

---

# 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.

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

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

For a level  $\beta > 0$  and a target volume  $\text{vol}_0 > 0$ , we define a  $\beta$ -sweep cut of  $p = (p_u)_{u \in V}$  as

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

For a set  $S \subseteq V$  with complement  $S^c = V \setminus S$ , the normalized cut is defined to be

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

where  $\text{cut}(S; G_{n,r}) = \sum_{u \in S, v \in S^c} \mathbf{A}_{uv}$  and  $\text{vol}(S; G_{n,r}) := \sum_{u \in S} \mathbf{D}_{uu}$ . Having computed sweep cuts  $S_\beta$  over a range  $\beta \in (\frac{1}{40}, \frac{1}{11})^1$ , we then output the cluster estimate  $\hat{C} = S_{\beta^*}$  which has minimum normalized cut  $\Phi(S_{\beta^*}; G_{n,r})$ . 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  $S_\beta$  as in (2).
- 4: Return  $\hat{C} = S_{\beta^*}$ , where

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


---

**Estimation of density clusters.** 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*. The symmetric set difference between estimated and empirical cluster is perhaps the most frequently used metric to quantify cluster estimation error [15, 18, 19].

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

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

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

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

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

with probability tending to 1.

---

<sup>1</sup>The 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].

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

- 63 1. In Section 2, we introduce a set of natural geometric conditions, formalize a measure of  
 64 difficulty based on these geometric conditions, and show that when properly initialized, the  
 65 symmetric set difference of Algorithm 1 is upper bounded by this difficulty measure.
- 66 2. We further show that if the density cluster  $\mathcal{C}$  is particularly well-conditioned, Algorithm  
 67 1 will perform consistent density cluster estimation in the sense of (5). In Corollary 1,  
 68 we establish that both metrics of cluster estimation can be bounded with respect to an  
 69 approximate form of PPR, which can be efficiently computed.
- 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 earlier, a few related lines of work are worth  
 79 highlighting. Similar in spirit to our results are the works [21, 20], who study the consistency of  
 80 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 ensure 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 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 measure  $\kappa(\mathcal{C})$  which encodes  
 96 the difficulty PPR will have estimating  $\mathcal{C}$ . We motivate this measure, and the underlying geometric  
 97 conditions, by giving density cluster estimation guarantees for Algorithm 1 in terms of  $\kappa(\mathcal{C})$ .

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

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

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

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

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

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

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

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

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

117 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-  
 118 connected, while being poorly connected to the rest of  $X$ . The cluster separation (A2) and low  
 119 noise density (A3) conditions guarantee low connectivity between  $\mathcal{C}_\sigma[X]$  and  $X \setminus \mathcal{C}_\sigma[X]$  in  $G_{n,r}$ ,  
 120 whereas (A1) and (A4) ensure high connectivity within  $\mathcal{C}_\sigma[X]$ . It may not be immediately obvious  
 121 how (A4) contributes to geometric conditioning. For now, we observe merely that random walks  
 122 will mix slowly over sets with large diameter, and comment on this condition in more detail in  
 123 Section 3. Finally, (A5) is a relatively harmless technical condition, merely excluding the case where  
 124  $\text{vol}(\mathcal{C}_\sigma[X]; G_{n,r}) > \text{vol}(X; G_{n,r})/2$ .

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

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

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

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

131 At first glance (7) may appear mysterious, but as will be shown in Section 3,  $\Phi_u(\theta)$  and  $\Psi_u(\theta)$  are  
 132 merely upper bounds on the normalized cut and inverse mixing time of  $\mathcal{C}_\sigma[X]$  in  $G_{n,r}$ . In [30],  
 133 building on the work of [4] and others, it is shown that the ratio of normalized cut to inverse mixing  
 134 time is a fundamental quantity governing the clustering performance of PPR on a general graph.  $\kappa(\mathcal{C})$   
 135 upper bounds this ratio for an empirical density cluster over the neighborhood graph  $G_{n,r}$ , and is  
 136 therefore a natural criterion to measure difficulty of the density clustering task.

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

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

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

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

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

142 **Definition 4.** If the input parameters to Algorithm 1 satisfy (8) for some well-conditioned density  
 143 cluster  $\mathcal{C}$ , we say the algorithm is well-initialized.

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

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

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

150 To translate (9) into meaningful bounds on the symmetric set difference  $\Delta(\mathcal{C}_\sigma[X], \widehat{C})$ , we wish to  
 151 preclude vertices  $x \in X$  from having arbitrarily small degree. To do so, we make some regularity  
 152 assumptions on  $\mathcal{X} := \text{supp}(f)$ .

