
Local Spectral Clustering of Density Upper Level Sets

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

13 1 Introduction

14 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
15 consider the problem of clustering: splitting the data into groups which satisfy some notion of
16 within-group similarity and between-group difference. We focus on spectral clustering methods, a
17 family of powerful nonparametric clustering algorithms. Roughly speaking, a spectral technique first
18 constructs a geometric graph G , where vertices are associated with samples, and edges correspond
19 to proximities between samples. It then learns a feature embedding based on the Laplacian of G ,
20 and applies a simple clustering technique (such as k-means clustering) in the embedded feature
21 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

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

22 where $\mathbf{I} \in \mathbb{R}^{n \times n}$ is the identity matrix. Classical global spectral methods take a eigendecomposition
23 $\mathbf{L} = \mathbf{U} \Sigma \mathbf{U}^T$, use some number of eigenvectors (columns in \mathbf{U}) as a feature representation for the
24 samples, and then run (say) k-means in this new feature space.

25 When applied to geometric graphs constructed from a large number of samples, global spectral
26 clustering methods can be computationally cumbersome and insensitive to the local geometry of
27 the underlying distribution [Leskovec et al., 2010, Mahoney et al., 2012]. This has led to recent
28 increased interest in local spectral algorithms, which leverage locally-biased spectra computed using

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

random walks around a user-specified seed node. A popular local clustering algorithm is Personalized PageRank (PPR), first introduced by [Haveliwala, 2003], and further developed by [Spielman and 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, Andersen et al., 2012, Gleich and Seshadhri, 2012, Mahoney et al., 2012, Wu et al., 2012], which has led many authors to develop supporting theory [Spielman and Teng, 2013, Andersen and Peres, 2009, Gharan and Trevisan, 2012, Zhu et al., 2013] that gives 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 $\{x \in \mathbb{R}^d : f(x) \geq \lambda\}$ for some threshold $\lambda > 0$ (a central object of central interest in the classical statistical literature on clustering, dating back to Hartigan [1981]).

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 and 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

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

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 the *relative pointwise mixing time* of $G[S]$ as

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

Now our internal connectivity parameter is simply the inverse mixing time,

$$\Psi(S; G) = \frac{1}{\tau_\infty(G[S])}. \quad (2)$$

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

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

60 1.2 PPR on a Neighborhood Graph

61 We now describe the clustering algorithm that will be our focus for the rest of the paper. We
 62 start with the geometric graph that we form based on the samples \mathbf{X} : for a radius $r > 0$, we
 63 consider the r -neighborhood graph of \mathbf{X} , denoted $G_{n,r} = (V, E)$, an unweighted graph with
 64 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
 65 Euclidean norm. Note that this is a special case of the general construction introduced above, with
 66 $K(u, v) = 1(\|x_u - x_v\| \leq r)$.

67 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
 68 teleportation parameter $\alpha \in [0, 1]$, to be the solution of the following linear system:

$$\mathbf{p} = \alpha \mathbf{e}_v + (1 - \alpha) \mathbf{p} \mathbf{W}, \quad (3)$$

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

72 For a level $\beta > 0$ and a target volume $\text{vol}_0 > 0$, we define a β -sweep cut of \mathbf{p} as

$$S_\beta = \{u \in V : \frac{p_u}{\mathbf{D}_{uu}} > \frac{\beta}{\text{vol}_0}\}. \quad (4)$$

73 Having computed sweep cuts over a range $\beta \in (\frac{1}{40}, \frac{1}{11})$,⁴ we output a cluster $\hat{C} = S_{\beta^*}$, based on the
 74 sweep cut S_{β^*} that minimizes the normalized cut $\Phi(S_{\beta^*}; G_{n,r})$ as defined in (1). For concreteness,
 75 we summarize this procedure in Algorithm 1.

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
 stationary volume $\pi_0 > 0$.

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

- 1: Form the neighborhood graph $G_{n,r}$.
- 2: Compute the PPR vector $\mathbf{p}(v, \alpha; G_{n,r})$ as in (3).
- 3: For $\beta \in (\frac{1}{40}, \frac{1}{11})$ compute sweep cuts S_β as in (4).
- 4: Return $\hat{C} = S_{\beta^*}$, where

$$\beta^* = \arg \min_{\beta \in (\frac{1}{40}, \frac{1}{11})} \Phi(S_\beta; G_{n,r}).$$

76 1.3 Summary of Results

77 Let $\mathbb{C}_f(\lambda)$ denote the connected components of the density upper level set $\{x \in \mathbb{R}^d : f(x) > \lambda\}$.
 78 For a given density cluster $\mathcal{C} \in \mathbb{C}_f(\lambda)$, we call $\mathcal{C}[\mathbf{X}] = \mathcal{C} \cap \mathbf{X}$ the *empirical density cluster*. Below
 79 we give two notions of performance of a density cluster estimate.

80 **Definition 1** (Misclassification rate). *For an estimator $\hat{C} \subseteq \mathbf{X}$ and set $\mathcal{S} \subseteq \mathbb{R}^d$, the misclassification*
 81 *rate of \mathcal{S} by \hat{C} is*

$$|\hat{C} \setminus (\mathcal{S} \cap \mathbf{X})| + |(\mathcal{S} \cap \mathbf{X}) \setminus \hat{C}|. \quad (5)$$

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 $(\frac{1}{40}, \frac{1}{11})$ is standard in the analysis of PPR algorithms, see, e.g., [Zhu et al., 2013].

Definition 2 (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 a (thickened) empirical density cluster $\mathcal{C}_\sigma[\mathbf{X}]$, under natural geometric conditions (precluding clusters that are arbitrarily thin).
2. Under **largely the same set of geometric conditions**, we derive in Theorem ?? a lower bound on the inverse mixing time of a random walk over $\mathcal{C}_\sigma[\mathbf{X}]$.
3. In Section 3, we show in Theorem 4 that the bounds on the cluster quality criteria established in Theorems 1 and ?? have algorithmic consequences for PPR: properly initialized, Algorithm 1 has low misclassification rate with respect to a small enlargement of the set \mathcal{C} , and if the density cluster \mathcal{C} is particularly well-conditioned, Algorithm 1 will perform consistent density cluster estimation in the sense of (6).
4. We show in Corollary 1 that these statements hold also with respect to an approximate form of PPR, which can be efficiently computed.
5. In Section 4, we empirically demonstrate the tightness of the bounds in Theorems 1 and ??, and provide examples showing how violations of the geometric conditions we require manifestly impact density cluster recovery by PPR.

