
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
2 tools. Recent works have proposed and analyzed *local* spectral methods, such as
3 Personalized PageRank (PPR), which extract clusters using locally-biased random
4 walks around a user-specified seed node. In contrast to existing results, we analyze
5 PPR in a traditional statistical learning setup, where we obtain samples from
6 an unknown distribution, and aim to identify connected regions of high-density
7 (density clusters). We prove that PPR, run on a neighborhood graph, extracts
8 sufficiently salient density clusters, and provide empirical support for 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.

18 When applied to geometric graphs constructed from a large number of samples, global spectral
19 clustering methods can be computationally cumbersome and insensitive to the local geometry of the
20 underlying distribution [14, 15]. This has led to recent increased interest in local spectral algorithms,
21 which leverage locally-biased spectra computed using random walks around a user-specified seed
22 node. A popular local clustering algorithm is Personalized PageRank (PPR), first introduced by
23 Haveliwala [11], and further developed in [21, 23, 4, 15, 28], among others.

24 Local spectral clustering techniques have been practically very successful [14, 5, 8, 15, 27], which
25 has led many authors to develop supporting theory [22, 3, 7, 28] that gives worst-case guarantees on
26 traditional graph-theoretic notions of cluster quality (like conductance). In this paper, we adopt a
27 more traditional statistical viewpoint, and examine what the output of a local clustering algorithm on
28 X reveals about the unknown density f . In particular, we examine the ability of the PPR algorithm
29 to recover *density clusters* of f , which are defined as the connected components of the upper level set
30 $\{x \in \mathbb{R}^d : f(x) \geq \lambda\}$ for some threshold $\lambda > 0$ (a central object of interest in the classical statistical
31 literature on clustering, dating back to Hartigan [10]).

32 **PPR on a neighborhood graph** We now describe the clustering algorithm that will be our focus
33 for the rest of the paper. We start with the geometric graph that we form based on the samples X : for
34 a radius $r > 0$, we consider the r -neighborhood graph of X , denoted $G_{n,r} = (V, E)$, an unweighted,
35 undirected graph with vertices $V = X$, and an edge $(x_i, x_j) \in E$ if and only if $\|x_i - x_j\| \leq r$,
36 where $\|\cdot\|$ denotes Euclidean norm. We denote by $A \in \mathbb{R}^{n \times n}$ the adjacency matrix, with entries

37 $A_{uv} = 1$ if $(u, v) \in E$ and 0 otherwise, by D the diagonal degree matrix, with $D_{uu} = \sum_{v \in V} A_{uv}$,
 38 and by I the $n \times n$ identity matrix.

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

$$p = \alpha e_v + (1 - \alpha)pW, \quad (1)$$

41 where $W = (I + D^{-1}A)/2$ is the lazy random walk matrix over $G_{n,r}$ and e_v denotes the indicator
 42 vector for node v (with a 1 in the v th position and 0 elsewhere).

43 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

$$S_\beta = \left\{ u \in V : \frac{p_u}{D_{uu}} > \frac{\beta}{\text{vol}_0} \right\}. \quad (2)$$

44 We will use normalized cut to determine which sweep cut S_β is the best cluster estimate. For a set
 45 $S \subseteq V$ with complement $S^c = V \setminus S$, we define the cut as $\text{cut}(S; G_{n,r}) := \sum_{u \in S, v \in S^c} A_{uv}$, the
 46 volume as $\text{vol}(S; G_{n,r}) := \sum_{u \in S} D_{uu}$, and the *normalized cut* as

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

47 Having computed sweep cuts S_β over a range $\beta \in (\frac{1}{40}, \frac{1}{11})^1$, we output the cluster estimate $\hat{C} = S_{\beta^*}$
 48 which has minimum normalized cut $\Phi(S_{\beta^*}; G_{n,r})$. For concreteness, we summarize this procedure
 49 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

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

50 **Estimation of density clusters** Let $\mathbb{C}_f(\lambda)$ denote the connected components of the density upper
 51 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
 52 *empirical density cluster*. The symmetric set difference between estimated and empirical cluster is
 53 perhaps the most frequently used metric to quantify cluster estimation error [13, 16, 17].

54 **Definition 1** (Symmetric set difference). *For an estimator $\hat{C} \subseteq X$ and set $\mathcal{S} \subseteq \mathbb{R}^d$, the symmetric
 55 set difference of \hat{C} and $\mathcal{S} \cap X = \mathcal{S}[X]$ is*

$$\Delta(\hat{C}, \mathcal{S}) := |\hat{C} \setminus \mathcal{S}[X] \cup \mathcal{S}[X] \setminus \hat{C}|. \quad (4)$$

56 However, the symmetric set difference does not account for the distance points in $\hat{C} \setminus \mathcal{S}[X]$ may be
 57 from \mathcal{S} [20]. We therefore give a second notion of cluster estimation, first introduced by Hartigan
 58 [10] and defined asymptotically, which measures whether \hat{C} can distinguish any two distinct elements
 59 $\mathcal{C}, \mathcal{C}' \in \mathbb{C}_f(\lambda)$.