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

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

157 **Theorem 1.** Fix  $\lambda > 0$ , let  $\mathcal{C} \in \mathbb{C}_f(\lambda)$  be a  $\kappa$ -well conditioned density cluster (with respect to some  
 158  $\theta$ ), and additionally assume  $f$  satisfies (A5). Then, there exists universal constant  $c_4 > 0$  such that  
 159 with probability tending to one as  $n \rightarrow \infty$ ,

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

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

163 Neither (9) nor Theorem 1 imply consistent density cluster estimation in the sense of (5). This notion  
 164 of consistency requires a uniform bound over  $p$  for all  $u \in \mathcal{C}, w \in \mathcal{C}'$

$$\frac{p_w}{\mathbf{D}_{ww}} \leq \frac{1}{40\text{vol}_0} < \frac{1}{11\text{vol}_0} \leq \frac{p_u}{\mathbf{D}_{uu}}. \quad (11)$$

165 so that any sweep cut  $S_\beta$  for  $\beta\text{vol}_0 \in [1/40, 1/11]$  (i.e. any sweep cut considered by Algorithm 1)  
 166 will fulfill both conditions laid out in (5). In Theorem 2, we show that a sufficiently small upper  
 167 bound on  $\kappa(\mathcal{C})$  ensures such a gap exists with probability one as  $n \rightarrow \infty$ , and therefore guarantees  $\widehat{C}$   
 168 will be a consistent estimator. As was the case before, we wish to preclude arbitrarily low degree  
 169 vertices, this time for points  $x \in \mathcal{C}'[X]$ .

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

172 **Theorem 2.** Fix  $\lambda > 0$ , let  $\mathcal{C} \in \mathbb{C}_f(\lambda)$  be a  $\kappa$ -well conditioned cluster (with respect to some  $\theta$ ), and  
 173 additionally assume (A6) holds. If Algorithm 1 is well-initialized, there exists universal constant  
 174  $c_5 > 0$  such that if

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

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

176 A few remarks are in order.

177 **Remark 1.** We note that the restriction on  $\kappa(\mathcal{C})$  imposed by (12) results in a misclassification rate on  
 178 the order of  $r^d$ . (See Theorem 1). In plain terms, we are able to recover a density cluster  $\mathcal{C}$  in the  
 179 sense of (5) only when we can guarantee a very small fraction of points will be misclassified. This  
 180 strong condition is the price we pay in order to obtain the uniform bound of 11.

<sup>2</sup>For 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.

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

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

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

192 Application of (13) within the proofs of Theorems 1 and 2 leads to analogous results which hold with  
 193 respect to  $p^{(\epsilon)}$ .

194 **Corollary 1.** Fix  $\lambda > 0$ , and let  $\mathcal{C} \in \mathbb{C}_f(\lambda)$  be a  $\kappa$ -well-conditioned cluster (with respect to some  
 195  $\theta$ ). Choose input parameters  $\alpha, r, \text{vol}_0, v$  to be well-initialized in the sense of (8), set  $\epsilon = \frac{1}{20\text{vol}_0}$ ,  
 196 and modify Algorithm 1 to compute the aPPR vector  $p^{(\epsilon)}$  rather than the exact PPR vector  $p$ , with  
 197 resulting output  $\hat{\mathcal{C}}$ .

- 198 1. If (A5) holds, then (10) is still a valid upper bound for the misclassification error of  $\hat{\mathcal{C}}$ .
- 199 2. If (A6) and (12) hold, then  $\hat{\mathcal{C}} \subseteq X$  is a consistent estimator for  $\mathcal{C}$ , in the sense of Definition  
 200 2.

### 201 3 Analysis

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

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

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

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

212 then

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

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

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

**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  $v$ . Also write  $\pi = (\pi_u)_{u \in S}$  for the stationary distribution of this random walk. (As  $\mathbf{W}_S$  is the transition matrix of a lazy random walk, it is well-known that a unique stationary distribution exists and is given by  $\pi_u = (\mathbf{D}_S)_{uu} / \text{vol}(S; G[S])$ , where  $\mathbf{D}_S$  is the degree matrix of  $G[S]$ .)

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

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

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  $\sigma, \lambda_\sigma, \Lambda_\sigma, D, L > 0$ . Then, for any  $0 < r < \sigma/2\sqrt{d}$ , with probability 1

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

for  $c_2, c_3 > 0$  universal constants.

