
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. In
4 contrast to existing works, we analyze PPR in a traditional statistical learning
5 setup, where we obtain samples from an unknown distribution, and aim to identify
6 connected regions of high-density (density clusters). We prove that PPR, run on
7 a neighborhood graph, extracts sufficiently salient density clusters, and provide
8 empirical support of our theory.

9 1 Introduction

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

To be more precise, let $G = (V, E, w)$ denote a weighted, undirected graph constructed from the
samples 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 (lazy) random walk transition probability matrix
and normalized¹ Laplacian matrix, respectively, which are defined as

$$\mathbf{W} = \frac{\mathbf{I} + \mathbf{D}^{-1}\mathbf{A}}{2}, \quad \mathbf{L} = \mathbf{I} - \mathbf{W},$$

18 where $\mathbf{I} \in \mathbb{R}^{n \times n}$ is the identity matrix. Classical global spectral methods take a eigendecomposition
19 $\mathbf{L} = \mathbf{U}\Sigma\mathbf{U}^T$, use some number of eigenvectors (columns in \mathbf{U}) as a feature representation for the
20 samples, and then run (say) k-means in this new feature space.

21 When applied to geometric graphs constructed from a large number of samples, global spectral
22 clustering methods can be computationally cumbersome and insensitive to the local geometry of the
23 underlying distribution [??]. This has led to recent increased interest in local spectral algorithms,
24 which leverage locally-biased spectra computed using random walks around a user-specified seed
25 node. A popular local clustering algorithm is Personalized PageRank (PPR), first introduced by [?],
26 and further developed by [????], among others.

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

Local spectral clustering techniques have been practically very successful [????], which has led many authors to develop supporting theory [????] 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 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 ?).

1.1 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 X : for a radius $r > 0$, we consider the r -neighborhood graph of X , denoted $G_{n,r} = (V, E)$, an unweighted graph with vertices $V = X$, and an edge $(x_i, x_j) \in E$ if and only if $\|x_i - x_j\| \leq r$, 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 $p = p(v, \alpha; G_{n,r})$, 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:

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

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 $\text{vol}_0 > 0$, we define a β -sweep cut of $p = (p_u)_{u \in V}$ as

$$S_\beta = \{u \in V : \frac{p_u}{\mathbf{D}_{uu}} > \frac{\beta}{\text{vol}_0}\}. \quad (2)$$

Having computed sweep cuts S_β over a range $\beta \in (\frac{1}{40}, \frac{1}{11})^2$, we then output a cluster estimate $\hat{C} = S_{\beta^*}$ to have minimum normalized cut $\Phi(S_{\beta^*}; G_{n,r})$, where for $S \cup S^c = G_{n,r}$, $\text{cut}(S; G_{n,r}) := |\{(u, v) \in E : u \in S, v \in S^c\}|$, $\text{vol}(S; G_{n,r}) := \sum_{u \in S} \mathbf{D}_{uu}$, and

$$\Phi(S; G_{n,r}) := \frac{\text{cut}(S; G_{n,r})}{\min\{\text{vol}(S; G_{n,r}), \text{vol}(S^c; G_{n,r})\}}. \quad (3)$$

For concreteness, we summarize this procedure in Algorithm 1.

Algorithm 1 PPR on a Neighborhood Graph

Input: data $X = \{x_1, \dots, x_n\}$, radius $r > 0$, teleportation parameter $\alpha \in [0, 1]$, seed $v \in X$, target stationary volume $\text{vol}_0 > 0$.

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

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

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

1.2 Summary of Results

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

²The choice of a specific range such as $(\frac{1}{40}, \frac{1}{11})$ is standard in the analysis of PPR algorithms, see, e.g., [?].

56 **Definition 1** (Misclassification error). For an estimator $\hat{C} \subseteq X$ and set $S \subseteq \mathbb{R}^d$, the misclassification
 57 error of S by \hat{C} is

$$|\hat{C} \setminus (S \cap X)| + |(S \cap X) \setminus \hat{C}|. \quad (4)$$

58 **Definition 2** (Consistent density cluster estimation). For an estimator $\hat{C} \subseteq X$ and cluster $C \in \mathbb{C}_f(\lambda)$,
 59 we say \hat{C} is a consistent estimator of C if for all $C' \in \mathbb{C}_f(\lambda)$ with $C \neq C'$ the following holds as
 60 $n \rightarrow \infty$:

$$C[X] \subseteq \hat{C} \quad \text{and} \quad \hat{C} \cap C'[X] = \emptyset, \quad (5)$$

61 with probability tending to 1.