60 **Definition 2** (Consistent density cluster estimation). *For an estimator $\hat{C} \subseteq X$ and cluster $\mathcal{C} \in \mathbb{C}_f(\lambda)$,
 61 we say \hat{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
 62 $n \rightarrow \infty$:*

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

63 *with probability tending to 1.*

¹The choice of a specific range such as $(\frac{1}{40}, \frac{1}{11})$ is standard in the analysis of PPR algorithms, see, e.g., [28].

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

- 66 1. In Section 2, we introduce a set of natural geometric conditions on the density cluster \mathcal{C} ,
 67 formalize a measure of difficulty based on these geometric conditions, and show that when
 68 Algorithm 1 is properly initialized, the symmetric set difference of $\widehat{\mathcal{C}}$ and a thickened version
 69 of the density cluster \mathcal{C}_σ is upper bounded by this difficulty measure.
- 70 2. We further show that if the density cluster \mathcal{C} is particularly well-conditioned, Algorithm 1
 71 will consistently estimate a density cluster in the sense of (5).
- 72 3. In Section 3, we detail some of the analysis required to prove our main results, and expose
 73 the part various geometric quantities play in the ultimate difficulty of the clustering problem.
- 74 4. In Section 4, we empirically demonstrate the tightness of our analysis, and provide examples
 75 showing how violations of the geometric conditions we require manifestly impact density
 76 cluster recovery by PPR.

77 Our main takeaway can be summarized as follows: PPR, run on a neighborhood graph, recovers
 78 geometrically compact high-density clusters.

79 **Related work** In addition to the background given previously, a few related lines of work are worth
 80 highlighting. Similar in spirit to our results are the works [19, 18], who study the consistency of
 81 spectral algorithms in recovering the latent labels in certain nonparametric mixture models. These
 82 results focus on global rather than local algorithms, and as such impose global rather than local
 83 conditions on the nature of the density. Moreover, they do not in general guarantee recovery of
 84 density clusters, which is the focus in our work. Perhaps most importantly, these works rely on
 85 general cluster saliency conditions, which implicitly depend on many distinct geometric aspects of
 86 the cluster \mathcal{C} under consideration. We make this dependence more explicit, and in doing so expose
 87 the role each geometric condition plays in the clustering problem.

88 Additionally, we note that density clustering and level set estimation is a well-studied problem.
 89 [16, 17] study density clustering under symmetric set difference, [25, 20] exhibit minimax optimal
 90 level set estimators under Hausdorff loss and [10, 6] consider consistent estimation of the cluster tree,
 91 to note but a few works on the subject. Our goal is not to improve on these results, or offer yet another
 92 algorithm for level set estimation; indeed, seen as a density clustering algorithm, PPR has none
 93 of the optimality guarantees of the previous works. Rather, our motivation is to better understand
 94 and characterize the distinctions between those density clusters which are well conditioned for local
 95 spectral algorithms, and those which are not.

96 2 Estimation of well-conditioned density clusters

97 We formalize some geometric conditions, before using these to define a condition number $\kappa(\mathcal{C})$ which
 98 measures the difficulty PPR will have in estimating \mathcal{C} . We motivate this measure, and the underlying
 99 geometric conditions, by giving density cluster estimation guarantees for Algorithm 1 in terms of
 100 $\kappa(\mathcal{C})$.

101 **Geometric conditions on density clusters** As mentioned previously, successful recovery of a
 102 density cluster by PPR requires the density cluster to be geometrically well-conditioned. At a
 103 minimum, we wish to avoid sets \mathcal{C} which contain arbitrarily thin bridges or spikes, and therefore as in
 104 [6] we introduce a buffer zone around \mathcal{C} . Let $B(x, r)$ be the closed ball of radius $r > 0$ centered at
 105 $x \in \mathbb{R}^d$. For $\lambda > 0$, consider a given cluster $\mathcal{C} \in \mathbb{C}_f(\lambda)$. We denote the distance between x and \mathcal{C} as
 106 $\text{dist}(x, \mathcal{C}) := \inf_{y \in \mathcal{C}} \|y - x\|$, and for a given $\sigma > 0$, we refer to $\mathcal{C}_\sigma := \{x \in \mathbb{R}^d : \text{dist}(x, \mathcal{C}) \leq \sigma\}$
 107 as the σ -expansion of \mathcal{C} . We now state our conditions with respect to \mathcal{C}_σ .