So far as we are aware, Theorem 4 is the first bound 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  $\sigma$  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 (17), achieved, for instance, when  $\mathcal{C}_\sigma$  is a dumbbell-like set consisting of two balls of diameter  $D$  linked by a cylinder of radius  $\sigma$ . [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 in 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 shows these theoretical bounds compared to the empirical quantities (3) and (16), as we vary the diameter  $D$  and thickness  $\sigma$  of the cluster  $\mathcal{C}$ . Panels (a) and (b) show the empirical clusters under consideration for two different values of  $D, \sigma$ .

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

<sup>3</sup>Note 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.

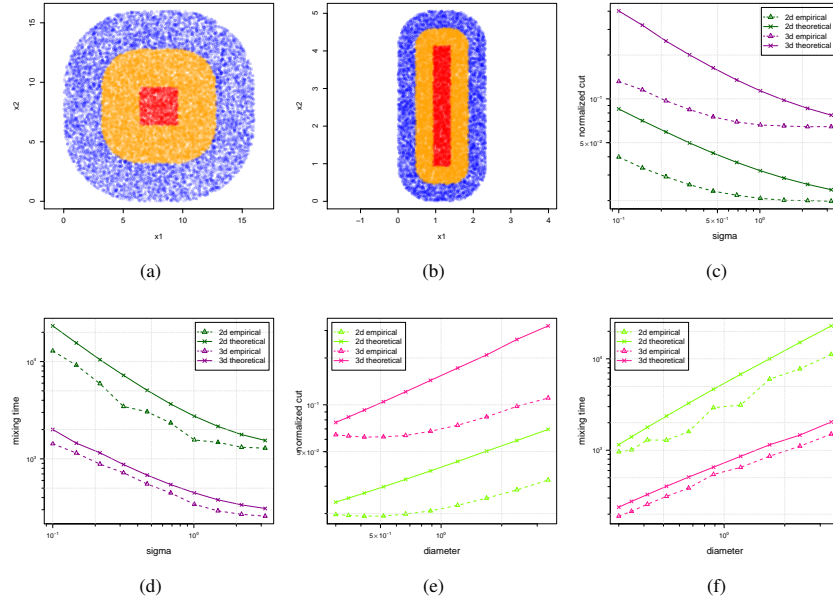


Figure 1: Samples, empirical results, and theoretical bounds for mixing time and normalized cut as diameter and thickness are varied. In (a) and (b), points in  $C$  are colored in red; points in  $C_\sigma \setminus C$  are colored in yellow; and remaining points in blue.