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

- 63 1. In Section 2, we introduce a set of natural geometric conditions. We formalize a measure of
 64 difficulty based on these geometric conditions, and show that when properly initialized, the
 65 misclassification error of Algorithm 1 is upper bounded by this difficulty measure.
- 66 2. We further show that if the density cluster C is particularly well-conditioned, Algorithm 1
 67 will perform consistent density cluster estimation in the sense of (5).
- 68 3. Corollary 1 establishes that these statements hold also with respect to an approximate form
 69 of PPR, which can be efficiently computed.
- 70 4. In [Section 3](#), we detail some of the main technical machinery required to prove our main
 71 results, highlighting the part various geometric quantities play in the ultimate difficulty of
 72 the clustering problem.
- 73 5. In Section 4, we empirically demonstrate the tightness of the bounds in Theorems 3 and
 74 4, and provide examples showing how violations of the geometric conditions we require
 75 manifestly impact density cluster recovery by PPR.

76 On the topic of conditions, it is worth mentioning that, as density clusters are inherently local,
 77 focusing on the PPR algorithm actually eases our analysis and allows us to require fewer global
 78 regularity conditions relative to those needed for more classical global spectral algorithms.

79 1.3 Related Work

80 In addition to the background given above, a few related lines of work are worth highlighting. Building
 81 on earlier work of [?], [??] studied the limiting behaviour of spectral clustering algorithms. These
 82 authors show that when samples are obtained from a distribution, and we appropriately construct a
 83 geometric graph, the spectrum of the Laplacian converges to that of the Laplace-Beltrami operator on
 84 the data-manifold. However, relating the partition obtained using the Laplace-Beltrami operator to
 85 the more intuitively defined high-density clusters can be challenging in general.

86 More similar to our results are the works [???], who study the consistency of spectral algorithms in
 87 recovering the latent labels in certain parametric and nonparametric mixture models. These results
 88 focus on global rather than local algorithms, and as such impose global rather than local conditions
 89 on the nature of the density. Moreover, they do not in general ensure recovery of density clusters,
 90 which is the focus in our work.

91 2 Estimation of Well-Conditioned Density Clusters.

92 2.1 Geometric Conditions on Density Clusters

93 As mentioned previously, successful recovery of a density cluster by PPR requires the density cluster
 94 to be geometrically well-conditioned. At a minimum, we wish to avoid dumbbell-like sets C which
 95 contain an arbitrarily thin bridge, and as in ? we therefore introduce a buffer zone around C . Letting
 96 $B(x, r)$ be the closed ball of radius $r > 0$ centered at $x \in \mathbb{R}^d$, for a given cluster $C \subseteq \mathbb{R}^d$ and $\sigma > 0$,
 97 we refer to $C_\sigma := \{y \in \mathbb{R}^d : \inf_{x \in C} \|y - x\| \leq \sigma\}$ as the σ -expansion of C , and state our conditions
 98 with respect to C_σ .

99 More generally, over the neighborhood graph $G_{n,r}$ we would like the empirical cluster $C_\sigma[X]$ to
 100 be **well connected** everywhere in its interior, and **poorly connected** to the rest of X . This intuition

motivates our required conditions, stated with respect to a density cluster $\mathcal{C} \in \mathbb{C}_f(\lambda)$ 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) *Cluster separation*: For all $\mathcal{C}' \in \mathbb{C}_f(\lambda)$ with $\mathcal{C}' \neq \mathcal{C}$,

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

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

(A3) *Low noise density*: There exists $\gamma, c_0 > 0$ such that for all $x \in \mathbb{R}^d$ with $0 < \text{dist}(x, \mathcal{C}_\sigma) \leq \sigma$,

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

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

(A4) *Lipschitz embedding*: \mathcal{C}_σ is the image of a convex set under a biLipschitz, measure preserving mapping. Formally, there exists $\mathcal{K} \subseteq \mathbb{R}^d$ convex, and $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that $\det(\nabla g(x)) = 1$ for all $x \in \mathcal{C}_\sigma$, and for some $L \geq 1$,

$$\frac{1}{L} \|x - y\| \leq \|g(x) - g(y)\| \leq L \|x - y\| \text{ for all } x, y \in \mathcal{C}_\sigma$$