108 (A1) *Bounded density within cluster:* There exist constants $\lambda_\sigma, \Lambda_\sigma$ such that $0 < \lambda_\sigma =$
 109 $\inf_{x \in \mathcal{C}_\sigma} f(x) \leq \sup_{x \in \mathcal{C}_\sigma} f(x) \leq \Lambda_\sigma < \infty$.

110 (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
 111 $\text{dist}(\mathcal{C}_\sigma, \mathcal{C}'_\sigma) := \inf_{x \in \mathcal{C}_\sigma} \text{dist}(x, \mathcal{C}'_\sigma)$.

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

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

112 (A4) *Lipschitz embedding*: There exists $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ which has the following properties: i)
 113 there exists a convex set $\mathcal{K} \subseteq \mathbb{R}^d$ with $\text{diam}(\mathcal{K}) = \sup_{x, y \in \mathcal{K}} \|x - y\| =: \rho < \infty$, such that
 114 $\mathcal{C}_\sigma = g(\mathcal{K})$, ii) $\det(\nabla g(x)) = 1$ for all $x \in \mathcal{C}_\sigma$, where $\nabla g(x)$ is the Jacobian of g evaluated
 115 at x , and iii) for some $L \geq 1$,

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

116 Simply put, \mathcal{C}_σ is the image of a convex set with finite diameter, under a measure preserving,
 117 biLipschitz transformation.

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

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

119 To motivate these conditions, for an arbitrary graph $G = (V, E)$ and subset of vertices $S \subseteq V$,
 120 consider the normalized cut $\Phi(S; G)$ (defined as in (3)), as well as the mixing time of a random walk
 121 over the induced subgraph $G[S]$ (defined in Section 3 by (16), and denoted by $\tau_\infty(G[S])$). In Zhu
 122 et al. [28], it is shown that for a constant $c > 0$, the PPR estimate \hat{C} of S satisfies

$$\text{vol}(\hat{C} \setminus S; G) + \text{vol}(S \setminus \hat{C}; G) \leq c(\Phi(S, G) \cdot \tau_\infty(G[S])) \text{vol}(S; G) \quad (6)$$

123 where the left hand side resembles a (degree-weighted) form of the symmetric set difference of
 124 (4).

125 As we will formally show in Section 3, the conditions (A1)-(A5) allow us to upper bound the
 126 normalized cut $\Phi(\mathcal{C}_\sigma[X]; G_{n,r})$, and the mixing time $\tau_\infty(G_{n,r}[\mathcal{C}_\sigma[X]])$. Intuitively, the low noise
 127 (A2) and cluster separation (A3) assumptions yield an upper bound on cut($\mathcal{C}_\sigma[X]; G_{n,r}$), the lower
 128 bound on density in (A1) yields a lower bound on $\text{vol}(\mathcal{C}_\sigma[X]; G_{n,r})$, and along with (A5), which
 129 ensures $\text{vol}(\mathcal{C}_\sigma[X]; G_{n,r}) \leq \text{vol}(\mathcal{C}_\sigma[X]^c; G_{n,r})$, these imply an upper bound on the normalized
 130 cut. (A1) and (A4) preclude bottlenecks in the induced subgraph $G_{n,r}[\mathcal{C}_\sigma[X]]$, and combined with
 131 the upper bound on diameter in (A4), they yield an upper bound on the mixing time over this
 132 subgraph.

133 **Condition number** We will return to the topic of conditions in Section 3. Now, we define the
 134 condition number, $\kappa(\mathcal{C})$, which reflects the difficulty of the local spectral clustering task. Motivated
 135 by (6), we will set $\kappa(\mathcal{C})$ to be an upper bound on $\Phi(\mathcal{C}_\sigma[X]; G_{n,r}) \cdot \tau_\infty(G_{n,r}[\mathcal{C}_\sigma[X]])$. The smaller
 136 $\kappa(\mathcal{C})$ is, the more success PPR will have in recovering \mathcal{C} . Let $\theta := (r, \sigma, \lambda, \lambda_\sigma, \Lambda_\sigma, \gamma, \rho, L)$ contain
 137 those geometric parameters detailed in (A1) - (A5).

138 **Definition 3** (Well-conditioned density clusters). *For $\lambda > 0$ and $\mathcal{C} \in \mathbb{C}_f(\lambda)$, let \mathcal{C} satisfy (A1) - (A5)
 139 for some θ . Then, for universal constants $c_1, c_2, c_3 > 0$ to be specified later, we set*

$$\Phi_u(\theta) := c_1 r \frac{d}{\sigma} \frac{\lambda}{\lambda_\sigma} \frac{(\lambda_\sigma - c_0 \frac{r^\gamma}{\gamma+1})}{\lambda_\sigma}, \tau_u(\theta) := c_2 \frac{\Lambda_\sigma^4 d^3 \rho^2 L^2}{\lambda_\sigma^4 r^2} \log^2 \left(\frac{1}{r} \right) + c_3 \quad (7)$$

140 and letting $\kappa(\mathcal{C}) := \Phi_u(\theta) \cdot \tau_u(\theta)$, we call \mathcal{C} a κ -well-conditioned density cluster.

