
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 results, 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 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 [16, 17]. 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 [23, 25, 4, 17, 30], among others.

24 Local spectral clustering techniques have been practically very successful [16, 5, 8, 17, 29], which
25 has led many authors to develop supporting theory [24, 3, 7, 30] 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 central interest in the classical
31 statistical 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 $\|x_i - x_j\| \leq r$, where

36 $\|\cdot\|$ denotes Euclidean norm. We denote by $A \in \mathbb{R}^{n \times n}$ the adjacency matrix, with entries $A_{uv} = 1$ if
 37 and only if $(u, v) \in E$, and by D the diagonal degree matrix, with $D_{uu} = \sum_{v \in V} A_{uv}$.

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

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

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

42 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

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

43 We need a metric to determine which sweep cut S_β is the best cluster estimate. For a set $S \subseteq V$ with
 44 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 volume to be
 45 $\text{vol}(S; G_{n,r}) := \sum_{u \in S} D_{uu}$, and the *normalized cut* as

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

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

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

49 **Estimation of density clusters** Let $\mathbb{C}_f(\lambda)$ denote the connected components of the density upper
 50 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
 51 *empirical density cluster*. The symmetric set difference between estimated and empirical cluster is
 52 perhaps the most frequently used metric to quantify cluster estimation error [15, 18, 19].

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

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

55 However, the symmetric set difference does not account for the distance points in $\hat{C} \setminus \mathcal{S}[X]$ may be
 56 from \mathcal{S} [22]. We therefore give a second notion of cluster estimation, first introduced by Hartigan
 57 [10] and defined asymptotically, which measures whether \hat{C} can distinguish any two distinct elements
 58 $\mathcal{C}, \mathcal{C}' \in \mathbb{C}_f(\lambda)$.

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

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

62 *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., [30].

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

- 65 1. In Section 2, we introduce a set of natural geometric conditions, formalize a measure of
 66 difficulty based on these geometric conditions, and show that when properly initialized, the
 67 symmetric set difference of Algorithm 1 is upper bounded by this difficulty measure.
- 68 2. We further show that if the density cluster \mathcal{C} is particularly well-conditioned, Algorithm 1
 69 consistently estimate a density cluster in the sense of (5).
- 70 3. In Section 3, we detail some of the main technical machinery required to prove our main
 71 results, and expose the part various geometric quantities play in the ultimate difficulty of the
 72 clustering problem.
- 73 4. 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 Our main takeaway can be summarized as follows: PPR, run on a neighborhood graph, recovers
 77 geometrically compact high-density clusters.

78 **Related Work.** In addition to the background given previously, a few related lines of work are
 79 worth highlighting. Similar in spirit to our results are the works [21, 20], who study the consistency
 80 of spectral algorithms in recovering the latent labels in certain parametric and nonparametric mixture
 81 models. These results focus on global rather than local algorithms, and as such impose global rather
 82 than local conditions on the nature of the density. Moreover, they do not in general guarantee recovery
 83 of density clusters, which is the focus in our work. Perhaps most importantly, these works rely on
 84 general cluster saliency conditions, which depend implicitly on many distinct geometric aspects of
 85 the cluster \mathcal{C} under consideration. We make this dependence explicit, and in so doing, expose the role
 86 each geometric condition plays in the clustering problem.

87 Additionally, we note that density clustering and level set estimation is a well-studied problem.
 88 [18, 19] study density clustering under symmetric set difference, [27, 22] prove minimax optimal
 89 level set estimators under Hausdorff loss and [10, 6] consider consistent estimation of the cluster
 90 tree, to note but a few works on the subject. Our goal is not to improve on these results, or offer yet
 91 another algorithm for level set estimation; indeed, seen as a density clustering algorithm, PPR has
 92 none of the optimality guarantees of the previous works. This is in fact a major point of our article:
 93 PPR can provably recover density clusters, but only under strong geometric conditions.