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 $x, x' \in \mathcal{K}$

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

(A5) *Bounded volume*: Let the neighborhood graph radius $0 < r \leq \sigma/2d$ be such that

$$\frac{\int_{\mathcal{C}_\sigma} \mathbb{P}(B(x, r)) f(x) dx}{\int_{\mathbb{R}^d} \mathbb{P}(B(x, r)) f(x) dx} \leq \frac{1}{2}. \quad (6)$$

The cluster separation (A2) and low noise density (A3) conditions guarantee **poor connectivity** between $\mathcal{C}_\sigma[X]$ and $X \setminus \mathcal{C}_\sigma[X]$, whereas (A1) and (A4) ensure high connectivity within $\mathcal{C}_\sigma[X]$. **It may not be immediately obvious how (A4) contributes to geometric conditioning. For now, we observe merely that random walks will mix slowly over sets with large diameter, and make some more detailed commentary in Section 3.** Finally, (A5) is a relatively harmless technical condition, merely excluding the case where \mathcal{C}_σ contains over half the total mass.

2.2 Well-Conditioned Density Clusters

We turn to formally defining a **condition number**, $\kappa(\mathcal{C})$, reflects the difficulty of the local spectral clustering task. The smaller $\kappa(\mathcal{C})$ is, the more success PPR will have in recovering \mathcal{C} . Let $\theta := (r, \sigma, \lambda, \lambda_\sigma, \Lambda_\sigma, \gamma, D, L)$ contain those geometric parameters detailed in 2.1.

Definition 3 (Well-conditioned density clusters). *For $\lambda > 0$ and $\mathcal{C} \in \mathbb{C}_f(\lambda)$, let \mathcal{C} satisfy (A1) - (A5) for some θ , and **additionally let \mathcal{C}_σ satisfy (6)**. Then, setting*

$$\begin{aligned} \Phi(\theta) &:= c_1 r \frac{d}{\sigma} \frac{\lambda}{\lambda_\sigma} \frac{(\lambda_\sigma - c_0 \frac{r^\gamma}{\gamma+1})}{\lambda_\sigma} \\ \Psi(\theta) &:= \left(c_2 \frac{\Lambda_\sigma^4 d^3 D^2 L^2}{\lambda_\sigma^4 r^2} \log^2 \left(\frac{1}{r} \right) + c_3 \log \left(\frac{\Lambda_\sigma}{\lambda_\sigma} \right) \right)^{-1} \end{aligned} \quad (7)$$

and

$$\kappa(\mathcal{C}) := \frac{\Phi(\theta)}{\Psi(\theta)} \quad (8)$$

we call \mathcal{C} a κ -well-conditioned density cluster.

At first glance (7) may appear mysterious, but as will be shown in [Section 3](#), these are merely upper bounds on the normalized cut and inverse mixing time of (the σ -expansion of) a given empirical density cluster $\mathcal{C}_\sigma[X]$ in $G_{n,r}$. In [?](#), building on the work of [?](#) and others, it is shown that the ratio of normalized cut to inverse mixing time is a fundamental quantity governing the performance of PPR over a general graph. $\kappa(\mathcal{C})$ upper bounds this ratio for an empirical density cluster over the neighborhood graph $G_{n,r}$, and is therefore a natural criterion to measure difficulty of the clustering task.

Well-initialized algorithms. As is typical in the local clustering literature, our algorithmic results will be stated with respect to specific choices or ranges of each of the user-specified parameters.

In particular, for a well-conditioned density cluster \mathcal{C} (with respect to some θ), we require

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

$$v \in \mathcal{C}_\sigma[X]^g, \text{vol}_0 \in [3/4, 5/4] \cdot n(n-1) \int_{\mathcal{C}_\sigma} \mathbb{P}(B(x, r)) f(x) dx \quad (9)$$

$\mathcal{C}_\sigma[X]^g \subseteq \mathcal{C}_\sigma[X]$ will be some large subset of $\mathcal{C}_\sigma[X]$, in particular $\text{vol}(\mathcal{C}_\sigma[X]^g; G_{n,r}) \geq \text{vol}(\mathcal{C}_\sigma[X]; G_{n,r})/2$.

Definition 4. If the input parameters to Algorithm 1 satisfy (9) for some well-conditioned density cluster \mathcal{C} , we say the algorithm is well-initialized.

