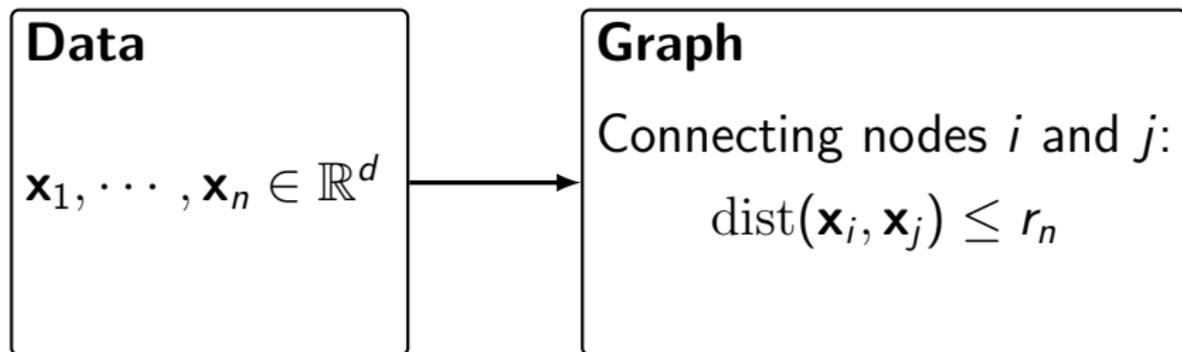


# Recent Developments in Random Geometric Graphs and Their Applications

Xiucui Ding

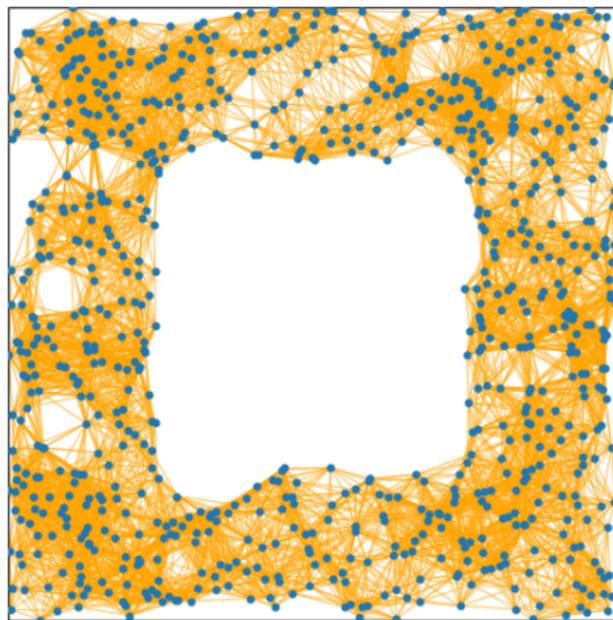
Department of Statistics  
University of California, Davis

Joint work with Yichen Hu (UC Davis)/Haixiao Wang (UW–Madison)

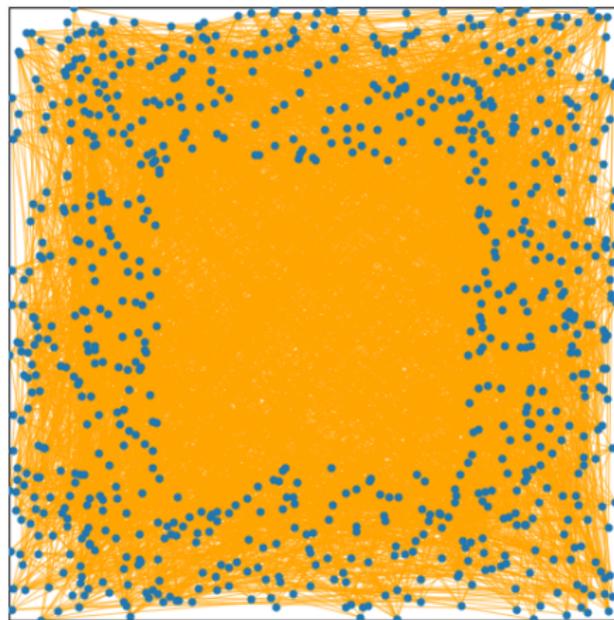


# Example

$\mathbf{x}_1, \dots, \mathbf{x}_n$  be i.i.d. uniform on  $[0, 1]^2 \setminus [1/4, 3/4]^2$



