

Manifold Clustering in the Generalized Random Dot Product Graph

John Koo

Department of YYY, University of XXX

April 27, 2022

Abstract

The text of your abstract. 200 or fewer words.

Keywords: block models, community detection, coordinate descent, latent structure models, manifold clustering, random dot product graph

1 Introduction

We define a *Bernoulli graph* as a random graph model for which edge probabilities are contained in an edge probability matrix $P \in [0, 1]^{n \times n}$, and an edge occurs between vertices i and j with probability P_{ij} . Common random graph models then impose structure on P , based on various assumptions about the way in which the data are generated, or to allow P to be estimated. One example is the Erdős-Rényi model, in which all edge probabilities are fixed, i.e., $P_{ij} = p$.

One common analysis task for graph and network data is community detection, which assumes that each vertex of a graph has a hidden community label. The goal of the analysis is then to retrieve these labels. In order to perform this analysis as a statistical inference task is to define a probability model with inherent community structure. We call such models *block models*: First, each vertex is assigned a label $z_1, \dots, z_n \in \{1, 2, \dots, K\}$ where $K \ll n$. Then each edge probability P_{ij} is said to depend on the labels z_i and z_j , possibly along with some other parameters. For example, the stochastic block model (SBM) sets a fixed edge probability for each pair of communities, i.e., $P_{ij} = \omega_{z_i, z_j}$. The degree-corrected block model (DCBM) assigns an additional parameter θ_i to each vertex by which edge probabilities are scaled, i.e., $P_{ij} = \theta_i \theta_j \omega_{z_i, z_j}$. The popularity adjusted block model (PABM) assigns K parameters to each vertex $\lambda_{i1}, \lambda_{i2}, \dots, \lambda_{iK}$ that describe that vertex's affinity toward each community; the edge probability between vertices i and j is then defined as the product of vertex i 's affinity toward vertex j 's community and vertex j 's affinity toward vertex i 's community, i.e., $P_{ij} = \lambda_{iz_j} \lambda_{jz_i}$.

The three block model types, as well as the Erdős-Rényi model, impose structure on P , including on the rank of P . P has rank 1 for the Erdős-Rényi model, rank K for the SBM and DCBM, and rank K^2 for the PABM. This provides the intuition behind another family of Bernoulli graphs called the *random dot product graph* (RDPG) and *generalized random dot product graph* (GRDPG). In the RDPG, each vertex has a corresponding latent vector in d -dimensional Euclidean space, where d is the rank of P and P is positive semidefinite. Then the edge probability between each pair of vertices is defined as the inner product between the corresponding latent vectors, i.e., $P_{ij} = x_i^\top x_j$. If the latent vectors are collected in a data matrix $X = [x_1 \mid \dots \mid x_n]^\top$, then the edge probability matrix for the RDPG is

$P = XX^\top$. Similarly, the edge probability between each pair of vertices for the GRDPG is defined as the indefinite inner product between the corresponding latent vectors, i.e., $P_{ij} = x_i^\top I_{p,q} x_j$, where $I_{p,q} = \text{blockdiag}(I_p, -I_q)$ and $p + q = d$. Then the edge probability matrix for the GRDPG is $P = XI_{p,q}X^\top$. This allows for a model similar to the RDPG for non-positive semidefinite P . While the RDPG and GRDPG do not necessarily have community structure, it has been shown that block models are specific cases of the RDPG or GRDPG in which latent vectors are organized by community. This includes the SBM, in which communities correspond to point masses, DCBM, in which communities correspond to line segments, and PABM, in which communities correspond to orthogonal subspaces. In this work, we extend this idea to communities organized into more general latent structures. In particular, we assume that each community corresponds to a manifold in the latent space.

2 Latent Structure Block Models

All block models are latent structure models.

Definition 1 (Manifold block model). Let $p, q \geq 0$, $d = p + q \geq 1$, $K \geq 1$, and $n \geq 1$ be integers. If there are manifolds $\mathcal{M}_1, \dots, \mathcal{M}_K \subset \mathcal{X}$ for $\mathcal{X} = \{x, y \in \mathbb{R}^d : x^\top I_{p,q} y \in [0, 1]\}$, and $X = \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}^\top$ such that each $x_i \in \mathcal{M}_k$ for some $k \in [K]$, then $A \sim \text{GRDPG}_{p,q}(X; \rho_n)$ is a *manifold block model*.