141 We note that $\Phi_u(\theta)$ and $\tau_u(\theta)$ are exactly the upper bounds on $\Phi(\mathcal{C}_\sigma[X]; G_{n,r})$ and $\tau_\infty(G_{n,r}[\mathcal{C}_\sigma[X]])$
 142 that we derive in Section 3.

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

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

$$0 < r \leq \frac{\sigma}{2d}, \alpha \in [1/10, 1/9] \cdot \frac{1}{\tau_u(\theta)},$$

$$v \in \mathcal{C}_\sigma[X]^g, \text{vol}_0 \in [3/4, 5/4] \cdot n(n-1) \int_{\mathcal{C}_\sigma} \mathbb{P}(B(x, r)) f(x) dx \quad (8)$$

where $\mathcal{C}_\sigma[X]^g \subseteq \mathcal{C}_\sigma[X]$ will be some large subset of $\mathcal{C}_\sigma[X]$. In particular, letting $\text{vol}_{n,r}(S) := \text{vol}(S; G_{n,r})$ for $S \subseteq X$, we have $\text{vol}_{n,r}(\mathcal{C}_\sigma[X]^g) \geq \text{vol}_{n,r}(\mathcal{C}_\sigma[X])/2$.

Definition 4. If the input parameters to Algorithm 1 satisfy (8) 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 minimum normalized cut; it is not obvious how this scheme would affect the performance of PPR in the density clustering context.

Density cluster estimation by PPR The results of Section 3, along with (6), give an upper bound on the volume of $\widehat{\mathcal{C}} \setminus \mathcal{C}_\sigma[X]$ and $\mathcal{C}_\sigma[X] \setminus \widehat{\mathcal{C}}$,

$$\text{vol}_{n,r}(\widehat{\mathcal{C}} \setminus \mathcal{C}_\sigma[X]) + \text{vol}_{n,r}(\mathcal{C}_\sigma[X] \setminus \widehat{\mathcal{C}}) \leq c\kappa(\mathcal{C})\text{vol}_{n,r}(\mathcal{C}_\sigma[X]). \quad (9)$$

To translate (9) into meaningful bounds on the symmetric set difference $\Delta(\mathcal{C}_\sigma[X], \widehat{\mathcal{C}})$, we wish to preclude vertices $x \in X$ from having arbitrarily small degree, and so we make some regularity assumptions on $\mathcal{X} := \text{supp}(f)$. Let ν denote the Lebesgue measure on \mathbb{R}^d , and $\nu_d := \nu(B)$ be the measure of the unit ball $B = B(0, 1)$.

(A6) *Regular support:* 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 $x \in \partial\mathcal{X}$, $\nu(B(x, r) \cap \mathcal{X}) \geq c\nu_d r^d$.

Note that the latter condition in (A6) will be satisfied if, for instance, the support \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 (A6). Then, there exists a universal constant $c_4 > 0$ such that with probability tending to 1 as $n \rightarrow \infty$,

$$\Delta(\mathcal{C}_\sigma[X], \widehat{\mathcal{C}}) \leq c_4\kappa(\mathcal{C}) \frac{\Lambda_\sigma}{\lambda_{\min}}. \quad (10)$$

The proof of Theorem 1, along with all other proofs in this paper, can be found in the supplementary material. We observe that the symmetric set difference $\Delta(\mathcal{C}_\sigma[X], \widehat{\mathcal{C}})$ is proportional to the difficulty of the clustering problem, as measured by the condition number $\kappa(\mathcal{C})$.

Neither (9) nor Theorem 1 imply consistent density cluster estimation in the sense of (5). This notion of consistency requires a uniform bound over p : namely, for all $\mathcal{C}' \in \mathbb{C}_f(\lambda)$, $\mathcal{C}' \neq \mathcal{C}$, and each $u \in \mathcal{C}$, $w \in \mathcal{C}'$,

$$\frac{p_w}{D_{ww}} \leq \frac{1}{40\text{vol}_0} < \frac{1}{11\text{vol}_0} \leq \frac{p_u}{D_{uu}}, \quad (11)$$

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 $\widehat{\mathcal{C}}$ will be a consistent estimator. As was the case before, we wish to preclude arbitrarily low degree vertices, this time for points $x \in \mathcal{C}'[X]$.

(A7) *Bounded density:* Letting σ, λ_σ be as in (A1), for each $\mathcal{C}' \in \mathbb{C}_f(\lambda)$ and for all $x \in \mathcal{C}'_\sigma$, $\lambda_\sigma \leq f(x)$.