(a) RGG



(b) ERG

# Application: manifold learning

Most manifold learning algorithms are based on graph structures, e.g., Laplacian eigenmap. It can be summarized in the following framework:

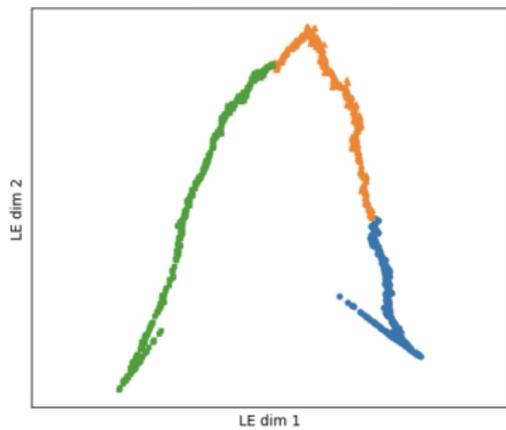
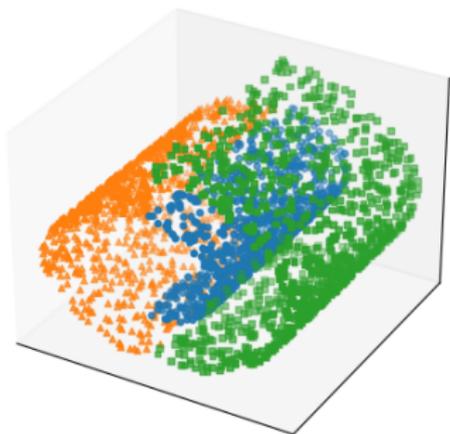
---

## Algorithm

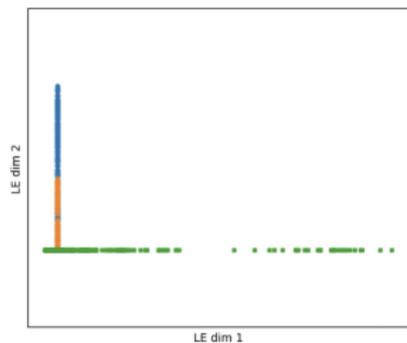
---

- 1: Input: nodes  $1, \dots, n$  associated with features  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$
  - 2: Connect node  $i$  and node  $j$  if  $\text{dist}(\mathbf{x}_i, \mathbf{x}_j) \leq r_n$ .
  - 3: Assign a weight  $g(\text{dist}(\mathbf{x}_i, \mathbf{x}_j)/r_n)$  if node  $i$  and node  $j$  are connected.
  - 4: Construct affinity matrix  $\mathbf{K}$  and degree matrix  $\mathbf{D}$ .
  - 5: Compute random-walk Laplacian matrix  $\mathbf{L}_{rw} = \mathbf{I} - \mathbf{D}^{-1}\mathbf{K}$  and its eigenvalues  $0 \leq \lambda_1 \leq \dots \leq \lambda_n$  with eigenvectors  $\nu_1, \dots, \nu_n$ .
  - 6: Embed  $\mathbf{x}_j \rightarrow (\nu_{2,j}, \dots, \nu_{M+1,j})^\top$ .
- 

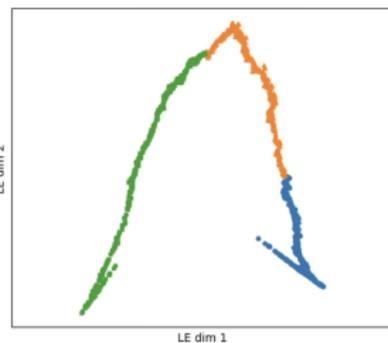
This uses the spectrum of an RGG via  $\mathbf{L}_{rw}$ .