On the topic of conditions, it is worth mentioning that, as density clusters are inherently local, focusing on the PPR algorithm 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 and Hoffman, 1973] and their performance is well-understood in this context (see, e.g., Tolliver and 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 and Rinaldo, 2015, Abbe, 2018] have studied the efficacy of spectral methods in successfully recovering the community structure in the stochastic block model and variants.

Building on earlier work of [Koltchinskii and Gine, 2000], [von Luxburg et al., 2008, Hein et al., 2005] 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 and 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 Bounds for Density Clusters

2.1 Geometric Conditions on Density Clusters

In order to provide meaningful bounds on the normalized cut and inverse mixing time of an empirical density cluster, we must introduce conditions on the density f . Let $B(x, r) = \{y \in \mathbb{R}^d : \|y - x\| \leq$

129 $r\}$ be the closed ball of radius $r > 0$, centered at $x \in \mathbb{R}^d$. Given a set $\mathcal{A} \subseteq \mathbb{R}^d$ and $\sigma > 0$, define
 130 $\mathcal{A}_\sigma = \mathcal{A} + B(0, \sigma) = \{y \in \mathbb{R}^d : \inf_{x \in \mathcal{A}} \|y - x\| \leq \sigma\}$, which we call the σ -expansion of \mathcal{A} .
 131 For a differentiable function $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$, write $\nabla g(x)$ to denote the Jacobian of g evaluated at
 132 $x \in \mathbb{R}^d$.

133 We are now ready to give our required conditions, stated with respect to a density cluster $\mathcal{C} \in \mathbb{C}_f(\lambda)$
 134 for some threshold $\lambda > 0$, and an expansion parameter $\sigma > 0$.

(A1) *Bounded density within cluster*: There are $0 < \lambda_\sigma < \Lambda_\sigma < \infty$ such that

$$\lambda_\sigma = \inf_{x \in \mathcal{C}_\sigma} f(x) \leq \sup_{x \in \mathcal{C}_\sigma} f(x) \leq \Lambda_\sigma.$$

(A2) *Low noise density*: There exists $\gamma, c_0 > 0$ such that for all $x \in \mathbb{R}^d$ with $0 < \text{dist}(x, \mathcal{C}_\sigma) \leq \sigma$,

$$\inf_{x' \in \mathcal{C}_\sigma} f(x') - f(x) \geq c_0 \text{dist}(x, \mathcal{C}_\sigma)^\gamma,$$

135 where $\text{dist}(x, \mathcal{A}) = \inf_{x_0 \in \mathcal{A}} \|x - x_0\|$.

(A3) *Cluster separation*: For all $\mathcal{C}' \in \mathbb{C}_f(\lambda)$ with $\mathcal{C}' \neq \mathcal{C}$,

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

136 where $\text{dist}(\mathcal{A}, \mathcal{A}') = \inf_{x \in \mathcal{A}} \text{dist}(x, \mathcal{A}')$.

137 (A4) *Lipschitz embedding*: There exists $\mathcal{K} \subseteq \mathbb{R}^d$ convex, and $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ satisfying

$$\det(\nabla g(x)) = 1, \frac{1}{L} \|x - y\| \leq \|g(x) - g(y)\| \leq L \|x - y\| \text{ for all } x, y \in \mathbb{R}^d$$

such that \mathcal{C}_σ is the image of \mathcal{K} by g , $\mathcal{C}_\sigma = g(\mathcal{K})$. Furthermore, there exists $D < \infty$ such
 that for all $x, x' \in \mathcal{K}$

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

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

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

146 We also note that while many of these geometric conditions are typical in the density clustering
 147 literature, the restrictions we will impose upon them in order to obtain meaningful implications for
 148 PPR will not be. This is natural. The spectral algorithm we consider is not specifically designed for
 149 the task of level set estimation, and in fact one should expect PPR to fail to recover – either in the
 150 sense of (6), or indeed any reasonable notion of cluster recovery – a density cluster of sufficiently
 151 large diameter or sufficiently small thickness (though we do not provide any lower bounds to this
 152 effect). Indeed, one of the primary motivations of this work was to better understand and characterize
 153 the distinctions between those level sets which are well conditioned for spectral algorithms, and those
 154 which are not.

155 In the next several subsections, we will derive bounds on the cluster quality criteria evaluated on
 156 $(\sigma$ -expansions of) density clusters. For notational simplicity, hereafter for $S \subseteq V$, we will abbreviate
 157 $\Phi(S; G_{n,r})$ by $\Phi_{n,r}(S)$, and similarly, $\Psi(S; G_{n,r})$ by $\Psi_{n,r}(S)$, and $\tau_\infty(G_{n,r}[S])$ by $\tau_{n,r}(S)$. We
 158 will also use ν for Lebesgue measure on \mathbb{R}^d , and $\nu_d = \nu(B)$ for the measure of the unit ball
 159 $B = B(0, 1)$.

160 2.2 Upper Bound on Normalized Cut

161 We start with an upper bound on the normalized cut (1) of $\mathcal{C}_\sigma[\mathbf{X}]$. (In Theorem 1, the upper bound on
 162 the density in Assumption (A1) will not actually be needed, so we omit the parameter $\Lambda_\sigma > 0$ from
 163 the theorem statement.) For $S \subseteq \mathbb{R}^d$ and $r > 0$, let

$$\pi_{\mathbb{P},r}(S) := \frac{\int_S \mathbb{P}(B(x, r)) f(x) dx}{\int_{\mathbb{R}^d} \mathbb{P}(B(x, r)) f(x) dx}.$$

Theorem 1. Fix $\lambda > 0$, and let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ satisfy Assumptions (A1), (A2), and (A3), for some $\sigma, \lambda_\sigma, c_0, \gamma > 0$. Let $0 < r \leq \sigma/4d$ be such that

$$\pi_{\mathbb{P},r}(\mathcal{C}_\sigma) \leq \frac{1}{2}. \quad (7)$$

Then for any $0 < \delta < 1$, $\epsilon > 0$, if

$$n \geq \frac{(2 + \epsilon)^2 \log(3/\delta)}{\epsilon^2} \left(\frac{25}{6\lambda_\sigma^2 \nu(\mathcal{C}_\sigma) \nu_d r^d} \right)^2, \quad (8)$$

then

$$\frac{\Phi_{n,r}(\mathcal{C}_\sigma[\mathbf{X}])}{r} \leq c \frac{d}{\sigma} \frac{\lambda}{\lambda_\sigma} \frac{(\lambda_\sigma - c_0 \frac{r^\gamma}{\gamma+1})}{\lambda_\sigma} + \epsilon, \quad (9)$$

with probability at least $1 - \delta$ (where $c_1 > 0$ is a universal constant).