94 2 Estimation of well-conditioned density clusters

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

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

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

109 (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
 110 $\text{dist}(\mathcal{C}_\sigma, \mathcal{C}'_\sigma) := \inf_{x \in \mathcal{C}_\sigma} \text{dist}(x, \mathcal{C}'_\sigma)$.

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

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

111 (A4) *Lipschitz embedding*: There exists $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ which has the following properties: i)
 112 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
 113 $\mathcal{C}_\sigma = g(\mathcal{K})$, ii) $\det(\nabla g(x)) = 1$ for all $x \in \mathcal{C}_\sigma$, where $\nabla g(x)$ is the Jacobian of g evaluated
 114 at x , and iii) for some $L \geq 1$,

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

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

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

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

118 Thinking of $\mathcal{C}_\sigma[X]$ as a subset of vertices in $G_{n,r}$, we would like $\mathcal{C}_\sigma[X]$ to be internally well-
 119 connected, while being poorly connected to the rest of X . The cluster separation (A2) and low
 120 noise density (A3) conditions guarantee low connectivity between $\mathcal{C}_\sigma[X]$ and $X \setminus \mathcal{C}_\sigma[X]$ in $G_{n,r}$,
 121 whereas (A1) and (A4) ensure high connectivity within $\mathcal{C}_\sigma[X]$. It may not be immediately obvious
 122 how (A4) contributes to geometric conditioning. For now, we observe merely that random walks
 123 will mix slowly over sets with large diameter, and comment on this condition in more detail in
 124 Section 3. Finally, (A5) is a relatively harmless technical condition, merely excluding the case where
 125 $\text{vol}(\mathcal{C}_\sigma[X]; G_{n,r}) > \text{vol}(X; G_{n,r})/2$.

126 We can now formally define the condition number, $\kappa(\mathcal{C})$, which reflects the difficulty of the local
 127 spectral clustering task. The smaller $\kappa(\mathcal{C})$ is, the more success PPR will have in recovering \mathcal{C} . Let
 128 $\theta := (r, \sigma, \lambda, \lambda_\sigma, \Lambda_\sigma, \gamma, \rho, L)$ contain those geometric parameters detailed in (A1) - (A5).

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

$$\Phi_u(\theta) := c_1 r \frac{d}{\sigma} \frac{\lambda}{\lambda_\sigma} \frac{(\lambda_\sigma - c_0 \frac{r^\gamma}{\gamma+1})}{\lambda_\sigma}, \quad \Psi_u(\theta) := \left(c_2 \frac{\Lambda_\sigma^4 d^3 D^2 L^2}{\lambda_\sigma^4 r^2} \log^2 \left(\frac{1}{r} \right) + c_3 \right)^{-1} \quad (6)$$

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

132 At first glance (6) may appear mysterious, but as will be shown in Section 3, $\Phi_u(\theta)$ and $\Psi_u(\theta)$ are
 133 merely upper bounds on the normalized cut (3) of $\mathcal{C}_\sigma[X]$ in $G_{n,r}$, and the inverse of the mixing time
 134 (defined in Section 3 by (15)) of $\mathcal{C}_\sigma[X]$ in $G_{n,r}$. In [30], building on the work of [4] and others,
 135 it is shown that the ratio of normalized cut to inverse mixing time (or equivalently, the product of
 136 normalized cut and mixing time) is a fundamental quantity governing the clustering performance of
 137 PPR on a general graph. The condition number $\kappa(\mathcal{C})$ is an asymptotic upper bound of this ratio for
 138 an empirical density cluster over the neighborhood graph $G_{n,r}$, and is therefore a natural criterion to
 139 measure difficulty of the density clustering task.

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

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

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

$$v \in \mathcal{C}_\sigma[X]^g, \text{vol}_0 \in [3/4, 5/4] \cdot n(n-1) \int_{\mathcal{C}_\sigma} \mathbb{P}(B(x, r)) f(x) dx \quad (7)$$

143 where $\mathcal{C}_\sigma[X]^g \subseteq \mathcal{C}_\sigma[X]$ will be some large subset of $\mathcal{C}_\sigma[X]$. In particular, letting $\text{vol}_{n,r}(S) :=$
 144 $\text{vol}(S; G_{n,r})$ for $S \subseteq X$, we have $\text{vol}_{n,r}(\mathcal{C}_\sigma[X]^g) \geq \text{vol}_{n,r}(\mathcal{C}_\sigma[X])/2$.