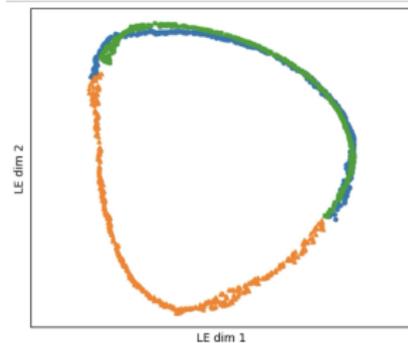
# Impact of $r_n$



(a)  $r_n = 0.05$



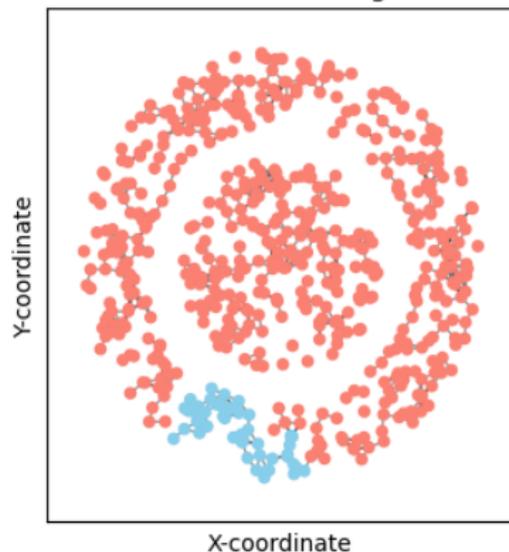
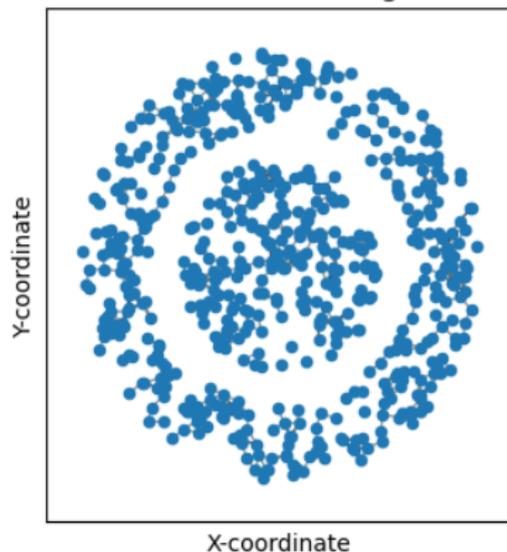
(b)  $r_n = 0.2$



(c)  $r_n = 0.5$

# Clustering

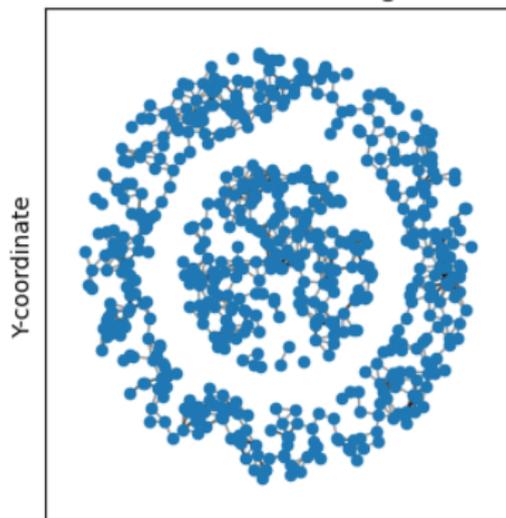
Graph is almost disconnected when  $\text{RADIUS} = 0.25$