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 f satisfies (A7). If Algorithm 1 is well-initialized, there exists a universal constant $c_5 > 0$ such that if

$$\kappa(\mathcal{C}) \leq c_5 \frac{\lambda_\sigma^2 r^d \nu_d}{\Lambda_\sigma \mathbb{P}(\mathcal{C}_\sigma)}, \quad (12)$$

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

Remark 1. We note that the restriction on $\kappa(\mathcal{C})$ imposed by (12) results in a symmetric set difference $\Delta(\mathcal{C}_\sigma[X], \widehat{\mathcal{C}})$ on the order of r^d . 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 will be misclassified. This strong condition is the price we pay in order to obtain the uniform bound of (11).

188 *Remark 2.* While taking the radius of the neighborhood graph $r \rightarrow 0$ as $n \rightarrow \infty$ —and thereby
 189 ensuring $G_{n,r}$ is sparse—is computationally attractive, the presence of a factor of $\frac{\log^2(1/r)}{r}$ in $\kappa(\mathcal{C})$
 190 unfortunately prevents us from making claims about the behavior of PPR in this regime. Although
 191 the restriction to a kernel function fixed in n is standard for theoretical analysis of spectral clustering
 192 [18, 26], it is an interesting question whether PPR exhibits some degeneracy over r -neighborhood
 193 graphs as $r \rightarrow 0$, or if this is merely looseness in our upper bounds.

194 **Approximate PPR vector** In practice, exactly solving (1) may be too computationally expensive.
 195 To address this limitation, Andersen et al. [4] introduced the ϵ -approximate PPR vector (aPPR),
 196 which we will denote $p^{(\epsilon)}$. We refer the curious reader to [4] for a formal algorithmic definition of
 197 the aPPR vector, and limit ourselves to highlighting a few salient points. Namely, the aPPR vector
 198 can be computed in order $\mathcal{O}(\frac{1}{\epsilon\alpha})$ time, while satisfying the following uniform error bound:

$$\text{for all } u \in V, \quad p(u) - \epsilon D_{uu} \leq p^{(\epsilon)}(u) \leq p(u). \quad (13)$$

199 Application of (13) within the proofs of Theorems 1 and 2 leads to analogous results which hold with
 200 respect to $p^{(\epsilon)}$. We formally state and prove this fact in the supplementary material.

201 3 Analysis

202 The primary technical contribution of our work is showing that the geometric conditions (A1) - (A5)
 203 translate to meaningful bounds on the normalized cut and mixing time of $\mathcal{C}_\sigma[X]$ in $G_{n,r}$. In doing
 204 so, we elaborate on how some of the geometric conditions introduced in Section 2 contribute to the
 205 difficulty of the clustering problem.

206 **Normalized cut** We start with a finite sample upper bound on the normalized cut (3) of $\mathcal{C}_\sigma[X]$. For
 207 simplicity, we write $\Phi_{n,r}(\mathcal{C}_\sigma[X]) := \Phi(\mathcal{C}_\sigma[X]; G_{n,r})$.

208 **Theorem 3.** Fix $\lambda > 0$, and let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ satisfy Assumptions (A1)-(A3), and (A5) for some
 209 $r, \sigma, \lambda_\sigma, c_0, \gamma > 0$ (no bound on maximum density is needed). Then for any $0 < \delta < 1$, $\epsilon > 0$, if

$$n \geq \frac{(2 + \epsilon)^2 \log(3/\delta)}{\epsilon^2} \left(\frac{25}{6\lambda_\sigma^2 \nu(\mathcal{C}_\sigma) \nu_d r^d} \right)^2, \quad (14)$$

210 then

$$\frac{\Phi_{n,r}(\mathcal{C}_\sigma[X])}{r} \leq c_1 \frac{d}{\sigma} \frac{\lambda}{\lambda_\sigma} \frac{(\lambda_\sigma - c_0 \frac{r^\gamma}{\gamma+1})}{\lambda_\sigma} + \epsilon, \quad (15)$$

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

212 *Remark 3.* Observe that the diameter ρ is absent from Theorem 3, in contrast to the difficulty function
 213 $\kappa(\mathcal{C})$, which worsens (increases) as ρ increases. This phenomenon reflects established wisdom
 214 regarding spectral partitioning algorithms more generally [9, 12], albeit newly applied to the density
 215 clustering setting. It suggests that if the diameter ρ is large, PPR may fail to recover $\mathcal{C}_\sigma[X]$ even
 216 when \mathcal{C} is sufficiently well-conditioned to ensure $\mathcal{C}_\sigma[X]$ has a small normalized cut in $G_{n,r}$. This
 217 intuition will be supported by simulations in Section 4.