Remark 1. The proof of Theorem 1, along with all other proofs in this paper, can be found in the supplementary document. The key idea is that for any $x \in \mathcal{C}$, the simple fact $B(x, \sigma) \subseteq \mathcal{C}_\sigma$ translates into the upper bound $\nu(\mathcal{C}_\sigma + rB) \leq (1 + 2dr/\sigma)\nu(\mathcal{C}_\sigma)$. We leverage (A2) to find a corresponding bound on the weighted volume, then apply standard concentration inequalities to convert from population- to sample-based results.

Remark 2. The inequality in (9) is tight in the case of $\mathcal{C} = \{0\}$. To see this, let $\mathcal{C}_\sigma = B(0, \sigma)$ and

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

Then, some simple calculations yield

$$\frac{\mathbb{P}((\mathcal{C}_\sigma + B(0, r)) \setminus \mathcal{C}_\sigma)}{\mathbb{P}(\mathcal{C}_\sigma)} \geq d \frac{(\lambda - r^\gamma)r}{\sigma},$$

and for some $c > 0$,

$$\mathbb{E}(\text{cut}(\mathcal{C}_\sigma[\mathbf{X}]; G_{n,r}) \setminus \mathcal{C}_\sigma) \geq c\lambda\nu_d r^d \mathbb{P}((\mathcal{C}_\sigma + B(0, r)) \setminus \mathcal{C}_\sigma), \quad \text{and} \quad \mathbb{E}(\text{vol}_{n,r}(\mathcal{C}_\sigma[\mathbf{X}]; G_{n,r})) \leq c\lambda\nu_d r^d \mathbb{P}(\mathcal{C}_\sigma)$$

so the ratio $\mathbb{E}(\text{cut}(\mathcal{C}_\sigma[\mathbf{X}]; G_{n,r}))/\mathbb{E}(\text{vol}_{n,r}(\mathcal{C}_\sigma[\mathbf{X}]; G_{n,r}))$ matches 9, up to constants.

2.3 Lower Bound on Inverse Mixing Time

Next we lower bound the inverse mixing time (2) of $\mathcal{C}_\sigma[\mathbf{X}]$, or equivalently, as $\Psi_{n,r}(\mathcal{C}_\sigma[\mathbf{X}]) = 1/\tau_{n,r}(\mathcal{C}_\sigma[\mathbf{X}])$, we upper bound the mixing time.

Theorem 2. Fix $\lambda > 0$, and let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ satisfy Assumptions (A1) and (A4) for some $\sigma, \lambda_\sigma, \Lambda_\sigma, D, K > 0$. Then, for any $0 < r < \sigma/2\sqrt{d}$, with probability one

$$\limsup_{n \rightarrow \infty} \tau_{n,r}(\mathcal{C}_\sigma[\mathbf{X}]) \leq c_2 \frac{\Lambda_\sigma^4 d^3 D^2 L^2}{\lambda_\sigma^4 r^2} \log^2 \left(\frac{1}{r} \right) + c_3 \log \left(\frac{\Lambda_\sigma}{\lambda_\sigma} \right)$$

for $c_2, c_3 > 0$ universal constants.

We will outline the main techniques used to prove Theorem 2, but first, we show that a straightforward and perhaps immediately obvious route to upper bounding the mixing time is insufficient for our purposes. Let $\tilde{\Phi}_{n,r}$ be the *conductance* of the induced subgraph $\tilde{G}_{n,r} = G_{n,r}[\mathcal{C}_\sigma[\mathbf{X}]]$, given by

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

and $\tilde{\mu}_2$ be the 2nd largest eigenvalue of $\mathbf{W}_{\mathcal{C}_\sigma[\mathbf{X}]}$, the random walk matrix over $\tilde{G}_{n,r}$. The famous Cheeger's inequality

$$\frac{\tilde{\Phi}_{n,r}^2}{2} \leq 1 - \tilde{\mu}_2 \leq \tilde{\Phi}_{n,r}$$

189 applied with well known theory on ∞ -mixing time yields

$$\left| \tilde{\pi}_{n,r}(u) - e_v \mathbf{W}_{\mathcal{C}_\sigma[\mathbf{X}]}^t e_u \right| \leq \left(1 - \frac{\tilde{\Phi}_{n,r}^2}{2} \right)^t \sqrt{\frac{\deg(u; \tilde{G}_{n,r})}{\min_{v \in \mathcal{C}_\sigma[\mathbf{X}]} \deg(v; \tilde{G}_{n,r})}} \quad (u, v \in \mathcal{C}_\sigma[\mathbf{X}])$$

190 where $\tilde{\pi}_{n,r}$ is the stationary distribution over $\tilde{G}_{n,r}$. Unfortunately, even if $\tilde{G}_{n,r}$ is an approximately
 191 degree-regular graph, (meaning we can ignore the term under the square root), this bound only yields
 192 ⁵

$$\tau_{n,r}(\mathcal{C}_\sigma[\mathbf{X}]) = \mathcal{O} \left(\frac{\log \left(\frac{1}{\tilde{\pi}_{\min}} \right)}{\tilde{\Phi}_{n,r}^2} \right).$$

193 where $\tilde{\pi}_{\min} = \min_{u \in \mathcal{C}_\sigma[\mathbf{X}]} \tilde{\pi}_{n,r}(u)$. As $\tilde{\pi}_{\min} = \mathcal{O}(1/n)$, the upper bound on the right hand side will
 194 be of order at least $\log(n)$, and the bound becomes trivial as $n \rightarrow \infty$.