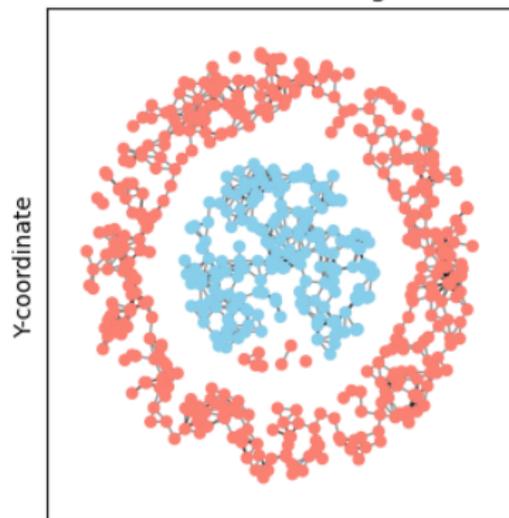
# Clustering

Exact recovery when  $\text{RADIUS} = 0.3$

Before Clustering



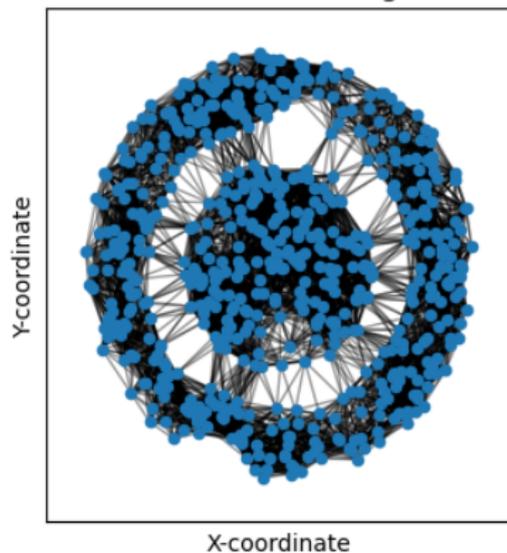
After Clustering



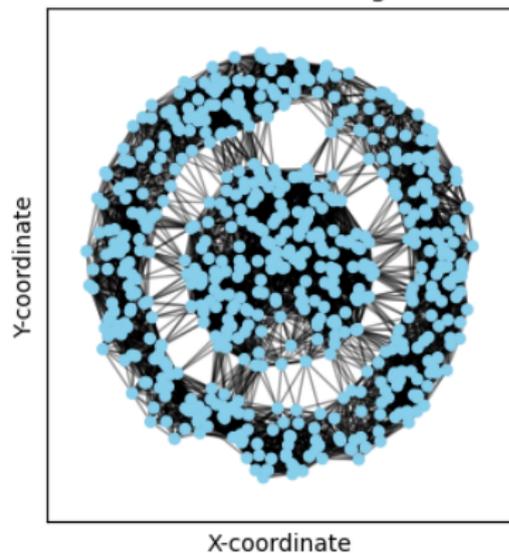
# Clustering

Only one connected component when  $\text{RADIUS} = 0.8$

Before Clustering



After Clustering



- When  $r_n \equiv r$ ,  $g \equiv 1$ ,  $\ell_\infty$  metric is used and the sampling distribution is uniform on  $[-1, 1]^d$ , the edge eigenvalues of  $\mathbf{L}_{rw}$  are studied by Adhikari et al. (2022)
- $d = 2, r = 1$  :

$$\lambda_1(\mathbf{L}_{rw}) = 1, \lambda_2(\mathbf{L}_{rw}) \approx 1/2, \lambda_3(\mathbf{L}_{rw}) \approx 1/2, \lambda_4(\mathbf{L}_{rw}) \approx 3/4.$$

- When  $r_n = o(1)$ ,  $\ell_2$  metric is used and the sampling distribution is compactly supported with density bounded from above and below, García Trillos et al. (2020) and García Trillos et al. (2021) show the (scaled) edges eigenvalues of  $\mathbf{L}_{rw}$  converge to those of a weight Laplace-Beltrami operator.

# Our results

- Gaussian samples  $\mathbf{x}_j \sim \mathcal{N}_d(0, \Sigma)$  with  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ ;
- $\ell_p$  distance, i.e.,  $\|\cdot\|_p$ , with  $1 \leq p \leq \infty$ ;
- radius  $r_n$ ,

$$n^{-\frac{1}{d+4} + \varepsilon} \ll r_n \ll n^{-\varepsilon},$$

for some small  $0 < \varepsilon < \frac{1}{2(d+4)}$ ;

- weight function  $g \in C^2([0, \infty))$  and bounded from above and below by positive constants on  $[0, 1]$

$$\mathbf{K}(i, j) = \mathbf{1}(\|\mathbf{x}_i - \mathbf{x}_j\|_p \leq r_n) g\left(\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_p}{r_n}\right).$$

- scaling constants

$$m_\ell = \int_{\|\mathbf{u}\|_p \leq 1} u_i^\ell g(\|\mathbf{u}\|_p) d\mathbf{u}$$