In practice, we often define each manifold \mathcal{M}_k by a continuous function $g_k : [0, 1]^r \rightarrow \mathcal{X}$, where $1 \leq r < d$ is the dimensionality of manifold \mathcal{M}_k . To complete the probability model, we define a probability distribution F with support $[0, 1]^r$ and draw the community memberships from a multinomial distribution.

The full mixture model can be described as follows:

1. Draw labels $z_1, \dots, z_n \stackrel{\text{iid}}{\sim} \text{Multinomial}(\alpha_1, \dots, \alpha_K)$.
2. Draw latent vectors by first drawing each $t_1, \dots, t_n \stackrel{\text{iid}}{\sim} F$ and then computing each $x_i = g_{z_i}(t_i)$.
3. Let $X = \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}^\top$, and draw $A \sim \text{RDPG}(X)$ or $A \sim \text{GRDPG}_{p,q}(X)$.

Example 1. Let $p = 2$, $q = 0$, $K = 2$, and $r = 1$. Define two one-dimensional manifolds by $f_1(t) = \begin{bmatrix} \cos(\frac{\pi}{3}t) & \sin(\frac{\pi}{3}t) \end{bmatrix}^\top$ and $f_2(t) = \begin{bmatrix} 1 - \cos(\frac{\pi}{3}t) & 1 - \sin(\frac{\pi}{3}t) \end{bmatrix}^\top$. Draw $t_1, \dots, t_n \stackrel{\text{iid}}{\sim}$

Uniform(0, 1) and $z_1, \dots, z_n \stackrel{\text{iid}}{\sim} \text{Multinomial}(\frac{1}{2}, \frac{1}{2})$, and compute latent vectors $x_i = f_{z_i}(t_i)$, which are collected in data matrix $X = \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}^\top$. Finally, let $A \sim \text{RDPG}(X)$. Fig 1 shows the latent configuration drawn from this latent distribution, a random dot product graph drawn from the latent configuration, and the adjacency spectral embedding of the graph.

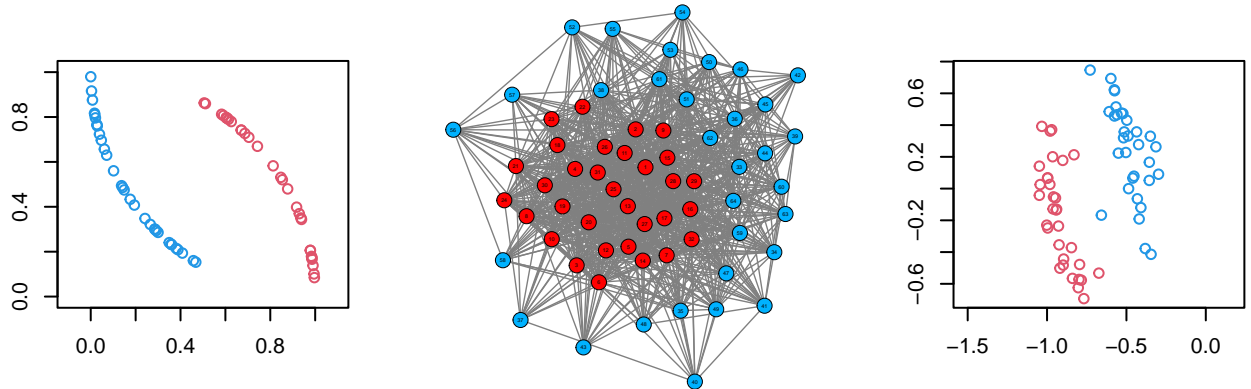


Figure 1: Manifold block model described in Example 1. The latent configuration is on the left, and a random dot product graph drawn from the latent configuration is on the middle, and the ASE is on the right.

3 Methods

Here, we provide two algorithms for MBM community detection. First, we consider the case in which communities correspond to manifolds in the latent space that do not intersect and are separated by some finite distance. In this scenario, we use the convergence of the ASE to show that single linkage clustering on the latent space produces a clustering such that the total number of misclustered vertices goes to zero, with high probability.

Next, we consider the case in which communities correspond to one-dimensional manifolds in the latent space and may or may not intersect. In this scenario, we propose an alternating coordinate descent algorithm that alternates between estimating the structure of the manifolds and the community labels, which we call K -curves clustering. We again use the convergence of the ASE to show that under certain conditions, K -curves clustering produces a clustering such that the proportion of misclustered vertices goes to zero, with

high probability.