195 **We therefore apply a more complex approach.**

196 3 Consistent Cluster Estimation

197 3.1 Well-Conditioned Density Clusters

198 For PPR to accurately estimate a set, the ratio of normalized cut to inverse mixing time should be
 199 small. Letting $\theta := (r, \sigma, \lambda, \lambda_\sigma, \Lambda_\sigma, \gamma, D, L)$ contain those parameters which govern the bounds
 200 given in Theorems 1 and 2, further abbreviate

$$\begin{aligned} \Phi(\theta) &:= c_1 r \frac{d}{\sigma} \frac{\lambda}{\lambda_\sigma} \frac{(\lambda_\sigma - c_0 \frac{r^\gamma}{\gamma+1})}{\lambda_\sigma} \\ \Psi(\theta) &:= \left(c_2 \frac{\Lambda_\sigma^4 d^3 D^2 L^2}{\lambda_\sigma^4 r^2} \log^2 \left(\frac{1}{r} \right) + c_3 \log \left(\frac{\Lambda_\sigma}{\lambda_\sigma} \right) \right)^{-1} \end{aligned}$$

201 for these bounds (where all constants $c_0, c_1, c_2, c_3 > 0$ are as in these theorems).

202 Well-conditioned density clusters satisfy all of the given assumptions, for parameters which results in
 203 ‘good’ values of $\Phi(\theta)$ and $\Psi(\theta)$.

204 **Definition 3** (Well-conditioned density clusters). *For $\lambda > 0$ and $\mathcal{C} \in \mathbb{C}_f(\lambda)$, let \mathcal{C} satisfy (A1) - (A4)
 205 for some $\sigma, \lambda, \lambda_\sigma, \Lambda_\sigma, \gamma, D, L > 0$, and additionally let \mathcal{C}_σ satisfy (7). Then, setting*

$$\kappa(\mathcal{C}) := \frac{\Phi(\theta)}{\Psi(\theta)}$$

206 *we call \mathcal{C} a κ -well-conditioned density cluster (with respect to θ).*

207 We focusing for a moment on the neighborhood graph radius r . While taking $r \rightarrow 0$ as $n \rightarrow \infty$ – and
 208 thereby ensuring $G_{n,r}$ is sparse – is computationally attractive, the presence of a factor of $\frac{1}{r}$ in $\kappa(\mathcal{C})$
 209 unfortunately prevents us from making claims about the behavior of PPR for this regime. While the
 210 restriction to a kernel function fixed in n is standard for theoretical analysis Schiebinger et al. [2015],
 211 von Luxburg et al. [2008], it is an interesting question whether PPR exhibits some degeneracy over
 212 r -neighborhood graphs as $r \rightarrow 0$, or if this is merely looseness in our upper bounds.

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

⁵For sequences of numbers $\{a_n\}_{n \in \mathbb{N}}, \{b_n\}_{n \in \mathbb{N}}$, we say $a_n = \mathcal{O}(b_n)$ when there exists c sufficiently large such that $a_n \leq cb_n$ for all n .

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

$$\alpha \in [1/10, 1/9] \cdot \Psi(\theta), r \leq \frac{\sigma}{2d},$$

$$\text{vol}_0 \in [3/4, 5/4] \cdot n(n-1) \int_{\mathcal{C}_\sigma} \mathbb{P}(B(x, r)) f(x) dx, v \in \mathcal{C}_\sigma[\mathbf{X}]^g \quad (10)$$

218 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
 219 $\text{vol}(\mathcal{C}_\sigma[\mathbf{X}]^g; G_{n,r}) \geq \text{vol}(\mathcal{C}_\sigma[\mathbf{X}]; G_{n,r})/2$. (Intuitively one can think of $\mathcal{C}_\sigma[\mathbf{X}]^g$ as consisting of
 220 the nodes sufficiently close to the center of $\mathcal{C}_\sigma[\mathbf{X}]$, although we provide no formal justification to this
 221 effect.)

222 **Definition 4.** If the input parameters to Algorithm 1 satisfy (10) with respect to some $\mathcal{C}_\sigma[\mathbf{X}]$, we say
 223 the algorithm is well-initialized.

224 In practice, a reasonable way to choose these hyperparameters is by tuning. For example, if one
 225 wanted to successfully recover a density cluster, one could vary each hyperparameter over a grid,
 226 retaining outputs $\hat{\mathcal{C}}$ of Algorithm 1 only if they recover $\mathcal{C}[\mathbf{X}]$ for some $\mathcal{C} \in \mathbb{C}_f(\lambda)$ and discarding
 227 them otherwise. Then simply return the minimum normalized cut set from those $\hat{\mathcal{C}}$ which were
 228 retained. Assuming there existed $\lambda > 0, \mathcal{C} \in \mathbb{C}_f(\lambda)$ such that \mathcal{C} satisfied the conditions of Theorem
 229 4, and moreover some combination of tuning parameters in the chosen grid satisfied (10), this scheme
 230 would inherit the consistency guarantees of Theorem 4.

231 **Misclassification rates for PPR.** In Zhu et al. [2013], building on the work of Andersen et al.
 232 [2006] and others, theory is developed which links algorithmic performance of PPR to the normalized
 233 cut and mixing time parameters. This work, combined with the results of Section 2, immediately
 234 implies a bound on the volume of $\hat{\mathcal{C}} \setminus \mathcal{C}_\sigma[\mathbf{X}]$ (and likewise $\mathcal{C}_\sigma[\mathbf{X}] \setminus \hat{\mathcal{C}}$),

$$\frac{\text{vol}_{n,r}(\hat{\mathcal{C}} \setminus \mathcal{C}_\sigma[\mathbf{X}])}{\text{vol}_{n,r}(\mathcal{C}_\sigma[\mathbf{X}])}, \frac{\text{vol}_{n,r}(\mathcal{C}_\sigma[\mathbf{X}] \setminus \hat{\mathcal{C}})}{\text{vol}_{n,r}(\mathcal{C}_\sigma[\mathbf{X}])} = \mathcal{O}(\kappa(\mathcal{C})). \quad (11)$$

235 To translate (11) into meaningful bounds on misclassification rate, we wish to preclude vertices
 236 $x \in \mathbf{X}$ from having arbitrarily small degree. To do so, we make some regularity assumptions on
 237 $\mathcal{X} = \text{supp}(f)$.