- our objective matrix is

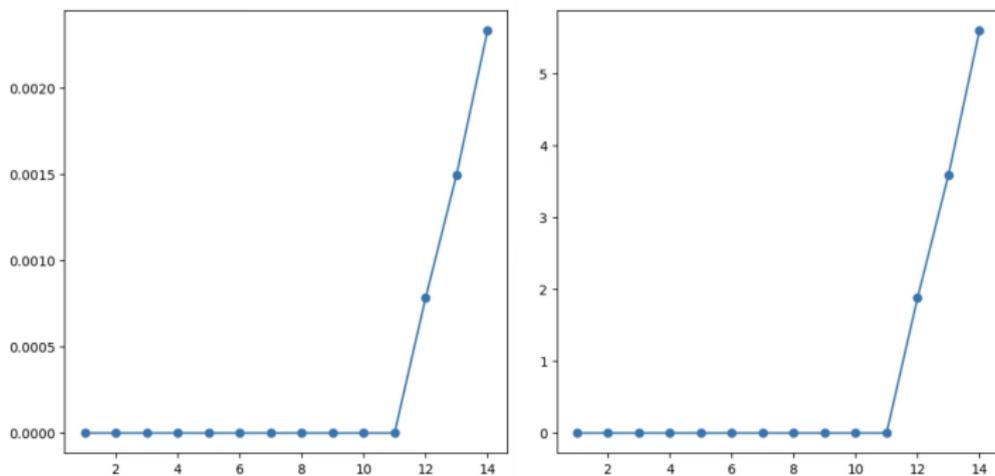
$$\mathbf{L} = \frac{2m_0}{m_2 r_n^2} (\mathbf{I} - \mathbf{D}^{-1} \mathbf{K}) = \frac{2m_0}{m_2 r_n^2} \mathbf{L}_{rw},$$

where  $\mathbf{K}(i, j) = \mathbf{1}(\|\mathbf{x}_i - \mathbf{x}_j\|_p \leq r_n) g\left(\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_p}{r_n}\right)$ .

- eigenvalues of  $\mathbf{L}$

$$\lambda_1(\mathbf{L}) \leq \lambda_2(\mathbf{L}) \leq \dots \leq \lambda_n(\mathbf{L})$$

# Setup



- exclude small eigenvalues

$$K_0 = \max_j \{j : \lambda_j(\mathbf{L}) \leq \delta\},$$

for some small  $0 < \delta < \min\{\sigma_1^{-2}, \dots, \sigma_d^{-2}\}$ .

- Hilbert space

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R} \mid \int_{\mathbb{R}^d} f^2(\mathbf{x}) \varrho^2(\mathbf{x}) d\mathbf{x} < \infty \right\}$$

with weight

$$\varrho(\mathbf{x}) = \exp\left(-\sum_{i=1}^d \frac{x_i^2}{2\sigma_i^2}\right)$$

and inner product

$$\langle f_1, f_2 \rangle = \int_{\mathbb{R}^d} f_1(\mathbf{x}) f_2(\mathbf{x}) \varrho^2(\mathbf{x}) d\mathbf{x}$$

- weighted Laplace–Beltrami operator

$$\Delta_\varrho f = -\frac{1}{\varrho^2} \operatorname{div}(\varrho^2 \nabla f) = \sum_{i=1}^d \frac{\partial^2 f}{\partial x_i^2} - \frac{2x_i}{\sigma_i^2} \frac{\partial f}{\partial x_i}$$

- eigenvalues of  $\Delta_\varrho$  on  $(\mathcal{F}, \langle \cdot, \cdot \rangle)$

$$\mu_1(\Delta_\varrho) \leq \mu_2(\Delta_\varrho) \leq \dots$$

# Main results (I)

## Theorem

For the weighted Laplace-Beltrami operator  $\Delta_\rho$  on  $(\mathcal{F}, \langle \cdot, \cdot \rangle)$ , its spectrum is discrete and consists of eigenvalues  $\sum_{i=1}^d 2(k_i - 1)/\sigma_i^2$  with corresponding orthonormal eigenfunctions

$$\prod_{i=1}^d \psi_{i, k_i - 1}(x_i), \quad k_1, \dots, k_d \in \mathbb{N}^+,$$

where

$$\psi_{i, k_i - 1}(x_i) = \left( \frac{1}{\sqrt{\pi} (k_i - 1)! 2^{k_i - 1} \sigma_i} \right)^{1/2} H_{k_i - 1}(x_i / \sigma_i),$$

and  $H_{k_i - 1}$  is the  $(k_i - 1)$ -th physicist's Hermite polynomial.