3.1 Nonintersecting Manifolds

In this section, we consider the following scenario: Suppose that each community is represented by a manifold \mathcal{M}_k , $k \in \{1, \dots, K\}$ in the latent space of a RDPG or GDRPG, and the manifolds do not intersect each other. Define $\delta = \min_{k \neq \ell} \min_{x \in \mathcal{M}_k, y \in \mathcal{M}_\ell} \|x - y\|$, the minimum distance between two manifolds. We assume that $\delta > 0$, i.e., the manifolds do not intersect.

In the noiseless setting, if the subsample on each manifold is sufficiently dense, it is possible to construct for each manifold an η_k -neighborhood graph for each manifold for some $\eta_k > 0$ such that the graph is connected. Then if $\max_k \eta_k = \eta < \delta$, an η -neighborhood graph for the entire sample will consist of K disconnected subgraphs that map onto each manifold. Equivalently, we can apply single-linkage clustering. The remainder of this section explores under which conditions these criteria are met for the latent configuration, in which latent vectors lie exactly on manifolds, as well as the ASE, which introduces noise.

Algorithm 1: ASE clustering for nonintersecting communities.

Data: Adjacency matrix A , number of communities K , embedding dimensions p and q .

Result: Community assignments $z_1, \dots, z_n \in \{1, \dots, K\}$.

- 1 Compute \hat{X} , the ASE of A using the p most positive and q most negative eigenvalues and their corresponding eigenvectors.
 - 2 Apply single linkage clustering with K communities on \hat{X}
-

Consider the specific case in which each $\mathcal{M}_k \subset \mathbb{R}^d$ is a one-dimensional manifold. Let F_k be a probability distribution with support \mathcal{M}_k . Then we define a mixture model as follows:

1. Draw labels $z_1, \dots, z_n \stackrel{\text{iid}}{\sim} \text{Multinomial}(\alpha_1, \dots, \alpha_K)$.
2. Draw latent vectors each as $x_i \stackrel{\text{ind}}{\sim} F_{z_i}$ for distributions F_1, \dots, F_K with respective supports $\mathcal{M}_1, \dots, \mathcal{M}_K$.
3. Let $X = \begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix}^\top$, and draw $A \sim \text{RDPG}(X)$ or $A \sim \text{GRDPG}_{p,q}(X)$.

Note that here, we redefine the model to ignore g_1, \dots, g_K , the parameterizations of each manifold. Instead, we sample points directly on the manifolds themselves. We will return to the parameterizations in Section 3.2.

In this setting, although each manifold is in d -dimensional space, since the manifolds are one-dimensional, it is possible to consider an ordering of points sampled on the manifold. In the original latent space in which there is no noise, the condition for which community detection is possible via η -neighborhood graph clustering or single linkage clustering is each pair of adjacent order statistics on each manifold being less than δ apart.

Theorem 1 (Community detection for one-dimensional nonintersecting manifolds without noise). *Let x_1, \dots, x_n be points sampled on K one-dimensional manifolds $\mathcal{M}_1, \dots, \mathcal{M}_K$. Suppose $\delta = \min_{k \neq \ell} \min_{x_i \in \mathcal{M}_k, x_j \in \mathcal{M}_\ell} \|x_i - x_j\| > 0$. Let A_n be the event such that $\max_k \max_i \|x_{(i+1)}^{(k)} - x_{(i)}^{(k)}\| < \delta$, where $x_{(i)}^{(k)}$ is the i^{th} order statistic on the k^{th} manifold. Then for any $\epsilon > 0$, there exists an $N \in \mathbb{N}$ such that when $n > N$, $P(A_n) > 1 - \epsilon$.*

Sampling a graph from this latent configuration and constructing an ASE from the graph introduces error relative to the original latent vectors. If the errors are small, which is guaranteed by the consistency property of the ASE for sufficiently large n , then sufficient separation between embedding points belonging to different manifolds is still guaranteed.