238 (A5) *Valid region:* $0 < \lambda_{\min} < f(x)$ for all $x \in \mathcal{X}$. Additionally, there exists some $c > 0$ such
 239 that for each $x \in \partial\mathcal{X}$, $\nu(B(x, r) \cap \mathcal{X}) \geq c\nu(B(x, r))$.

240 Note that the latter condition in (A5) will be satisfied if, for instance, \mathcal{X} is a σ -expansion.

241 **Theorem 3.** Fix $\lambda > 0$, let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ be a κ -well conditioned density cluster (with respect to some
 242 θ), and additionally assume f satisfies (A5). Then, with probability tending to one as $n \rightarrow \infty$,

$$\frac{|\mathcal{C}_\sigma[\mathbf{X}] \setminus \hat{\mathcal{C}}|}{|\mathcal{C}_\sigma[\mathbf{X}]|} \leq c_5 \kappa(\mathcal{C}) \frac{\Lambda_\sigma}{\lambda_\sigma} \quad (12)$$

243 and

$$\frac{|\hat{\mathcal{C}} \setminus \mathcal{C}_\sigma[\mathbf{X}]|}{|\mathcal{C}_\sigma[\mathbf{X}]|} \leq c_6 \kappa(\mathcal{C}) \frac{\Lambda_\sigma}{\lambda_{\min}}. \quad (13)$$

244 for universal constants $c_5, c_6 > 0$.

245 **Remark 3.** A notable implication of our theory is that as the diameter D increases, our upper bounded
 246 on the normalized cut of \mathcal{C}_σ remains unchanged, but $\kappa(\mathcal{C})$, and therefore the misclassification rate,
 247 worsens (increases). This phenomenon reflects established wisdom regarding spectral partitioning
 248 algorithms more generally Guattery and Miller [1995], Hein and Bühler [2010], albeit newly applied
 249 to the density clustering setting. It suggests that PPR may fail to recover $\mathcal{C}_\sigma[\mathbf{X}]$ even when \mathcal{C} is
 250 sufficiently well-conditioned to ensure $\mathcal{C}_\sigma[\mathbf{X}]$ has a small normalized cut in $G_{n,r}$. This intuition will
 251 be supported by simulations in Section 4.

252 **Consistent density cluster estimation.** Neither (11) nor Theorem 3 imply consistent density
 253 cluster estimation in the sense of (6). This notion of consistency requires a uniform bound over \mathbf{p} for
 254 all $u \in \mathcal{C}, u' \in \mathcal{C}'$, of the form

$$\frac{\mathbf{p}(u')}{\mathbf{D}_{uu}} \leq \frac{c}{\text{vol}_0} < \frac{C}{\text{vol}_0} \leq \frac{\mathbf{p}(u)}{\mathbf{D}_{uu}}. \quad (14)$$

255 so that a sweep cut S_β for $\beta \text{vol}_0 \in [c, C]$ would fulfill both conditions laid out in (6). Theorem 4
 256 provides sufficient conditions under which this gap will appear.

257 As before, we wish to preclude arbitrarily low degree vertices, this time for points $x \in \mathcal{C}'[\mathbf{X}]$.

258 (A6) \mathcal{C}' -bounded density : For each $\mathcal{C}' \in \mathbb{C}_f(\lambda), \mathcal{C}' \neq \mathcal{C}$, for all $x \in \mathcal{C}' + \sigma B$,

$$\lambda_\sigma \leq f(x)$$

259 where σ, λ_σ are as in (A1).

260 **Theorem 4.** Fix $\lambda > 0$, let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ be a κ -well conditioned cluster (with respect to some θ), and
 261 additionally assume (A6) holds. If Algorithm 1 is well-initialized, there exists universal constant
 262 $c_7 > 0$ such that if

$$\kappa(\mathcal{C}) \leq c_7 \frac{\lambda_\sigma^2 r^d \nu_d}{\Lambda_\sigma \mathbb{P}(\mathcal{C}_\sigma)}, \quad (15)$$

263 then the output set $\hat{\mathcal{C}} \subseteq \mathbf{X}$ is a consistent estimator for \mathcal{C} , in the sense of Definition 2.

264 **Cluster estimation with the approximate PPR vector.** As mentioned previously, in practice
 265 solving (3) may be too computationally expensive. To address this limitation, Andersen et al. [2006]
 266 introduced the ϵ -approximate PPR vector (aPPR), which we will denote $\mathbf{p}^{(\epsilon)}$. We refer the curious
 267 reader to Andersen et al. [2006] for a formal algorithmic definition of the aPPR vector, and limit
 268 ourselves to highlighting a few salient points. Namely, the aPPR vector can be computed in $\mathcal{O}(\frac{1}{\epsilon\alpha})$
 269 time, while satisfying the following uniform error bound:

$$\text{for all } x \in \mathbf{X}, \quad \mathbf{p}(x) - \epsilon \deg_{n,r}(x) \leq \mathbf{p}^{(\epsilon)}(x) \leq \mathbf{p}(x) \quad (16)$$

270 Application of (16) within the proofs of Theorems 3 and 4 leads to analogous results which hold with
 271 respect to $\mathbf{p}^{(\epsilon)}$.

272 **Corollary 1.** Fix $\lambda > 0$, let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ be a κ -well-conditioned cluster (with respect to some θ).
 273 Choose input parameters $\alpha, r, \text{vol}_0, v$ to be well-initialized in the sense of (10). Set $\epsilon = \frac{1}{20\text{vol}_0}$, and
 274 modify Algorithm 1 to compute the aPPR vector $\mathbf{p}^{(\epsilon)}$ rather than the exact PPR vector \mathbf{p} , with
 275 resulting output $\hat{\mathcal{C}}$.

276 1. Assume (A5) holds. Then (12) and (13) are still valid upper bounds.