In practice it is clearly not feasible to set hyperparameters based on the underlying (unknown) density f . Typically, one tunes PPR over a range of hyperparameters and optimizes for some criterion such as normalized cut; it is unclear how this scheme would affect the performance of PPR in the density clustering context.

Density cluster estimation by PPR. Theorem 1 of [?](#), combined with the results of [Section 3](#), immediately implies a bound on the volume of $\widehat{\mathcal{C}} \setminus \mathcal{C}_\sigma[X]$ (and likewise $\mathcal{C}_\sigma[X] \setminus \widehat{\mathcal{C}}$),

$$\text{vol}_{n,r}(\widehat{\mathcal{C}} \setminus \mathcal{C}_\sigma[X]), \text{vol}_{n,r}(\mathcal{C}_\sigma[X] \setminus \widehat{\mathcal{C}}) \lesssim \kappa(\mathcal{C}) \text{vol}_{n,r}(\mathcal{C}_\sigma[X]). \quad (10)$$

To translate (10) into meaningful bounds on misclassification error, we wish to preclude vertices $x \in X$ from having arbitrarily small degree. To do so, we make some regularity assumptions on $\mathcal{X} := \text{supp}(f)$.

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

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

Theorem 1. Fix $\lambda > 0$, let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ be a κ -well conditioned density cluster (with respect to some θ), and additionally assume f satisfies (A5). Then, with probability tending to one as $n \rightarrow \infty$,

$$\frac{|\mathcal{C}_\sigma[X] \setminus \widehat{\mathcal{C}}|}{|\mathcal{C}_\sigma[X]|} \leq c_5 \kappa(\mathcal{C}) \frac{\Lambda_\sigma}{\lambda_\sigma}, \quad \text{and} \quad \frac{|\widehat{\mathcal{C}} \setminus \mathcal{C}_\sigma[X]|}{|\mathcal{C}_\sigma[X]|} \leq c_6 \kappa(\mathcal{C}) \frac{\Lambda_\sigma}{\lambda_{\min}}. \quad (11)$$

for universal constants $c_4, c_5 > 0$.

The proof of Theorem 1, along with all other proofs in this paper, can be found in the supplementary material. We observe that the misclassification error is proportional to the difficulty of the clustering problem, as measured by the [condition number](#).

Neither (10) nor Theorem 1 imply consistent density cluster estimation in the sense of (5). This notion of consistency requires a uniform bound over p for all $u \in \mathcal{C}, u' \in \mathcal{C}'$

$$\frac{p_{u'}}{\mathbf{D}_{uu}} \leq \frac{1}{40\text{vol}_0} < \frac{1}{11\text{vol}_0} \leq \frac{p_u}{\mathbf{D}_{uu}}. \quad (12)$$

so that any sweep cut S_β for $\beta\text{vol}_0 \in [1/40, 1/11]$ (i.e. any sweep cut considered by Algorithm 1) will fulfill both conditions laid out in (5). In Theorem 2, we show that a sufficiently small upper

bound on $\kappa(\mathcal{C})$ ensures such a gap exists with probability one as $n \rightarrow \infty$, and therefore guarantees \hat{C} will be a consistent estimator. As was the case before, we wish to preclude arbitrarily low degree vertices, this time for points $x \in \mathcal{C}'[X]$.

(A6) \mathcal{C}' -bounded density : For each $\mathcal{C}' \in \mathbb{C}_f(\lambda)$, $\mathcal{C}' \neq \mathcal{C}$ and for all $x \in \mathcal{C}' + \sigma B$, $\lambda_\sigma \leq f(x)$ where σ, λ_σ are as in (A1).

Theorem 2. Fix $\lambda > 0$, let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ be a κ -well conditioned cluster (with respect to some θ), and additionally assume (A6) holds. If Algorithm 1 is well-initialized, there exists universal constant $c_7 > 0$ such that if

$$\kappa(\mathcal{C}) \leq c_7 \frac{\lambda_\sigma^2 r^d \nu_d}{\Lambda_\sigma \mathbb{P}(\mathcal{C}_\sigma)}, \quad (13)$$

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

A few remarks are in order.

Remark 1. We note that the restriction on $\kappa(\mathcal{C})$ imposed by (13) results in a misclassification rate on the order of r^d . (See Theorem 1). In plain terms, we are able to recover a density cluster \mathcal{C} in the sense of (5) only when we can guarantee a very small fraction of points are misclassified. This strong condition is the price we pay in order to obtain the uniform bound of 12.