145 **Definition 4.** If the input parameters to Algorithm 1 satisfy (7) for some well-conditioned density
146 cluster \mathcal{C} , we say the algorithm is well-initialized.

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

151 **Density cluster estimation by PPR** Theorem 1 of [30], combined with the results of Section 3,
152 immediately implies a bound on the volume of $\widehat{C} \setminus \mathcal{C}_\sigma[X]$ (and likewise $\mathcal{C}_\sigma[X] \setminus \widehat{C}$),²

$$\text{vol}_{n,r}(\widehat{C} \setminus \mathcal{C}_\sigma[X]), \text{vol}_{n,r}(\mathcal{C}_\sigma[X] \setminus \widehat{C}) \lesssim \kappa(\mathcal{C}) \text{vol}_{n,r}(\mathcal{C}_\sigma[X]). \quad (8)$$

153 To translate (8) into meaningful bounds on the symmetric set difference $\Delta(\mathcal{C}_\sigma[X], \widehat{C})$, we wish to
154 preclude vertices $x \in X$ from having arbitrarily small degree. To do so, we make some regularity
155 assumptions on $\mathcal{X} := \text{supp}(f)$. Let ν denote the Lebesgue measure on \mathbb{R}^d , and $\nu_d = \nu(B)$ be the
156 measure of the unit ball $B = B(0, 1)$.

157 (A6) *Regular support:* There exists some number $\lambda_{\min} > 0$ such that $\lambda_{\min} < f(x)$ for all $x \in \mathcal{X}$.
158 Additionally, there exists some $c > 0$ such that for each $x \in \partial\mathcal{X}$, $\nu(B(x, r) \cap \mathcal{X}) \geq c\nu_d r^d$.

159 Note that the latter condition in (A6) will be satisfied if, for instance, the support \mathcal{X} is a σ -expanded
160 set.

161 **Theorem 1.** Fix $\lambda > 0$, let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ be a κ -well conditioned density cluster (with respect to some
162 θ), and additionally assume f satisfies (A6). Then, there exists a universal constant $c_4 > 0$ such that
163 with probability tending to one as $n \rightarrow \infty$,

$$\Delta(\mathcal{C}_\sigma[X], \widehat{C}) \leq c_4 \kappa(\mathcal{C}) \frac{\Lambda_\sigma}{\lambda_{\min}}. \quad (9)$$

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

167 Neither (8) nor Theorem 1 imply consistent density cluster estimation in the sense of (5). This notion
168 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
169 $u \in \mathcal{C}, w \in \mathcal{C}'$,

$$\frac{p_w}{D_{ww}} \leq \frac{1}{40\text{vol}_0} < \frac{1}{11\text{vol}_0} \leq \frac{p_u}{D_{uu}}, \quad (10)$$

170 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)
171 will fulfill both conditions laid out in (5). In Theorem 2, we show that a sufficiently small upper
172 bound on $\kappa(\mathcal{C})$ ensures such a gap exists with probability one as $n \rightarrow \infty$, and therefore guarantees \widehat{C}
173 will be a consistent estimator. As was the case before, we wish to preclude arbitrarily low degree
174 vertices, this time for points $x \in \mathcal{C}'[X]$.

175 (A7) *Bounded density:* Letting σ, λ_σ be as in (A1), for each $\mathcal{C}' \in \mathbb{C}_f(\lambda)$ and for all $x \in \mathcal{C}'_\sigma$,
176 $\lambda_\sigma \leq f(x)$.