277 2. Assume (A6) holds. If

$$\kappa(\mathcal{C}) \leq c_7 \frac{\lambda_\sigma^2 r^d \nu_d}{\Lambda_\sigma^2 \nu(\mathcal{C}_\sigma)}$$

278 then $\hat{\mathcal{C}} \subseteq \mathbf{X}$ is a consistent estimator for \mathcal{C} , in the sense of Definition 2.

279 4 Experiments

280 4.1 Validating Theoretical Bounds

281 As we do not provide any theoretical lower bounds, we validate the tightness of Theorems 1 and 2
 282 via simulation. We sample points according to the density function q , where for $x \in \mathbb{R}^d$

$$q(x) := \begin{cases} \lambda, & x \in [0, \sigma] \times D^{d-1} =: \mathcal{C}, \\ \lambda - \text{dist}(x, \mathcal{C})\eta, & x \in \mathcal{C}_\sigma \setminus \mathcal{C}, \\ (\lambda - \sigma\eta) - \text{dist}(x, \mathcal{C}_\sigma)^\gamma, & x \in (\mathcal{C}_\sigma + \sigma B) \setminus \mathcal{C}_\sigma, \\ 0, & \text{otherwise,} \end{cases} \quad (17)$$

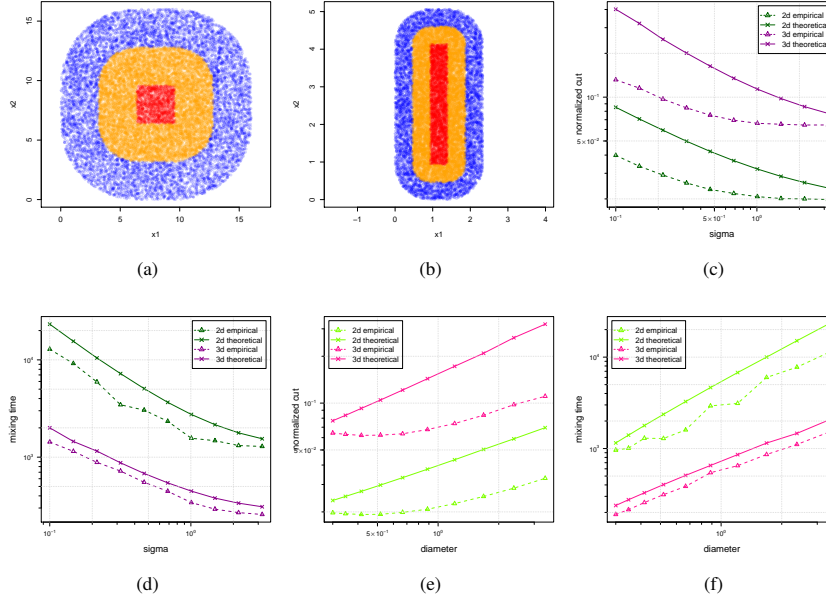


Figure 1

where $\lambda = \frac{150}{81}\sigma^\gamma$ and $\eta = \frac{15}{81}\sigma^{\gamma-1}$. Panels (a) and (b) in Figure 1 show 20,000 samples from two parameterizations of q . In (a), $\sigma = D = 3.2$, while in (b) $\sigma = .1$ and $D = 3.2$. (For both, $d = 2$).

Panels (c) – (f) in Figure 1 show the change in normalized cut and mixing time, respectively, as the parameters σ ((c) and (d)) and D ((e) and (f)) are varied. In panels (c) and (d) $\sigma = .1 \cdot \sqrt{2}^j$, $j = 1, \dots, 10$, and D is fixed at 3.2. In panels (e) and (f), $D = .1 \cdot \sqrt{2}^j$, $j = 1, \dots, 10$ and σ is fixed at .1. For each panel, the solid lines show, up to constants⁶, the theoretical upper bound, given by Theorem 1 for panels (c) and (e) and Theorem 2 for panels (d) and (f). The dashed lines show the computed empirical value, averaged over m trials ($m = 100$ for the normalized cut, dashed lines in panels (c) and (e), and $m = 20$ for the mixing time, dashed lines in panel (d) and (f)). For each trial across all parameters, r , the neighborhood graph radius, is set throughout to be as small as possible such that the resulting graph is connected, for computational efficiency. Green lines correspond to dimension $d = 2$, whereas purple/pink lines correspond to $d = 3$.

Panels (d) and (f) show the solid lines tracking closely to the dashed lines, in both 2 and 3 dimensions. This provides empirical evidence that the upper bound on mixing time given by Theorem 2 has the right dependency on both thickness parameter σ and diameter D .

The story in panels (c) and (e) is less obvious. We note that while, broadly speaking, the trends do not appear to match, this gap between theory and empirical results seems largest when $\sigma \approx D$; this is the right hand side on panel (c) and the left hand side on panel (e). It is in these regions that the slopes of the dashed and solid lines appear most distinct. As the ratio D/σ grows, we see the slopes of the empirical curves becoming more similar to those predicted by theory. The takeaway message is that while the dependency in (9) on σ and D is loose for clusters with diameter close to thickness, it becomes tighter as D/σ grows.

4.2 Empirical PPR, normalized cut, and density clustering comparison

To drive home the main implications of Theorems 3 and 4, we show the behavior of PPR, normalized cut, and the density clustering algorithm of [Chaudhuri and Dasgupta, 2010] on (a variant of) the famous ‘two moons’ dataset, considered a prototypical success story for spectral clustering algorithms.

⁶Note that we have rescaled all values of theoretical upper bounds by a constant, in order to mask the effect of large universal constants in these bounds. Therefore only comparison of slopes, rather than intercepts, is meaningful.