# Main results (II)

## Theorem

Let  $M$  be any fixed integer. For sufficiently large  $n$ , with  $1 - o(1)$  probability

$$K_0 \leq n^{1 - \frac{2}{(d+2)^2 + \eta}} \quad \text{for some small } \eta > 0.$$

Moreover,

$$\lambda_{K_0+k}(\mathbf{L}) = \mu_{k+1}(\Delta_\varrho) + o(1), \quad \text{for } k = 1, \dots, M.$$

- scaling  $2m_0/(m_2 r_n^2)$
- range of  $r_n$  as  $n^{-\frac{1}{d+4} + \varepsilon} \ll r_n \ll n^{-\varepsilon}$
- bound of  $K_0$

## Main results (III)

For  $k_1, k_2, \dots, k_d \in \mathbb{N}^+$ , denote a sequence

$$\left\{ \sum_{i=1}^d \frac{2(k_i - 1)}{\sigma_i^2} \right\},$$

ordered by  $a_1 \leq a_2 \leq \dots$ .

### Proposition

Under mild assumptions,  $\mu_j(\Delta_\rho) = a_j$ . Consequently, for sufficiently large  $n$  and any fixed  $M$ , with probability  $1 - o(1)$ ,

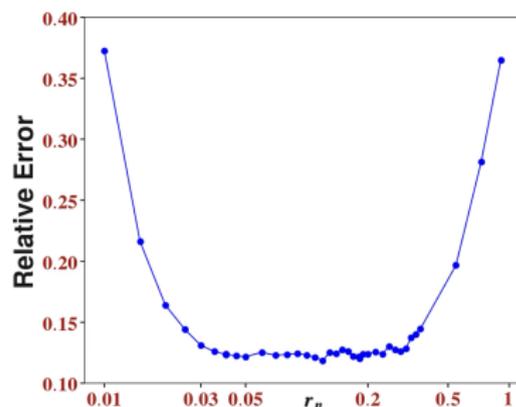
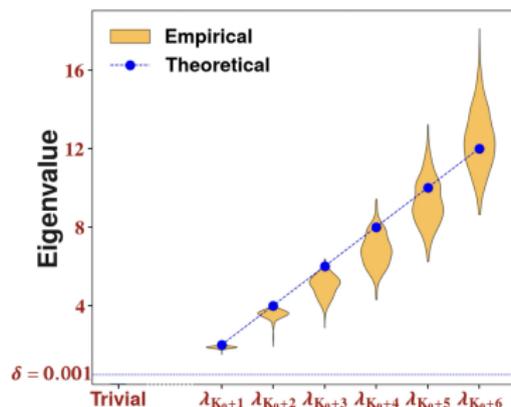
$$\lambda_{K_0+k}(\mathbf{L}) = a_{k+1} + o(1), \quad \text{for } k = 1, \dots, M.$$

# Examples

- $d = 1, \sigma_1 \equiv \sigma$ . For  $1 \leq k \leq M$ ,

$$\lambda_{K_0+k}(\mathbf{L}) = \frac{2k}{\sigma^2} + o_{\mathbb{P}}(1).$$

- $\sigma = 1, g \equiv 1, n = 5000$



$$\text{RE} = \frac{1}{M} \sum_{k=1}^M |\lambda_{K_0+k}(\mathbf{L}) - \mu_{k+1}(\Delta_{\varrho})| / \mu_{k+1}(\Delta_{\varrho})$$

