d \nu_d}{\nu(\mathcal{C}_\sigma)}. \quad (14)$$

and Algorithm 1 is well-initialized, the output set $\hat{\mathcal{C}} \subseteq \mathbf{X}$ is a consistent estimator for \mathcal{C} , in the sense of Definition 1.

Remark 8. We note that larger constants than $40 \cdot 36$ allow for wider range of the parameters in (13).

Approximate cluster recovery via PPR. In (Zhu et al., 2013), building on the work of (Andersen et al., 2006) and others, theory is developed which links algorithmic performance of PPR to the normalized cut and mixing time parameters. Although not the primary focus of our work, it is perhaps worth noting that these results, coupled with Theorems 1-3, translate immediately into bounds on the normalized cut and symmetric set difference of $\hat{\mathcal{C}}$.

We collect some of the main results of (Zhu et al., 2013) in Lemma 1.

For $G = (V, E)$ consider some $A \subseteq V$, and let $\Phi(A; G)$ and $\Psi(A; G)$ be defined as in (1) and (3), respectively.

Lemma 1 (PPR clustering). *There exists a set $A^g \subseteq A$ with $\text{vol}(A^g; G) \geq \text{vol}(A; G)/2$ such that the following statement holds: Choose any $v \in A^g$, fix $\alpha = 9/10\Psi(A; G)$, and compute the page rank vector $\mathbf{p}(v, \alpha; G)$. Letting*

$$\hat{\mathcal{C}} = \arg \min_{\beta \in [\frac{1}{8}, \frac{1}{2}]} \Phi(S_\beta; G)$$

the following guarantees hold:

$$\text{vol}(\hat{\mathcal{C}} \setminus A) \leq \frac{24\Phi(A; G)}{\Psi(A; G)} \text{vol}(A)$$

$$\text{vol}(A \setminus \hat{\mathcal{C}}) \leq \frac{30\Phi(A; G)}{\Psi(A; G)} \text{vol}(A)$$

$$\Phi(\hat{\mathcal{C}}; G) = O\left(\frac{\Phi(A; G)}{\sqrt{\Psi(A; G)}}\right)$$

Corollary 1 immediately follows.

Corollary 1. Fix $\lambda > 0$, and let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ be a (κ_1, κ_2) -well conditioned cluster (with respect to some $\sigma, \lambda_\sigma, \gamma, \Lambda_\sigma$ and D). Then, if Algorithm 1 is well-initialized (in the sense that the choices of input parameters satisfy (13)), the following guarantees hold for output set $\hat{\mathcal{C}} \subseteq \mathbf{X}$:

$$\text{vol}(\hat{\mathcal{C}} \setminus \mathcal{C}_\sigma[\mathbf{X}]), \text{vol}(\mathcal{C}_\sigma[\mathbf{X}] \setminus \hat{\mathcal{C}}) \leq 30\kappa_1(\mathcal{C}) \text{vol}(\mathcal{C}_\sigma[\mathbf{X}])$$

$$\Phi_{n,r}(\hat{\mathcal{C}}) = O(\kappa_2(\mathcal{C}))$$

4. Examples

Example 1 is intended to show how the machinery developed above translates in a specific, common mixture model, and the extent to which bounds are (or are not) tight. Example 2 will try to delve into some of the details of how PPR interpolates the conductance and density cut, and will show a case where a poorly conditioned density cluster is not recovered by PPR. Example 3 will emphasize finite sample cluster recovery, for a well-conditioned but non-convex mixture model

Examples 1 and 2 should be thought of as shedding light on the population performance of PPR, whereas Example 3 shows performance on a finite sample.

1. *Gaussian Mixture Model:* We will compute optimal Φ and Ψ for given λ , and show the following

- A graph comparing Φ to $\Phi_{n,r}$ as the value of λ changes.
- A graph comparing Ψ to $\Psi_{n,r}$ as the value of λ changes.
- That for some values of λ , the conditions required for Theorem 4 hold.