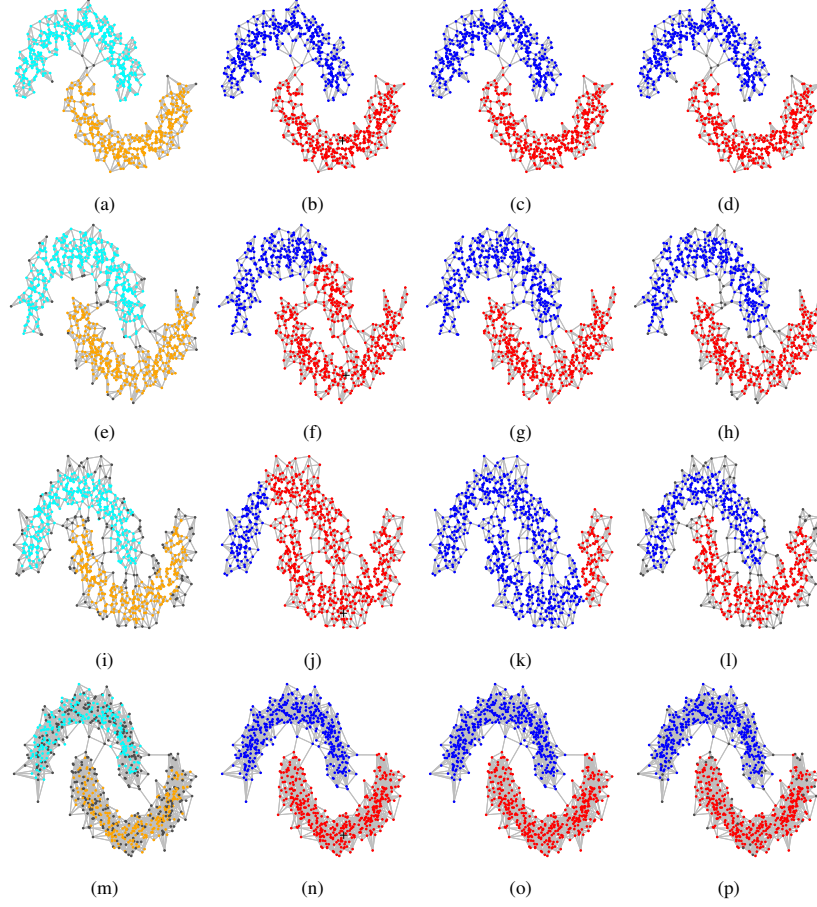


Figure 2: True density (column 1), PPR (column 2), normalized cut (column 3) and estimated density (column 4) clusters for 4 different simulated data sets. Seed node for PPR denoted by a black cross.

To form each of the four rows in Figure 2, 800 points are independently sampled following a ‘two moons plus Gaussian noise model’. Formally, the (respective) generative models for the data are

$$Z \sim \text{Bern}(1/2), \theta \sim \text{Unif}(0, \pi) \quad (18)$$

$$X(Z, \theta) = \begin{cases} \mu_1 + (r \cos(\theta), r \sin(\theta)) + \sigma \epsilon, & \text{if } Z = 1 \\ \mu_2 + (r \cos(\theta), -r \sin(\theta)) + \sigma \epsilon, & \text{if } Z = 0 \end{cases} \quad (19)$$

where

$$\mu_1 = (-.5, 0), \mu_2 = (0, 0), \epsilon \sim N(0, I_2) \quad (\text{row 1})$$

$$\mu_1 = (-.5, -.07), \mu_2 = (0, .07), \epsilon \sim N(0, I_2) \quad (\text{row 2})$$

$$\mu_1 = (-.5, -.125), \mu_2 = (0, .125), \epsilon \sim N(0, I_2) \quad (\text{row 3})$$

$$\mu_1 = (-.5, -.025), \mu_2 = (0, .025), \epsilon \sim N(0, I_{10}) \quad (\text{row 4})$$

for I_d the $d \times d$ identity matrix. The first column consists of the empirical density clusters C_n and C'_n for a particular threshold λ of the density function; the second column shows the PPR plus minimum normalized sweep cut cluster, with hyperparameter α and all sweep cuts considered; the third column shows the global minimum normalized cut, computed according to the algorithm of Szlam and Bresson [2010]; and the last column shows a cut of the cluster tree estimator of Chaudhuri and Dasgupta [2010].

Rows 1-3 show the degrading ability of PPR to recover density clusters as the two moons become less salient. In the first row, the normalized cut conforms to the density cluster, and PPR recovers both.

In the second row, the normalized cut still conforms to the density cluster, but because the internal connectivity of the lower moon is low, PPR fails to recover the normalized cut. In the third row, the moons have such low saliency that even the normalized cut fails to recover the lower moon; we also see from (k) that PPR does not somehow save us in this situation. Note that this is not a function of the finite sample: the 4th column shows us that a well-designed density clustering algorithm can recover the true density cluster.

The fourth row illustrates the effect of dimension. The gray dots in (m) (as in (a) , (e) and (i)) are observations in low-density regions. While the PPR sweep cut (n) has relatively high symmetric set difference with the chosen density cut, it still recovers C_n in the sense of Definition 2.

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 Figures ?? and ?? that poorly conditioned density clusters are not consistently estimated by PPR.

The initial motivation for this article was based on the intuition that density level sets, in the right conditions, will have small normalized cut. As a result, algorithms with normalized cut based guarantees (such as PPR) seemed likely to have density cluster recovery guarantees as well. However, the second-order behavior of PPR when failing to recover the conductance cut is also of interest. Are there situations when the conductance cut and density cut differ, yet PPR still recovers the latter? This is an open question.