# Examples

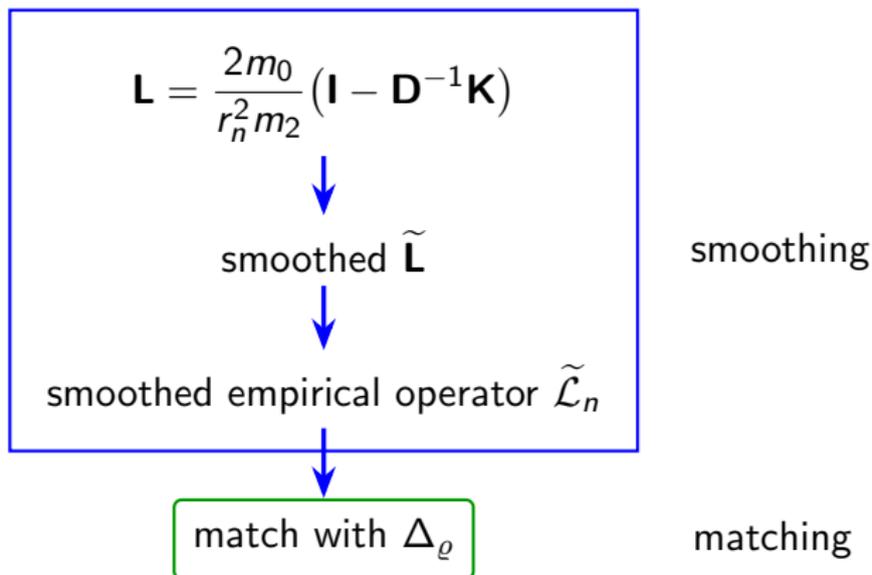
- $d = 2, \sigma_1 = \sigma_2 = \sigma$ .

$$\lambda_{K_0+k}(\mathbf{L}) = \frac{2}{\sigma^2} \min \left\{ l : \sum_{s=1}^l s \geq k+1, l \in \mathbb{N} \right\} - \frac{2}{\sigma^2} + o_{\mathbb{P}}(1).$$

- $d = 2, \sigma_1^2 = \sigma^2$  and  $\sigma_2^2 = 2\sigma^2$ .

$$\lambda_{K_0+k}(\mathbf{L}) = \frac{1}{\sigma^2} \min \left\{ l : \sum_{s=0}^l (\lfloor s/2 \rfloor + 1) \geq k+1, l \in \mathbb{N} \right\} + o_{\mathbb{P}}(1).$$

- smoothing-matching



- $\mathbf{L} \rightarrow \tilde{\mathbf{L}}$

$$\mathbf{K}(i, j) = \mathbf{1}(\|\mathbf{x}_i - \mathbf{x}_j\|_p \leq r_n) g(\|\mathbf{x}_i - \mathbf{x}_j\|_p / r_n)$$

$$\tilde{\mathbf{K}}(i, j) = s(\alpha_n(r_n^2 - \|\mathbf{x}_i - \mathbf{x}_j\|_p^2)) g^*(\|\mathbf{x}_i - \mathbf{x}_j\|_p / r_n)$$

where  $s(t) = 1/(1 + \exp(-t))$ ,  $\alpha_n \gg n^{7/2}$  and  $g^*(t)$  is some extension of  $g(t)$  from  $[0, 1]$  to  $[0, \infty)$ .

- $\|\mathbf{L} - \tilde{\mathbf{L}}\| = o_{\mathbb{P}}(1)$  under  $\alpha_n \gg n^{7/2}$

$$\lambda_j(\mathbf{L}) = \lambda_j(\tilde{\mathbf{L}}) + o_{\mathbb{P}}(1), \quad 1 \leq j \leq n.$$

- $\tilde{\mathbf{L}} \rightarrow \tilde{\mathcal{L}}_n$

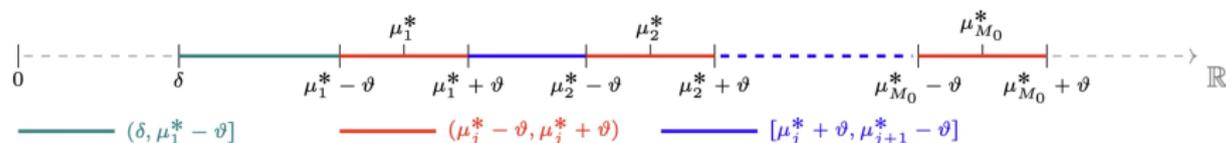
$$\tilde{\mathcal{L}}_n f(\mathbf{x}) = \frac{2m_0}{r_n^2 m_2} \frac{\frac{1}{n} \sum_{j=1}^n s(\alpha_n(r_n^2 - \|\mathbf{x} - \mathbf{x}_j\|_\rho^2)) g^*(\|\mathbf{x} - \mathbf{x}_j\|_\rho / r_n) (f(\mathbf{x}) - f(\mathbf{x}_j))}{\frac{1}{n} \sum_{j=1}^n s(\alpha_n(r_n^2 - \|\mathbf{x} - \mathbf{x}_j\|_\rho^2)) g^*(\|\mathbf{x} - \mathbf{x}_j\|_\rho / r_n)}$$