Inverse mixing time For $S \subseteq V$, denote by $G[S] = (S, E_S)$ the subgraph induced by S (where the
 edges are $E_S = E \cap (S \times S)$). Let W_S be the (lazy) random walk matrix over $G[S]$, and write

$$q_v^{(t)}(u) = e_v W_S^t e_u$$

218 for the t -step transition probability of the lazy random walk over $G[S]$ originating at $v \in V$. Also
 219 write $\pi = (\pi(u))_{u \in S}$ for the stationary distribution of this random walk. (As W_S is the transition
 220 matrix of a lazy random walk, it is well-known that a unique stationary distribution exists and is given
 221 by $\pi(u) = (D_S)_{uu} / \text{vol}(S; G[S])$, where D_S is the degree matrix of $G[S]$.)

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

$$\tau_\infty(G[S]) = \min \left\{ t : \frac{\pi(u) - q_v^{(t)}(u)}{\pi(u)} \leq \frac{1}{4}, \text{ for } u, v \in V \right\}. \quad (16)$$

223 In the following theorem, we give an asymptotic (in the number of vertices n) upper bound on
 224 $\tau_\infty(G_{n,r}[\mathcal{C}_\sigma[X]])$.

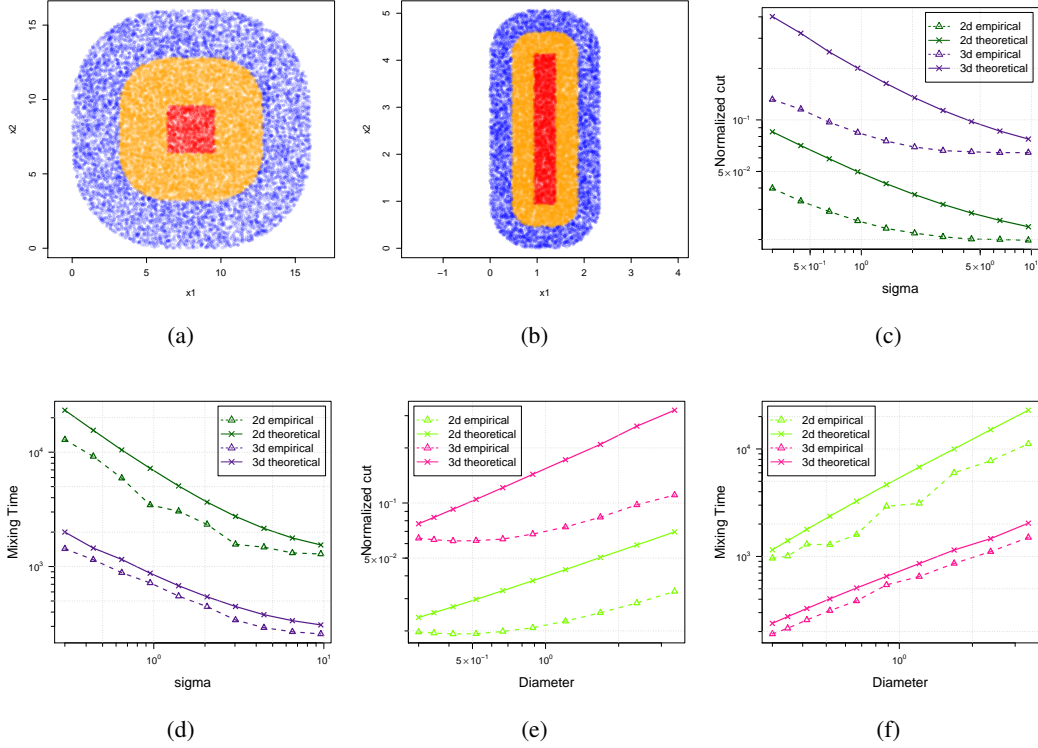


Figure 1: Samples, empirical results, and theoretical bounds for mixing time and normalized cut as diameter and thickness are varied. In (a) and (b), points in \mathcal{C} are colored in red; points in $\mathcal{C}_\sigma \setminus \mathcal{C}$ are colored in yellow; and remaining points in blue.

Theorem 4. Fix $\lambda > 0$, and let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ satisfy Assumptions (A1) and (A4) for some $\sigma, \lambda_\sigma, \Lambda_\sigma, \rho, L > 0$. Then, for any $0 < r < \sigma/2\sqrt{d}$, with probability 1

$$\limsup_{n \rightarrow \infty} \tau_\infty(G_{n,r}[\mathcal{C}_\sigma[X]]) \leq c_2 \frac{\Lambda_\sigma^4 d^3 \rho^2 L^2}{\lambda_\sigma^4 r^2} \log^2 \left(\frac{1}{r} \right) + c_3 \quad (17)$$

for $c_2, c_3 > 0$ universal constants.

To the best of our knowledge, Theorem 4 is the first bound, albeit asymptotic, on the mixing time of random walks over neighborhood graphs which is independent of n , the number of vertices.