2. *Thin and long parallel clusters, with ϵ -uniform noise:*

We will (try to) show that the set outputted by PPR interpolates between the minimum normalized cut solution (fatter) and the density cluster (thinner). The conductance is

$$\Phi^*(G_{n,r}) := \min_{C \subseteq \mathbf{X}} \Phi_{n,r}(C)$$

and the conductance cut is $C^* \subseteq \mathbf{X}$ which achieves the minimum.

We will show that

- For sufficiently small ϵ , all three of the conductance cut, PPR cut, and density cut agree.
- For an intermediate value of ϵ , the conductance cut and the density cut disagree. The PPR cut interpolates between the two.
- For a sufficiently large value of ϵ , the PPR cut fails to recover the density cut, and draws closer to the conductance cut.

3. *Non-convex mixture model:* We will show that, for well-conditioned non-convex mixture model, and a finite sample size n , cluster recovery is achieved with high probability over repeated simulations.

5. Discussion

For a clustering algorithm and a given object (such as a graph or set of points), there are an almost limitless number of ways to define what the 'right' clustering is. We have considered a few such ways – density level sets, and the bicriteria of normalized cut, inverse mixing time – and shown that under the right conditions, the latter agree with the former, with resulting algorithmic consequences.

There are still many directions worth pursuing in this area. Concretely, we might wish to generalize our results to hold over a wider range of kernel functions, and hyperparameter inputs to the PPR algorithm. More broadly, we do not provide any sort of theoretical lower bound, although we give empirical evidence in Example 2 that poorly conditioned density clusters are not consistently estimated by PPR. Example 2 also hints at a way of understanding local spectral algorithms – or, at least, PPR – as interpolating between normalized and density cut. Exploring this connection is an avenue for future work.