177 **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
178 additionally assume (A7) holds. If Algorithm 1 is well-initialized, there exists a universal constant
179 $c_5 > 0$ such that if

$$\kappa(\mathcal{C}) \leq c_5 \frac{\lambda_\sigma^2 r^d \nu_d}{\Lambda_\sigma \mathbb{P}(\mathcal{C}_\sigma)}, \quad (11)$$

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

181 *Remark 1.* We note that the restriction on $\kappa(\mathcal{C})$ imposed by (11) results in a misclassification rate on
182 the order of r^d . (See Theorem 1). In plain terms, we are able to recover a density cluster \mathcal{C} in the
183 sense of (5) only when we can guarantee a very small fraction of points will be misclassified. This
184 strong condition is the price we pay in order to obtain the uniform bound of (10).

²For sequences a_n and b_n , we write $a_n \lesssim b_n$ if there exists constant c such that $a_n \leq cb_n$ for all n sufficiently large.

185 *Remark 2.* While taking the radius of the neighborhood graph $r \rightarrow 0$ as $n \rightarrow \infty$ —and thereby
 186 ensuring $G_{n,r}$ is sparse—is computationally attractive, the presence of a factor of $\frac{\log^2(1/r)}{r}$ in $\kappa(\mathcal{C})$
 187 unfortunately prevents us from making claims about the behavior of PPR in this regime. Although
 188 the restriction to a kernel function fixed in n is standard for theoretical analysis of spectral clustering
 189 [20, 28], it is an interesting question whether PPR exhibits some degeneracy over r -neighborhood
 190 graphs as $r \rightarrow 0$, or if this is merely looseness in our upper bounds.

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

$$\text{for all } u \in V, \quad p(u) - \epsilon D_{uu} \leq p^{(\epsilon)}(u) \leq p(u). \quad (12)$$

196 Application of (12) within the proofs of Theorems 1 and 2 leads to analogous results which hold with
 197 respect to $p^{(\epsilon)}$. We formally state and prove this fact in the supplementary material.

198 3 Analysis

199 Given an arbitrary graph $G = (V, E)$ and subset $S \subseteq G$, Zhu et al. [30] bound the volume of $\widehat{C} \setminus S$
 200 and $S \setminus \widehat{C}$ in terms of the normalized cut and inverse mixing time of S . The key to deriving the
 201 algorithmic results of the previous section is therefore to show that the geometric conditions (A1) -
 202 (A4) translate to meaningful bounds on the normalized cut and inverse mixing time of $\mathcal{C}_\sigma[X]$ in $G_{n,r}$.
 203 In doing so, we expose how some of the geometric conditions introduced in Section 2 contribute to
 204 the difficulty of the clustering problem.

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

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

$$n \geq \frac{(2 + \epsilon)^2 \log(3/\delta)}{\epsilon^2} \left(\frac{25}{6\lambda_\sigma^2 \nu(\mathcal{C}_\sigma) \nu_d r^d} \right)^2, \quad (13)$$

209 then

$$\frac{\Phi_{n,r}(\mathcal{C}_\sigma[X])}{r} \leq c_1 \frac{d}{\sigma} \frac{\lambda}{\lambda_\sigma} \frac{(\lambda_\sigma - c_0 \frac{r^\gamma}{\gamma+1})}{\lambda_\sigma} + \epsilon, \quad (14)$$

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

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

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

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

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

$$\tau_\infty(G[S]) = \min \left\{ t : \frac{\pi(u) - q_v^{(t)}(u)}{\pi(u)} \leq \frac{1}{4}, \text{ for } u, v \in V \right\}. \quad (15)$$

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

224 **Theorem 4.** Fix $\lambda > 0$, and let $\mathcal{C} \in \mathbb{C}_f(\lambda)$ satisfy Assumptions (A1) and (A4) for some
225 $\sigma, \lambda_\sigma, \Lambda_\sigma, \rho, L > 0$. Then, for any $0 < r < \sigma/2\sqrt{d}$, with probability 1

$$\limsup_{n \rightarrow \infty} \tau_\infty(G_{n,r}[\mathcal{C}_\sigma[X]]) \leq c_2 \frac{\Lambda_\sigma^4 d^3 \rho^2 L^2}{\lambda_\sigma^4 r^2} \log^2 \left(\frac{1}{r} \right) + c_3 \quad (16)$$

226 for $c_2, c_3 > 0$ universal constants.