Remark 2. While taking the radius of the neighborhood graph $r \rightarrow 0$ as $n \rightarrow \infty$ —and thereby ensuring $G_{n,r}$ is sparse—is computationally attractive, the presence of a factor of $\frac{\log^2(1/r)}{r}$ in $\kappa(\mathcal{C})$ unfortunately prevents us from making claims about the behavior of PPR in this regime. Although the restriction to a kernel function fixed in n is standard for theoretical analysis of spectral clustering ??, it is an interesting question whether PPR exhibits some degeneracy over r -neighborhood graphs as $r \rightarrow 0$, or if this is merely looseness in our upper bounds.

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

$$\text{for all } x \in X, \quad p(x) - \epsilon \deg_{n,r}(x) \leq p^{(\epsilon)}(x) \leq p(x) \quad (14)$$

Application of (14) within the proofs of Theorems 1 and 2 leads to analogous results which hold with respect to $p^{(\epsilon)}$.

Corollary 1. Fix $\lambda > 0$, and let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ be a κ -well-conditioned cluster (with respect to some θ). Choose input parameters $\alpha, r, \text{vol}_0, v$ to be well-initialized in the sense of (9), set $\epsilon = \frac{1}{20\text{vol}_0}$, and modify Algorithm 1 to compute the aPPR vector $p^{(\epsilon)}$ rather than the exact PPR vector p , with resulting output \hat{C} .

1. Assume (A5) holds. Then (11) is still a valid upper bound for the misclassification error of \hat{C} .

2. Assume (A6) holds. If

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

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

3 Analysis

Given an arbitrary graph $G = (V, E)$ and candidate cluster $S \subseteq G$, ? bound the volume of $\hat{C} \setminus S$ and $S \setminus \hat{C}$ in terms of the normalized cut and inverse mixing time of S . The key to deriving the algorithmic results of the previous section is therefore to show that the geometric conditions (A1) - (A4) translate to meaningful bounds on the normalized cut and inverse mixing time of $\mathcal{C}_\sigma[X]$ in $G_{n,r}$. Doing so constitutes the bulk of our technical effort.

3.1 Upper Bound on Normalized Cut

We start with an upper bound on the normalized cut (3) of $\mathcal{C}_\sigma[X]$. (In Theorem 3, the upper bound on the density in Assumption (A1) will not actually be needed, so we omit the parameter $\Lambda_\sigma > 0$ from the theorem statement.) For simplicity, we write $\Phi_{n,r}(\mathcal{C}_\sigma[X]) := \Phi(\mathcal{C}_\sigma[X]; G_{n,r})$.

Theorem 3. Fix $\lambda > 0$, and let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ satisfy Assumptions (A1)-(A3), and (A5) for some $r, \sigma, \lambda_\sigma, c_0, \gamma > 0$. Then for any $0 < \delta < 1$, $\epsilon > 0$, if

$$n \geq \frac{(2 + \epsilon)^2 \log(3/\delta)}{\epsilon^2} \left(\frac{25}{6\lambda_\sigma^2 \nu(\mathcal{C}_\sigma) \nu_d r^d} \right)^2, \quad (15)$$

then

$$\frac{\Phi_{n,r}(\mathcal{C}_\sigma[X])}{r} \leq c_1 \frac{d}{\sigma} \frac{\lambda}{\lambda_\sigma} \frac{(\lambda_\sigma - c_0 \frac{r^\gamma}{\gamma+1})}{\lambda_\sigma} + \epsilon, \quad (16)$$

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

Remark 3. Observe that the diameter D is absent from Theorem 3, in contrast to the difficulty function $\kappa(\mathcal{C})$, which worsens (increases) as D increases. This phenomenon reflects established wisdom regarding spectral partitioning algorithms more generally ??, albeit newly applied to the density clustering setting. It suggests that PPR may fail to recover $\mathcal{C}_\sigma[X]$ even when \mathcal{C} is sufficiently well-conditioned to ensure $\mathcal{C}_\sigma[X]$ has a small normalized cut in $G_{n,r}$, if the diameter D is large. This intuition will be supported by simulations in Section 4.

3.2 Lower Bound on 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{W}_S be the (lazy) random walk matrix over $G[S]$, and write

$$q_v^{(t)}(u) = e_v \mathbf{W}_S^t e_u$$

for the t -step transition probability of a random walk over $G[S]$ originating at v .³ Also write $\pi = (\pi(u))_{u \in S}$ for the stationary distribution of this random walk. (Given the definition of \mathbf{W}_S , it is well-known that a unique stationary distribution exists and is given by $\pi(u) = \deg(u; G[S]) / \text{vol}(S; G[S])$.)