- well-defined
- Eigenvalues of  $\tilde{\mathbf{L}}$  and  $\tilde{\mathcal{L}}_n$  coincide for all  $\lambda_j(\tilde{\mathcal{L}}_n) < 2m_0/(r_n^2 m_2)$ .

$$\lambda_j(\mathbf{L}) = \lambda_j(\tilde{\mathcal{L}}_n) + o_{\mathbb{P}}(1), \quad \text{for } j \text{ such that } \lambda_j(\tilde{\mathcal{L}}_n) < \frac{2m_0}{r_n^2 m_2}.$$

# Matching

- $0 < \mu_1^* < \mu_2^* < \dots < \mu_{M_0}^*$  with multiplicities  $n_1, n_2, \dots, n_{M_0}$ ,  
 $\sum_{j=1}^{M_0} n_j \geq M$  and  $\mu_{j+1}^* - \mu_j^*$  lower bounded
- partition of real line



- counting regions
- empty regions

# Matching

- locating

$\tilde{\mathcal{L}}_n$  has at least one eigenvalue in each counting region

- counting

comparing  $-\frac{1}{2\pi i} \oint_{\Gamma} \text{Tr}(R(z, \tilde{\mathcal{L}}_n)) dz$  and  $-\frac{1}{2\pi i} \oint_{\Gamma} \text{Tr}(R(z, \Delta_{\varrho})) dz$

- empty regions

$$\lambda_j(\tilde{\mathcal{L}}_n) = \mu_{\ell}^* + o_{\mathbb{P}}(1), \text{ for } K_0 + 1 + \sum_{s=1}^{j-1} n_s \leq j \leq K_0 + \sum_{s=1}^j n_s, 1 \leq \ell \leq M_0.$$

# Key ingredient

- closeness by eigenfunctions  $\psi_j$  of  $\Delta_\varrho$

$$\left\| (\Delta_\varrho - \tilde{\mathcal{L}}_n)\psi_j \right\|_{\mathcal{F}} = o_{\mathbb{P}}(1)$$

- intermediate operator

$$\tilde{\mathcal{T}}_n f(\mathbf{x}) = \frac{2}{r_n^{d+2} m_2 \varrho(\mathbf{x})} \int_{\|\mathbf{x} - \mathbf{y}\|_\rho \leq r_n} g(\|\mathbf{x} - \mathbf{y}\|_\rho / r_n) (f(\mathbf{x}) - f(\mathbf{y})) \varrho(\mathbf{y}) d\mathbf{y}$$

- decomposition of  $(\Delta_\varrho - \tilde{\mathcal{L}}_n)\psi_j$

$$\left\| (\Delta_\varrho - \tilde{\mathcal{L}}_n)\psi_j \right\|_{\mathcal{F}} \leq \left\| (\Delta_\varrho - \tilde{\mathcal{T}}_n)\psi_j \right\|_{\mathcal{F}} + \left\| (\tilde{\mathcal{T}}_n - \tilde{\mathcal{L}}_n)\psi_j \right\|_{\mathcal{F}}$$

# Applications to manifold data clustering

- $\mathbf{x}_j = \mathbf{s}_j + \mathbf{z}_j \in \mathbb{R}^P$
- $\mathbf{s}_j$  are sampled from manifolds
- $\varrho(\mathbf{s}_j) = \sum_{k=1}^K \mathbb{P}(Y_j = k) \varrho_k(\mathbf{s}_j | k)$
- Distinct vertices  $i, j \in \mathcal{V}$  are connected if and only if  $\|\mathbf{x}_i - \mathbf{x}_j\| \leq r_n$ , where the threshold radius  $r_n > 0$  is such that

$$\mathbb{P}(\|\mathbf{s}_i - \mathbf{s}_j\| \leq r_n) = Q_{y(i)y(j)},$$

where  $Q \in [0, 1]^{K \times K}$  is some symmetric probability matrix

- Data  $\rightarrow$  Random Geometric Graph (RGG) construction with parameter selection based on spectral analysis  $\rightarrow$  Eigenvector-based embedding  $\rightarrow$  Clustering