227 So far as we are aware, Theorem 4 is the first bound, albeit asymptotic, on the mixing time of random
228 walks over neighborhood graphs which is independent of n , the number of vertices.

229 *Remark 4.* The embedding assumption (A4) and Lipschitz parameter L play an important role
230 in proving the upper bound of Theorem 4. There is some interdependence between L and other
231 geometric parameters σ and ρ , which might lead one to hope that (A4) is non-essential. However, it
232 is not possible to eliminate this condition without incurring an additional factor of at least $(\rho/\sigma)^d$
233 in (16), achieved, for instance, when \mathcal{C}_σ is a dumbbell-like set consisting of two balls of diameter ρ
234 linked by a cylinder of radius σ . [2, 1] develop theory regarding biLipschitz deformations of convex
235 sets, wherein it is observed that star-shaped sets as well as half-moon shapes of the type we consider
236 in Section 4 both satisfy (A4) for reasonably small values of L .

237 4 Experiments

238 We provide numerical experiments to investigate the tightness of our bounds in on normalized cut
239 and mixing time of $\mathcal{C}_\sigma[X]$, and examine the performance of PPR on the 'two moons' dataset. For
240 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 (15), as we vary the diameter ρ and thickness σ of a cluster \mathcal{C} .
244 Panels (a) and (b) show the resulting empirical clusters for two different values of ρ and σ .

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.³ 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 ρ . The story in panels (c) and (e) is less obvious. We note that while, broadly speaking, the
249 trends do not appear to match, this gap between theory and empirical results seems largest when σ
250 and ρ are approximately equal. As the ratio ρ/σ grows, we see the slopes of the empirical curves
251 becoming more similar to those predicted by theory.

252 **Empirical behavior of PPR** To drive home the main implications of Theorems 1 and 2, in Figure
253 2 we show the behavior of PPR, normalized cut, and the density clustering algorithm of [6] on the
254 well known 'two moons' dataset (with added 2d Gaussian noise), considered a prototypical success
255 story for spectral clustering algorithms. The first column consists of the empirical density clusters
256 C_n and C'_n for a particular threshold λ of the density function; the second column shows the cluster
257 recovered by PPR; the third column shows the global minimum normalized cut, computed according
258 to the algorithm of [26]; and the last column shows a cut of the density cluster tree estimator of
259 [6].

260 Figure 2 shows the degrading ability of PPR to recover density clusters as the two moons become
261 less well-separated. Of particular interest is the fact that PPR fails to recover one of the moons even

³Note 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.