Then, the relative pointwise mixing time of $G[S]$ is

$$\tau_\infty(G[S]) = \min \left\{ t : \frac{\pi(u) - q_v^{(t)}(u)}{\pi(u)} \leq \frac{1}{4}, \text{ for } u, v \in V \right\}. \quad (17)$$

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

Theorem 4. 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

$$\limsup_{n \rightarrow \infty} \tau_\infty(\mathcal{C}_\sigma[X]) \leq c_2 \frac{\Lambda_\sigma^4 d^3 D^2 L^2}{\lambda_\sigma^4 r^2} \log^2 \left(\frac{1}{r} \right) + c_3 \log \left(\frac{\Lambda_\sigma}{\lambda_\sigma} \right) \quad (18)$$

for $c_2, c_3 > 0$ universal constants.

So far as we are aware, Theorem 4 is a **novel bound** on the mixing time of random walks over neighborhood graphs.

Remark 4. The embedding assumption (A4) and Lipschitz parameter L play an important role in proving the upper bound of Theorem 4. There is some interdependence between L and other geometric parameters σ and D , which might lead one to hope that (A4) is non-essential. However, it is not possible to eliminate this condition without incurring an additional factor of at least $(D/\sigma)^d$ in (18), achieved, for instance, when \mathcal{C}_σ is a dumbbell-like set consisting of two balls of diameter D linked by a cylinder of radius σ . [??] develop theory regarding biLipschitz deformations of convex sets, wherein it is observed that star-shaped sets as well as half-moon shapes of the type we consider in Section 4 both satisfy (A4) for reasonably small values of L .

³Given 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.

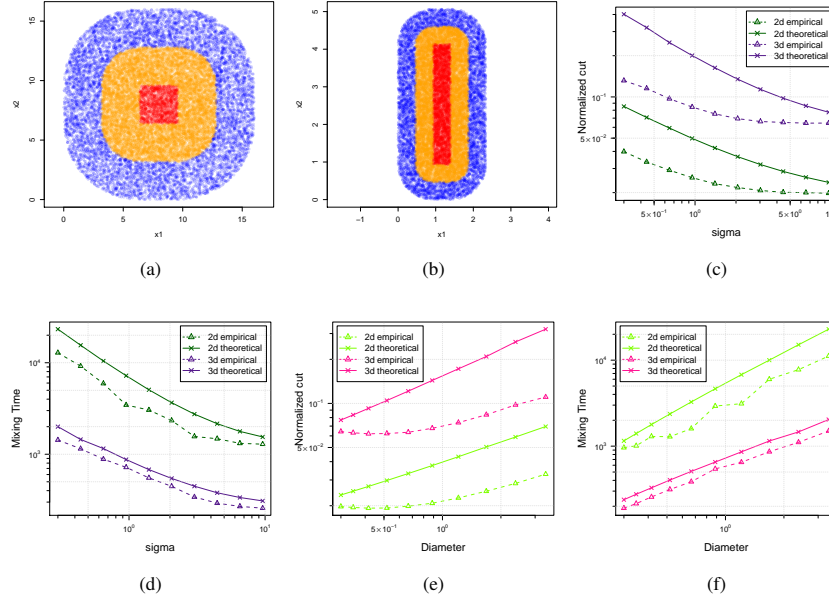


Figure 1: Samples, empirical results, and theoretical bounds for mixing time and normalized cut as diameter and thickness are varied. In (a) and (b), points in \mathcal{C} are colored in red; points in $\mathcal{C}_\sigma \setminus \mathcal{C}$ are colored in yellow; and remaining points in blue.

237 4 Experiments

238 We provide numerical experiments to investigate the tightness of our bounds on the cluster quality
 239 criteria normalized cut and mixing time, and examine the performance of PPR on the ‘two moons’
 240 dataset. For space reasons, we defer details of the experimental settings to the supplement.

241 **Validating Theoretical Bounds.** As we do not provide any theoretical lower bounds, we investigate
 242 the tightness of Theorems 3 and 4 via simulation. Figure 1 shows these theoretical bounds compared
 243 to the empirical quantities (3) and (17), as we vary the diameter D and thickness σ of the cluster
 244 \mathcal{C} .