References

- Abbasi-Yadkori, Y., Bartlett, P., Gabillon, V., and Malek, A. Hit-and-run for sampling and planning in non-convex spaces. In *Artificial Intelligence and Statistics*, pp. 888–895, 2017.
- Abbe, E. Community detection and stochastic block models. *Foundations and Trends® in Communications and Information Theory*, 14(1-2):1–162, 2018.
- Andersen, R., Chung, F., and Lang, K. Local graph partitioning using pagerank vectors. In *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science*, pp. 475–486, 2006.
- Balakrishnan, S., Xu, M., Krishnamurthy, A., and Singh, A. Noise thresholds for spectral clustering. In *Advances in Neural Information Processing Systems 24*. Curran Associates, Inc., 2011.
- Chaudhuri, K. and Dasgupta, S. Rates of convergence for the cluster tree. In Lafferty, J. D., Williams, C. K. I., Shawe-Taylor, J., Zemel, R. S., and Culotta, A. (eds.), *Advances in Neural Information Processing Systems 23*, pp. 343–351. Curran Associates, Inc., 2010.
- Chaudhuri, K., Graham, F. C., and Tsiatas, A. Spectral clustering of graphs with general degrees in the extended planted partition model. In *COLT*, volume 23, pp. 35.1–35.23, 2012.
- Donath, W. E. and Hoffman, A. J. Lower bounds for the partitioning of graphs. *IBM J. Res. Dev.*, 17(5):420–425, September 1973.
- Ery, A.-C., Pelletier, B., and Pudlo, P. The normalized graph cut and cheeger constant: from discrete to continuous. *Advances in Applied Probability*, 44(4):907–937, 12 2012.
- Fiedler, M. Algebraic connectivity of graphs. *Czechoslovak Mathematical Journal*, 23(2):298–305, 1973.
- Fortunato, S. Community detection in graphs. *Physics Reports*, 486(3):75 – 174, 2010. ISSN 0370-1573.
- Hartigan, J. A. Consistency of single-linkage for high-density clusters. *Journal of the American Statistical Association*, 1981.
- Haveliwala, T. H. Topic-sensitive pagerank: A context-sensitive ranking algorithm for web search. *IEEE transactions on knowledge and data engineering*, 15(4):784–796, 2003.
- Hein, M., Audibert, J.-Y., and von Luxburg, U. From graphs to manifolds – weak and strong pointwise consistency of graph laplacians. In *Learning Theory*, 2005.
- Jerrum, M. and Sinclair, A. Approximating the permanent. *SIAM journal on computing*, 18(6):1149–1178, 1989.
- Kannan, R., Vempala, S., and Vetta, A. On clusterings: Good, bad and spectral. *J. ACM*, 51(3):497–515, May 2004. ISSN 0004-5411.
- Koltchinskii, V. and Gine, E. Random matrix approximation of spectra of integral operators. *Bernoulli*, 6(1):113–167, 02 2000.
- Lei, J. and Rinaldo, A. Consistency of spectral clustering in stochastic block models. *Ann. Statist.*, 43(1):215–237, 02 2015.
- Leskovec, J., Lang, K. J., and Mahoney, M. Empirical comparison of algorithms for network community detection. In *Proceedings of the 19th International Conference on World Wide Web*, 2010.
- Mahoney, M. W., Orecchia, L., and Vishnoi, N. K. A local spectral method for graphs: with applications to improving graph partitions and exploring data graphs locally. *Journal of Machine Learning Research*, 13:2339–2365, 2012.
- Maier, M., von Luxburg, U., and Hein, M. How the result of graph clustering methods depends on the construction of the graph. *CoRR*, abs/1102.2075, 2011.
- McSherry, F. Spectral partitioning of random graphs. In *FOCS*, pp. 529–537, 2001.
- Rohe, K., Chatterjee, S., and Yu, B. Spectral clustering and the high-dimensional stochastic blockmodel. *Ann. Statist.*, 39(4):1878–1915, 08 2011.
- Schiebinger, G., Wainwright, M. J., and Yu, B. The geometry of kernelized spectral clustering. *Ann. Statist.*, 43(2): 819–846, 04 2015.
- Shi, T., Belkin, M., and Yu, B. Data spectroscopy: Eigenspaces of convolution operators and clustering. *Ann. Statist.*, 37(6B):3960–3984, 12 2009.
- Singh, A., Scott, C., and Nowak, R. Adaptive hausdorff estimation of density level sets. *Ann. Statist.*, 37(5B): 2760–2782, 10 2009.
- Spielman, D. A. and Teng, S.-H. Spectral sparsification of graphs. *SIAM Journal on Computing*, 40(4):981–1025, 2011.
- Spielman, D. A. and Teng, S.-H. A local clustering algorithm for massive graphs and its application to nearly linear time graph partitioning. *SIAM Journal on Computing*, 42(1):1–26, 2013.

- Spielman, D. A. and Teng, S.-H. Nearly linear time algorithms for preconditioning and solving symmetric, diagonally dominant linear systems. *SIAM Journal on Matrix Analysis and Applications*, 35(3):835–885, 2014.
- Tolliver, D. and Miller, G. L. Graph partitioning by spectral rounding: Applications in image segmentation and clustering. In *Computer Vision and Pattern Recognition, CVPR*, volume 1, pp. 1053–1060, 2006.
- Trillos, N. G., Slepčev, D., Von Brecht, J., Laurent, T., and Bresson, X. Consistency of cheeger and ratio graph cuts. *Journal of Machine Learning Research*, 17(1):6268–6313, January 2016.
- Vempala, S. and Wang, G. A spectral algorithm for learning mixture models. *Journal of Computer and System Sciences*, 68(4):841 – 860, 2004.
- von Luxburg, U. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, December 2007.
- von Luxburg, U., Belkin, M., and Bousquet, O. Consistency of spectral clustering. *Ann. Statist.*, 36(2):555–586, 04 2008.
- Yang, J. and Leskovec, J. Defining and evaluating network communities based on ground-truth. *Knowledge and Information Systems*, 42(1):181–213, Jan 2015.
- Zhu, Z. A., Lattanzi, S., and Mirrokni, V. S. A local algorithm for finding well-connected clusters. In *ICML (3)*, pp. 396–404, 2013.