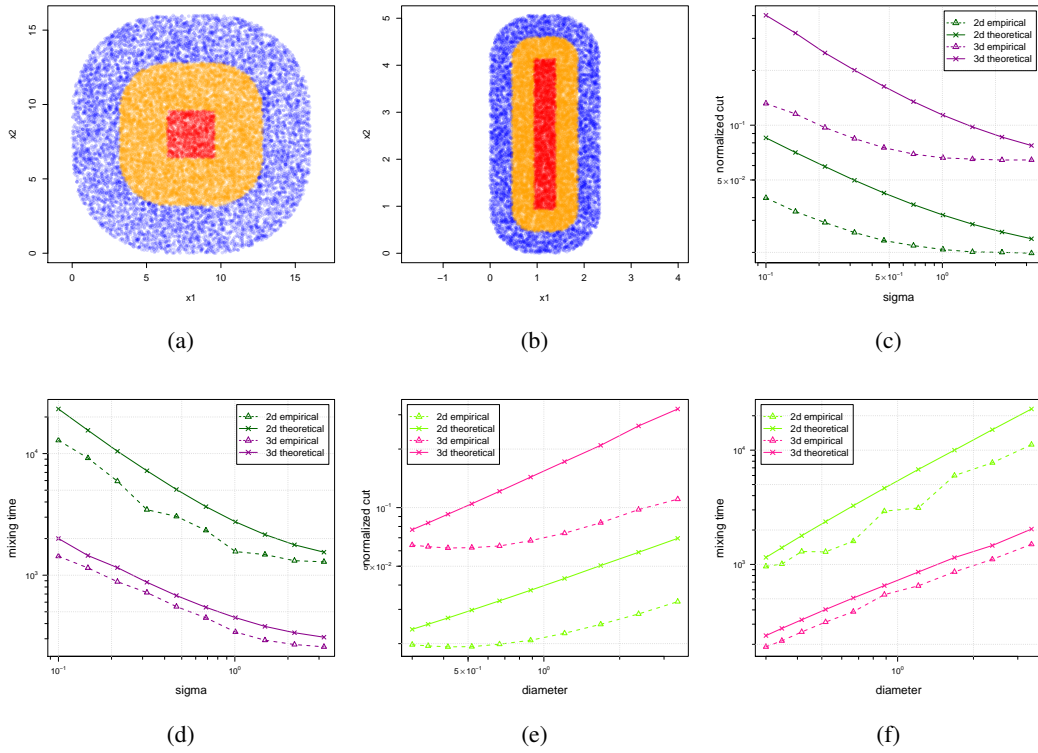


Figure 1: Samples, empirical results, and theoretical bounds for mixing time and normalized cut as diameter and thickness are varied. In (a) and (b), points in \mathcal{C} are colored in red; points in $\mathcal{C}_\sigma \setminus \mathcal{C}$ are colored in yellow; and remaining points in blue.

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

5 Discussion

For given data, there are an almost limitless number of ways to define what the 'right' clustering is. We have considered one such notion – density upper level sets – and have detailed a set of natural geometric criteria which, when appropriately satisfied, translate to provable bounds on estimation of the cluster by PPR. 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.