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 ρ , 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 $(\rho/\sigma)^d$ in (17), achieved, for instance, when \mathcal{C}_σ is a dumbbell-like set consisting of two balls of diameter ρ linked by a cylinder of radius σ . [2, 1] 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 .

4 Experiments

We provide numerical experiments to investigate the tightness of our bounds on normalized cut and mixing time of $\mathcal{C}_\sigma[X]$, and examine the performance of PPR on the “two moons” dataset. For space reasons, we defer details of the experimental settings to the supplement.

Validating theoretical bounds As we do not provide any theoretical lower bounds, we investigate the tightness of Theorems 3 and 4 via simulation. Figure 1 compares these theoretical bounds with

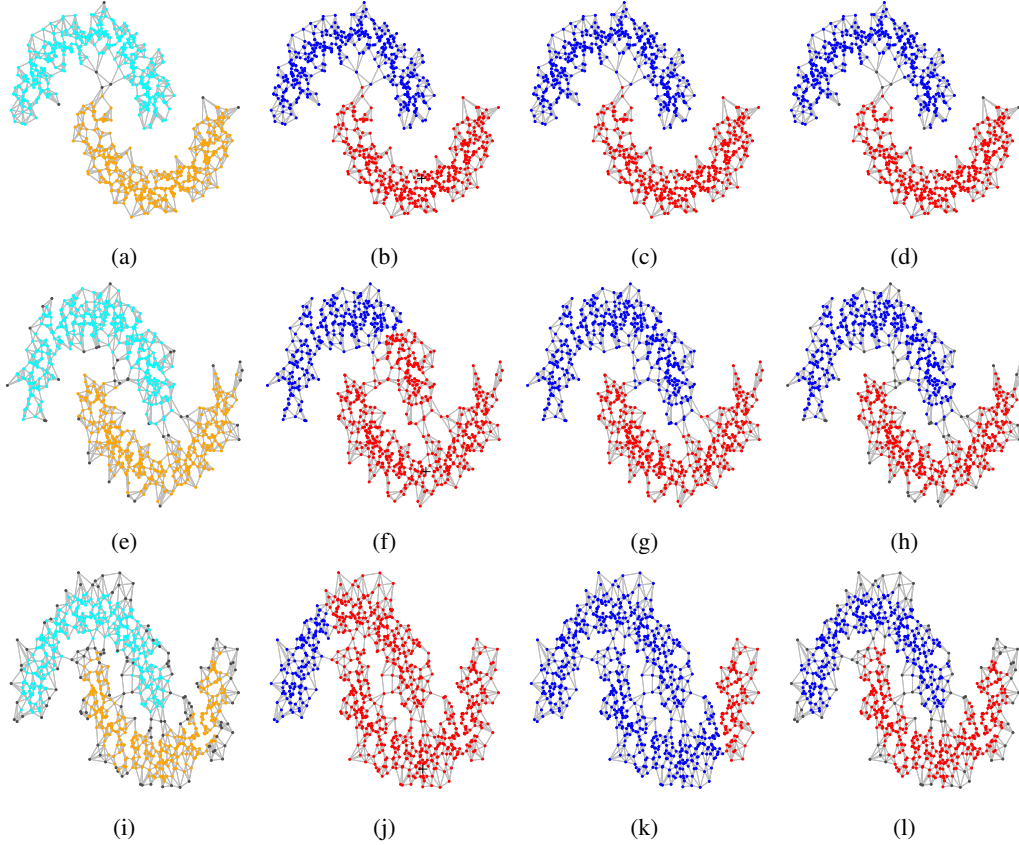


Figure 2: True density (column 1), PPR (column 2), normalized cut (column 3) and estimated density (column 4) clusters for 3 different simulated data sets. Seed node for PPR denoted by a black cross.

the empirical quantities (3) and (16), as we vary the diameter ρ and thickness σ of a cluster \mathcal{C} . Panels (a) and (b) show the resulting empirical clusters for two different values of ρ and σ .

Panels (d) and (f) show our theoretical bounds on mixing time tracking closely with empirical mixing time, in both 2 and 3 dimensions.² This provides empirical evidence that the upper bound on mixing time given by Theorem 4 has the right dependency on both expansion parameter σ and diameter ρ . 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 σ and ρ are approximately equal. As the ratio ρ/σ grows, we see the slopes of the empirical curves become more similar to those predicted by theory.

Empirical behavior of PPR To drive home the main implications of Theorems 1 and 2, in Figure 2 we show the behavior of PPR, normalized cut, and the density clustering algorithm of [6] on the well known “two moons” dataset (with added 2d Gaussian noise), considered a prototypical success story for spectral clustering algorithms. 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 cluster recovered by PPR; the third column shows the global minimum normalized cut, computed according to the algorithm of [24]; and the last column shows a cut of the density cluster tree estimator of [6].