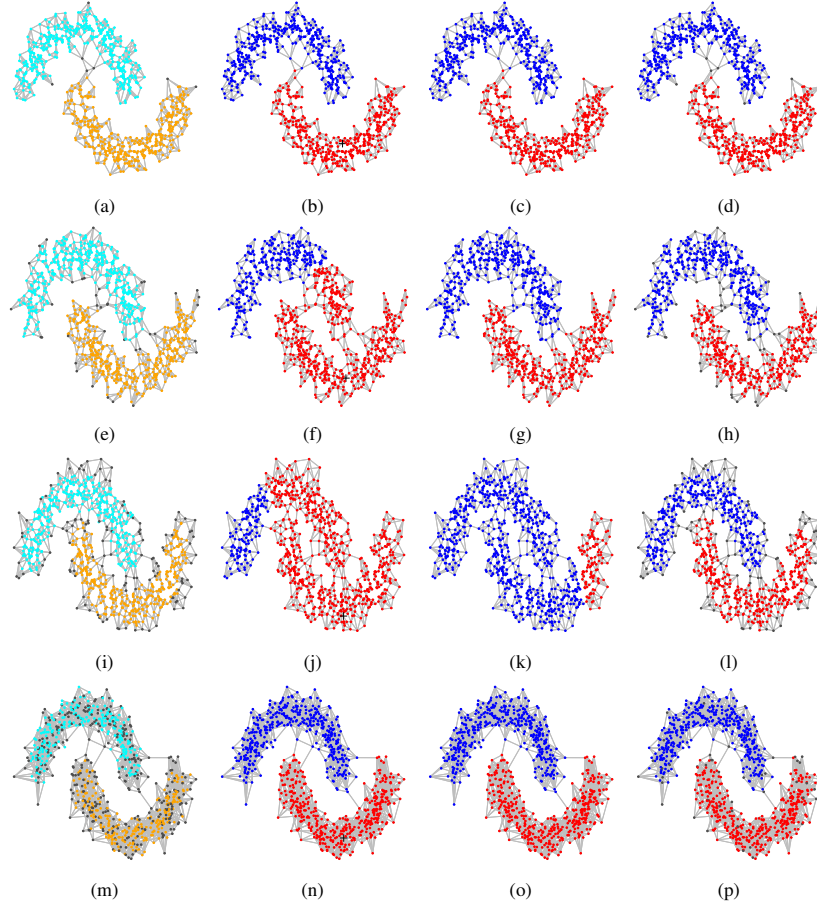


Figure 2: True density (column 1), PPR (column 2), normalized cut (column 3) and estimated density (column 4) clusters for 4 different simulated data sets. Seed node for PPR denoted by a black cross.



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

263 Rows 1-3 show the degrading ability of PPR to recover density clusters as the two moons become  
 264 less salient. Of particular interest is the fact that PPR fails to recover one of the moons even when  
 265 normalized cut still succeeds in doing so, and that a density clustering algorithm recovers a moon  
 266 even when both PPR and normalized cut fail. In the fourth row,  $10d$  Gaussian noise was added. The  
 267 gray dots in  $(m)$  (as in  $(a)$ ,  $(e)$  and  $(i)$ ) are observations in low-density regions. While the PPR sweep  
 268 cut  $(n)$  has relatively high symmetric set difference with the chosen density cut, it still recovers  $C_n$   
 269 in the sense of Definition 2.

## 270 5 Discussion

271 For given data, there are an almost limitless number of ways to define what the 'right' clustering is.  
 272 We have considered one such notion – density upper level sets – and have detailed a set of natural  
 273 geometric criteria which, when appropriately satisfied, translate to provable bounds on estimation  
 274 of the cluster by PPR. We do not provide a theoretical lower bound showing that our geometric  
 275 conditions are required for successful recovery on an upper level set. Although we investigate the  
 276 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] 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.

- 324 [19] Philippe Rigollet and Régis Vert. Optimal rates for plug-in estimators of density level sets.  
325 *Bernoulli*, 15(4):1154–1178, 2009.
- 326 [20] Geoffrey Schiebinger, Martin J. Wainwright, and Bin Yu. The geometry of kernelized spectral  
327 clustering. *Ann. Statist.*, 43(2):819–846, 04 2015.
- 328 [21] Tao Shi, Mikhail Belkin, and Bin Yu. Data spectroscopy: Eigenspaces of convolution operators  
329 and clustering. *Ann. Statist.*, 37(6B):3960–3984, 12 2009.
- 330 [22] Aarti Singh, Clayton Scott, and Robert Nowak. Adaptive hausdorff estimation of density level  
331 sets. *Ann. Statist.*, 37(5B):2760–2782, 10 2009.
- 332 [23] Daniel A Spielman and Shang-Hua Teng. Spectral sparsification of graphs. *SIAM Journal on*  
333 *Computing*, 40(4):981–1025, 2011.
- 334 [24] Daniel A Spielman and Shang-Hua Teng. A local clustering algorithm for massive graphs and  
335 its application to nearly linear time graph partitioning. *SIAM Journal on Computing*, 42(1):  
336 1–26, 2013.
- 337 [25] Daniel A Spielman and Shang-Hua Teng. Nearly linear time algorithms for preconditioning  
338 and solving symmetric, diagonally dominant linear systems. *SIAM Journal on Matrix Analysis*  
339 *and Applications*, 35(3):835–885, 2014.
- 340 [26] Arthur Szlam and Xavier Bresson. Total variation, cheeger cuts. In *ICML*, pages 1039–1046,  
341 2010.
- 342 [27] Alexandre B Tsybakov. On nonparametric estimation of density level sets. *The Annals of*  
343 *Statistics*, 25(3):948–969, 1997.
- 344 [28] Ulrike von Luxburg, Mikhail Belkin, and Olivier Bousquet. Consistency of spectral clustering.  
345 *Ann. Statist.*, 36(2):555–586, 04 2008.
- 346 [29] Xiao-Ming Wu, Zhenguo Li, Anthony M. So, John Wright, and Shih fu Chang. Learning  
347 with partially absorbing random walks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q.  
348 Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 3077–3085.  
349 Curran Associates, Inc., 2012.
- 350 [30] Zeyuan Allen Zhu, Silvio Lattanzi, and Vahab S Mirrokni. A local algorithm for finding  
351 well-connected clusters. In *ICML (3)*, pages 396–404, 2013.