References

- Emmanuel Abbe. Community detection and stochastic block models. *Foundations and Trends® in Communications and Information Theory*, 14(1-2):1–162, 2018.
- Reid Andersen and Yuval Peres. Finding sparse cuts locally using evolving sets. In *Proceedings of the Forty-first Annual ACM Symposium on Theory of Computing*, STOC '09, pages 235–244, New York, NY, USA, 2009. ACM. ISBN 978-1-60558-506-2. doi: 10.1145/1536414.1536449.
- Reid Andersen, Fan Chung, and Kevin Lang. Local graph partitioning using pagerank vectors. In *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science*, pages 475–486, 2006.
- Reid Andersen, David F Gleich, and Vahab Mirrokni. Overlapping clusters for distributed computation. In *Proceedings of the fifth ACM international conference on Web search and data mining*, pages 273–282. ACM, 2012.
- Sivaraman Balakrishnan, Min Xu, Akshay Krishnamurthy, and Aarti Singh. Noise thresholds for spectral clustering. In *Advances in Neural Information Processing Systems 24*. Curran Associates, Inc., 2011.
- Kamalika Chaudhuri and Sanjoy Dasgupta. Rates of convergence for the cluster tree. In J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems 23*, pages 343–351. Curran Associates, Inc., 2010.
- Kamalika Chaudhuri, Fan Chung Graham, and Alexander Tsiatas. Spectral clustering of graphs with general degrees in the extended planted partition model. In *COLT*, volume 23, pages 35.1–35.23, 2012.
- W. E. Donath and A. J. Hoffman. Lower bounds for the partitioning of graphs. *IBM J. Res. Dev.*, 17(5):420–425, September 1973.
- Miroslav Fiedler. Algebraic connectivity of graphs. *Czechoslovak Mathematical Journal*, 23(2): 298–305, 1973.
- Santo Fortunato. Community detection in graphs. *Physics Reports*, 486(3):75–174, 2010.
- Shayan Oveis Gharan and Luca Trevisan. Approximating the expansion profile and almost optimal local graph clustering. In *Foundations of Computer Science (FOCS), 2012 IEEE 53rd Annual Symposium on*, pages 187–196. IEEE, 2012.
- David F Gleich and C Seshadhri. Vertex neighborhoods, low conductance cuts, and good seeds for local community methods. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 597–605. ACM, 2012.
- Stephen Guattery and Gary L Miller. On the performance of spectral graph partitioning methods. In *SODA*, volume 95, pages 233–242, 1995.
- John A. Hartigan. Consistency of single-linkage for high-density clusters. *Journal of the American Statistical Association*, 1981.
- Taher H Haveliwala. Topic-sensitive pagerank: A context-sensitive ranking algorithm for web search. *IEEE transactions on knowledge and data engineering*, 15(4):784–796, 2003.
- Matthias Hein and Thomas Bühler. An inverse power method for nonlinear eigenproblems with applications in 1-spectral clustering and sparse pca. In *Advances in Neural Information Processing Systems*, pages 847–855, 2010.
- Matthias Hein, Jean-Yves Audibert, and Ulrike von Luxburg. From graphs to manifolds – weak and strong pointwise consistency of graph laplacians. In *Learning Theory*, 2005.
- Ravi Kannan, Santosh Vempala, and Adrian Vetta. On clusterings: Good, bad and spectral. *J. ACM*, 51(3):497–515, May 2004.

392 Vladimir Koltchinskii and Evarist Giné. Random matrix approximation of spectra of integral operators.
393 *Bernoulli*, 6(1):113–167, 02 2000.

394 Jing Lei and Alessandro Rinaldo. Consistency of spectral clustering in stochastic block models. *Ann.*
395 *Statist.*, 43(1):215–237, 02 2015.

396 Jure Leskovec, Kevin J. Lang, and Michael Mahoney. Empirical comparison of algorithms for
397 network community detection. In *Proceedings of the 19th International Conference on World Wide*
398 *Web*, 2010.

399 Michael W. Mahoney, Lorenzo Orecchia, and Nisheeth K. Vishnoi. A local spectral method for
400 graphs: with applications to improving graph partitions and exploring data graphs locally. *Journal*
401 *of Machine Learning Research*, 13:2339–2365, 2012.

402 Frank McSherry. Spectral partitioning of random graphs. In *FOCS*, pages 529–537, 2001.

403 Karl Rohe, Sourav Chatterjee, and Bin Yu. Spectral clustering and the high-dimensional stochastic
404 blockmodel. *Ann. Statist.*, 39(4):1878–1915, 08 2011.

405 Geoffrey Schiebinger, Martin J. Wainwright, and Bin Yu. The geometry of kernelized spectral
406 clustering. *Ann. Statist.*, 43(2):819–846, 04 2015.

407 Tao Shi, Mikhail Belkin, and Bin Yu. Data spectroscopy: Eigenspaces of convolution operators and
408 clustering. *Ann. Statist.*, 37(6B):3960–3984, 12 2009.

409 Aarti Singh, Clayton Scott, and Robert Nowak. Adaptive hausdorff estimation of density level sets.
410 *Ann. Statist.*, 37(5B):2760–2782, 10 2009.

411 Daniel A Spielman and Shang-Hua Teng. Spectral sparsification of graphs. *SIAM Journal on*
412 *Computing*, 40(4):981–1025, 2011.

413 Daniel A Spielman and Shang-Hua Teng. A local clustering algorithm for massive graphs and its
414 application to nearly linear time graph partitioning. *SIAM Journal on Computing*, 42(1):1–26,
415 2013.

416 Daniel A Spielman and Shang-Hua Teng. Nearly linear time algorithms for preconditioning and
417 solving symmetric, diagonally dominant linear systems. *SIAM Journal on Matrix Analysis and*
418 *Applications*, 35(3):835–885, 2014.

419 Arthur Szlam and Xavier Bresson. Total variation, cheeger cuts. In *ICML*, pages 1039–1046, 2010.

420 David Tolliver and Gary L. Miller. Graph partitioning by spectral rounding: Applications in image
421 segmentation and clustering. In *Computer Vision and Pattern Recognition, CVPR*, volume 1, pages
422 1053–1060, 2006.

423 Santosh Vempala and Grant Wang. A spectral algorithm for learning mixture models. *Journal of*
424 *Computer and System Sciences*, 68(4):841 – 860, 2004.

425 Ulrike von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416,
426 December 2007.

427 Ulrike von Luxburg, Mikhail Belkin, and Olivier Bousquet. Consistency of spectral clustering. *Ann.*
428 *Statist.*, 36(2):555–586, 04 2008.

429 Xiao-Ming Wu, Zhenguo Li, Anthony M. So, John Wright, and Shih fu Chang. Learning with partially
430 absorbing random walks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors,
431 *Advances in Neural Information Processing Systems 25*, pages 3077–3085. Curran Associates, Inc.,
432 2012.

433 Jaewon Yang and Jure Leskovec. Defining and evaluating network communities based on ground-truth.
434 *Knowledge and Information Systems*, 42(1):181–213, Jan 2015.

435 Zeyuan Allen Zhu, Silvio Lattanzi, and Vahab S Mirrokni. A local algorithm for finding well-
436 connected clusters. In *ICML (3)*, pages 396–404, 2013.