Figure 2 shows the degrading ability of PPR to recover density clusters as the two moons become less well-separated. Of particular interest is the fact that PPR fails to recover one of the moons even when normalized cut still succeeds in doing so, supporting our claim from Remark 3. Additionally,

²Note 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.

264 we note that a density clustering algorithm recovers a moon even when both PPR and normalized
265 cut fail, lending empirical weight to our overall message that PPR recovers only geometrically
266 well-conditioned density clusters.

267 **5 Discussion**

268 For given data, there are an almost limitless number of ways to define what the “right” clustering is.
269 We have considered one such notion – density upper level sets – and have detailed a set of natural
270 geometric criteria which, when appropriately satisfied, translate to provable bounds on estimation of
271 the cluster by PPR. We do not, however, provide a theoretical lower bound showing that our geometric
272 conditions are required for successful recovery on an upper level set. Although we investigate the
273 matter empirically, this is a direction for future work.

References

- [1] Yasin Abbasi-Yadkori. Fast mixing random walks and regularity of incompressible vector fields. *arXiv preprint arXiv:1611.09252*, 2016.
- [2] Yasin Abbasi-Yadkori, Peter Bartlett, Victor Gabillon, and Alan Malek. Hit-and-Run for Sampling and Planning in Non-Convex Spaces. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, pages 888–895, 2017.
- [3] 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.
- [4] 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.
- [5] 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.
- [6] 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.
- [7] 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.
- [8] 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.
- [9] Stephen Guattery and Gary L Miller. On the performance of spectral graph partitioning methods. In *SODA*, volume 95, pages 233–242, 1995.
- [10] John A. Hartigan. Consistency of single-linkage for high-density clusters. *Journal of the American Statistical Association*, 1981.
- [11] 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.
- [12] 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.
- [13] Aleksandr Petrovich Korostelev and Alexandre B Tsybakov. *Minimax theory of image reconstruction*, volume 82. 2012.
- [14] Jure Leskovec, Kevin J. Lang, and Michael Mahoney. Empirical comparison of algorithms for network community detection. In *Proceedings of the 19th International Conference on World Wide Web*, 2010.
- [15] Michael W. Mahoney, Lorenzo Orecchia, and Nisheeth K. Vishnoi. 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.
- [16] Wolfgang Polonik. Measuring mass concentrations and estimating density contour clusters-an excess mass approach. *The Annals of Statistics*, 23(3):855–881, 1995.
- [17] Philippe Rigollet and Régis Vert. Optimal rates for plug-in estimators of density level sets. *Bernoulli*, 15(4):1154–1178, 2009.
- [18] Geoffrey Schiebinger, Martin J. Wainwright, and Bin Yu. The geometry of kernelized spectral clustering. *Ann. Statist.*, 43(2):819–846, 04 2015.

- 321 [19] Tao Shi, Mikhail Belkin, and Bin Yu. Data spectroscopy: Eigenspaces of convolution operators
322 and clustering. *Ann. Statist.*, 37(6B):3960–3984, 12 2009.
- 323 [20] Aarti Singh, Clayton Scott, and Robert Nowak. Adaptive hausdorff estimation of density level
324 sets. *Ann. Statist.*, 37(5B):2760–2782, 10 2009.
- 325 [21] Daniel A Spielman and Shang-Hua Teng. Spectral sparsification of graphs. *SIAM Journal on*
326 *Computing*, 40(4):981–1025, 2011.
- 327 [22] Daniel A Spielman and Shang-Hua Teng. A local clustering algorithm for massive graphs and
328 its application to nearly linear time graph partitioning. *SIAM Journal on Computing*, 42(1):
329 1–26, 2013.
- 330 [23] Daniel A Spielman and Shang-Hua Teng. Nearly linear time algorithms for preconditioning
331 and solving symmetric, diagonally dominant linear systems. *SIAM Journal on Matrix Analysis*
332 *and Applications*, 35(3):835–885, 2014.
- 333 [24] Arthur Szlam and Xavier Bresson. Total variation, cheeger cuts. In *ICML*, pages 1039–1046,
334 2010.
- 335 [25] Alexandre B Tsybakov. On nonparametric estimation of density level sets. *The Annals of*
336 *Statistics*, 25(3):948–969, 1997.
- 337 [26] Ulrike von Luxburg, Mikhail Belkin, and Olivier Bousquet. Consistency of spectral clustering.
338 *Ann. Statist.*, 36(2):555–586, 04 2008.
- 339 [27] Xiao-Ming Wu, Zhenguo Li, Anthony M. So, John Wright, and Shih fu Chang. Learning
340 with partially absorbing random walks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q.
341 Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 3077–3085.
342 Curran Associates, Inc., 2012.
- 343 [28] Zeyuan Allen Zhu, Silvio Lattanzi, and Vahab S Mirrokni. A local algorithm for finding
344 well-connected clusters. In *ICML (3)*, pages 396–404, 2013.