245 Panels (d) and (f) show our theoretical bounds on mixing time tracking closely with empirical
 246 mixing time, in both 2 and 3 dimensions.⁴ This provides empirical evidence that the upper bound
 247 on mixing time given by Theorem 4 has the right dependency on both expansion parameter σ and
 248 diameter D . The story in panels (c) and (e) is less obvious. We note that while, broadly speaking,
 249 the trends do not appear to match, this gap between theory and empirical results seems largest when
 250 $\sigma \approx D$. As the ratio D/σ grows, we see the slopes of the empirical curves becoming more similar to
 251 those predicted by theory.

252 **PPR , normalized cut, and density clustering comparison.** To drive home the main implications
 253 of Theorems 1 and 2, in Figure 2 we show the behavior of PPR, normalized cut, and the density
 254 clustering algorithm of [?] on (a variant of) the famous ‘two moons’ dataset, considered a prototypical
 255 success story for spectral clustering algorithms. The first column consists of the empirical density
 256 clusters C_n and C'_n for a particular threshold λ of the density function; the second column shows the
 257 cluster recovered by PPR; the third column shows the global minimum normalized cut, computed
 258 according to the algorithm of ?; and the last column shows a cut of the density cluster tree estimator
 259 of ?.

260 Rows 1-3 show the degrading ability of PPR to recover density clusters as the two moons become
 261 less salient. Of particular interest is the fact that PPR fails to recover one of the moons even when

⁴Note 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.

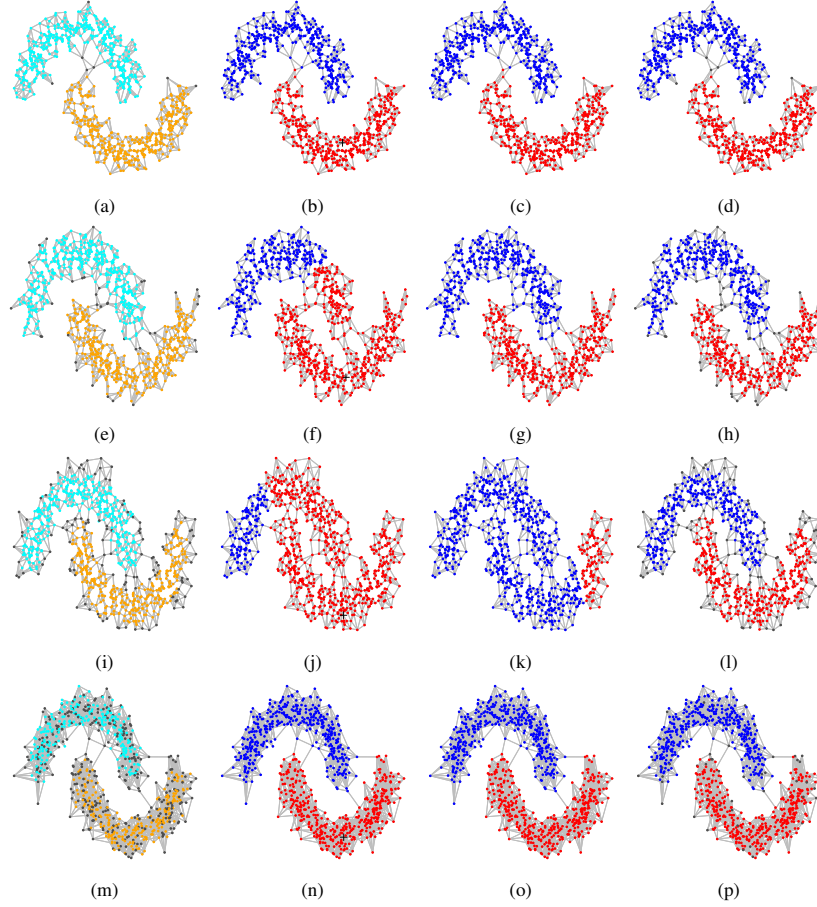


Figure 2: True density (column 1), PPR (column 2), normalized cut (column 3) and estimated density (column 4) clusters for 4 different simulated data sets. Seed node for PPR denoted by a black cross.

normalized cut still succeeds in doing so, and that a density clustering algorithm recovers a moon even when both PPR and normalized cut fail.

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.

We do not provide a theoretical lower bound showing that our geometric conditions are required for successful recovery on an upper level set. Although we investigate the matter empirically, this is a direction for future work.