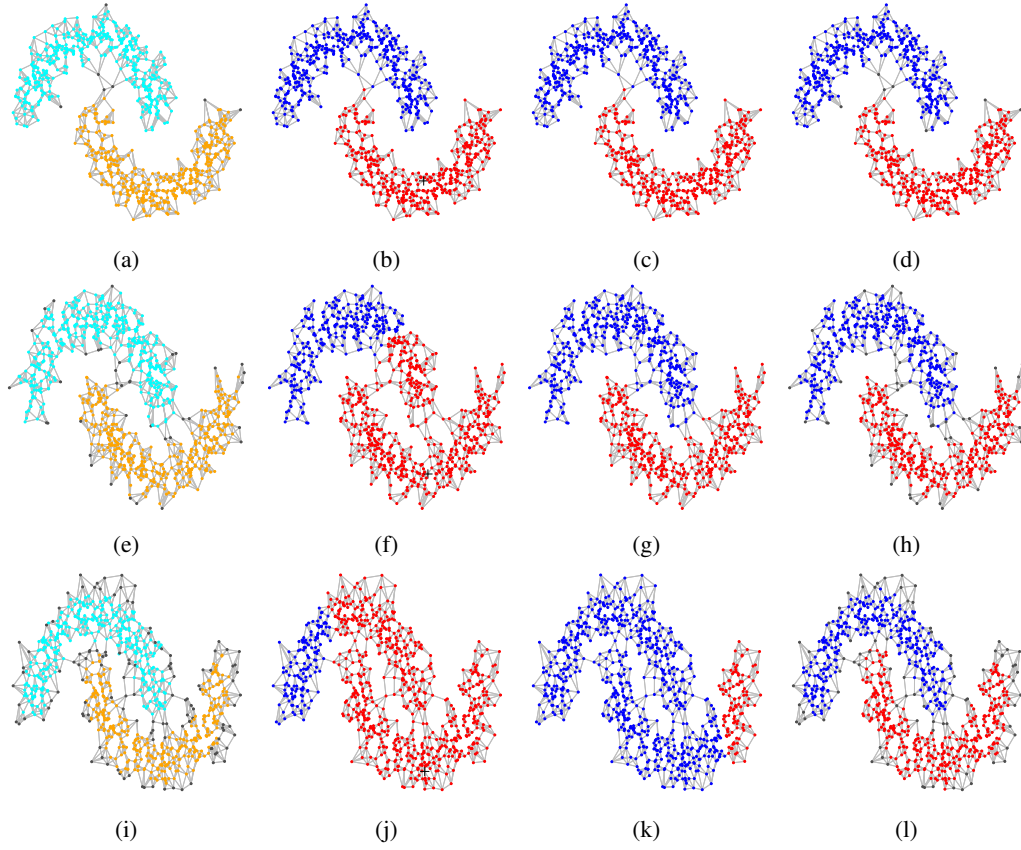


Figure 2: True density (column 1), PPR (column 2), normalized cut (column 3) and estimated density (column 4) clusters for 3 different simulated data sets. Seed node for PPR denoted by a black cross.

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] 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.
- [14] Vladimir Koltchinskii and Evarist Giné. Random matrix approximation of spectra of integral operators. *Bernoulli*, 6(1):113–167, 02 2000.
- [15] Aleksandr Petrovich Korostelev and Alexandre B Tsybakov. *Minimax theory of image reconstruction*, volume 82. 2012.
- [16] 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.
- [17] 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.
- [18] Wolfgang Polonik. Measuring mass concentrations and estimating density contour clusters—an excess mass approach. *The Annals of Statistics*, 23(3):855–881, 1995.

- 318 [19] Philippe Rigollet and Régis Vert. Optimal rates for plug-in estimators of density level sets.
319 *Bernoulli*, 15(4):1154–1178, 2009.
- 320 [20] Geoffrey Schiebinger, Martin J. Wainwright, and Bin Yu. The geometry of kernelized spectral
321 clustering. *Ann. Statist.*, 43(2):819–846, 04 2015.
- 322 [21] Tao Shi, Mikhail Belkin, and Bin Yu. Data spectroscopy: Eigenspaces of convolution operators
323 and clustering. *Ann. Statist.*, 37(6B):3960–3984, 12 2009.
- 324 [22] Aarti Singh, Clayton Scott, and Robert Nowak. Adaptive hausdorff estimation of density level
325 sets. *Ann. Statist.*, 37(5B):2760–2782, 10 2009.
- 326 [23] Daniel A Spielman and Shang-Hua Teng. Spectral sparsification of graphs. *SIAM Journal on*
327 *Computing*, 40(4):981–1025, 2011.
- 328 [24] Daniel A Spielman and Shang-Hua Teng. A local clustering algorithm for massive graphs and
329 its application to nearly linear time graph partitioning. *SIAM Journal on Computing*, 42(1):
330 1–26, 2013.
- 331 [25] Daniel A Spielman and Shang-Hua Teng. Nearly linear time algorithms for preconditioning
332 and solving symmetric, diagonally dominant linear systems. *SIAM Journal on Matrix Analysis*
333 *and Applications*, 35(3):835–885, 2014.
- 334 [26] Arthur Szlam and Xavier Bresson. Total variation, cheeger cuts. In *ICML*, pages 1039–1046,
335 2010.
- 336 [27] Alexandre B Tsybakov. On nonparametric estimation of density level sets. *The Annals of*
337 *Statistics*, 25(3):948–969, 1997.
- 338 [28] Ulrike von Luxburg, Mikhail Belkin, and Olivier Bousquet. Consistency of spectral clustering.
339 *Ann. Statist.*, 36(2):555–586, 04 2008.
- 340 [29] Xiao-Ming Wu, Zhenguo Li, Anthony M. So, John Wright, and Shih fu Chang. Learning
341 with partially absorbing random walks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q.
342 Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 3077–3085.
343 Curran Associates, Inc., 2012.
- 344 [30] Zeyuan Allen Zhu, Silvio Lattanzi, and Vahab S Mirrokni. A local algorithm for finding
345 well-connected clusters. In *ICML (3)*, pages 396–404, 2013.