Theorem 2 (Community detection for one-dimensional nonintersecting manifolds with noise). *Suppose the same setup as in theorem 1. Define $y_i = x_i + e_i$ and B_n as the event such that $\|y_{(i+1)}^{(k)} - y_{(i)}^{(k)}\| < \|y_{(i)}^{(k)} - y_{(j)}^{(\ell)}\|$ for any $k \neq \ell$. Here, $y_{(i)}^{(k)}$ is the i^{th} order statistic on the k^{th} manifold plus its corresponding noise. Then if $\max_i \|e_i\| < \delta/3$, for any $\epsilon > 0$, there exists an $N \in \mathbb{N}$ such that when $n > N$, $P(B_n) > 1 - \epsilon$.*

Theorem 3 (Community detection for one-dimensional nonintersecting manifold block models). *Suppose we have the same setup as in theorem 1. Let $p, q \geq 0$ such that $p + q = d$ where d is the dimension of the range of each g_k . Suppose we draw $A \sim \text{GRDPG}_{p,q}(X; \rho_n)$ for some sparsity parameter $\rho_n \in (0, 1]$, and let \hat{X} be the ASE of A . Then single linkage clustering produces zero community detection error with high probability as $n \rightarrow \infty$.*

3.2 Intersecting Manifolds

In this section, we again consider the setting for the RDPG or GRDPG in which each community lies on a manifold in the latent space. However, this time, we do not assume that the manifolds are nonintersecting. We also restrict this case to one-dimensional manifolds which are each described by $g_k : [0, 1] \rightarrow \mathcal{X}$. Then we define a mixture model as follows:

1. Draw $t_1, \dots, t_n \stackrel{\text{iid}}{\sim} F$ for probability distribution F with support $[0, 1]$.
2. Draw $z_1, \dots, z_n \stackrel{\text{iid}}{\sim} \text{Multinomial}(\alpha_1, \dots, \alpha_K)$, the community labels.
3. Let each $x_i = g_{z_i}(t_i)$ be the latent vector for vertex v_i , and collect the latent vectors into matrix $X = \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}^\top$.
4. Draw $A \sim \text{RDPG}(X)$ or $A \sim \text{GRDPG}_{p,q}(X)$.

Algorithm 2: K -curves clustering.

Data: Adjacency matrix A , number of communities K , embedding dimensions p , q , stopping criterion ϵ

Result: Community assignments $1, \dots, K$, curves g_1, \dots, g_K

- 1 Compute X , the ASE of A using the p most positive and q most negative eigenvalues and their corresponding eigenvectors.
 - 2 Initialize community labels z_1, \dots, z_n .
 - 3 **repeat**
 - 4 **for** $k = 1, \dots, K$ **do**
 - 5 Define X_k as the rows of X for which $z_i = k$.
 - 6 Fit curve g_k and positions t_{k_i} to X_k by minimizing $\sum_{k_i} \|x_{k_i} - g_k(t_{k_i})\|^2$.
 - 7 **end**
 - 8 **for** $k = 1, \dots, K$ **do**
 - 9 Assign $z_i \leftarrow \arg \min_\ell \|x_i - g_\ell(t_i)\|^2$.
 - 10 **end**
 - 11 **until** the change in $\sum_k \sum_{i \in C_k} \|x_i - g_k(t_i)\|^2$ is less than ϵ
-

Theorem 4. Let each g_k be smooth. Then K -curves clustering converges to a stationary point of the objective, $\sum_k \sum_{i \in C_k} \|x_i - g_k(t_i)\|^2$.

Proof. K -curves clustering is a batch coordinate descent algorithm. Thus, in order to show that it converges to a stationary point, it is sufficient to show that each descent step decreases the objective function. \square

K -curves clustering assumes that the functional form of g_k is known. The choice of g_k affects the difficulty of the algorithm. As a balance between flexibility and ease of estimation, we consider the case where each g_k is a Bezier polynomial of degree R with coefficients p_k . Then we have $g_k(t) = g(t; p_k) = \sum_{r=0}^R p_k^{(r)} \binom{R}{r} (1-t)^{R-r} t^r$.

Given $\{t_i\}$ and $\{z_i\}$, it is straightforward to obtain $\hat{p}_k = \arg \min_p \sum_{k_i} \|x_{k_i} - g_k(t_{k_i}; p)\|^2$

$$\hat{p}_k = (T_k^\top T_k)^{-1} T_k^\top X_k,$$

where T_k is an $n_k \times (R+1)$ matrix with rows $\left[(1-t_{k_i})^R \quad (1-t_{k_i})^{R-1} t_{k_i} \quad \dots \quad (1-t_{k_i}) t_{k_i}^{R-1} \quad t_{k_i}^R \right]$. Estimation of $\{t_i\}$ given $\{z_i\}$ and $\{p_k\}$ is more difficult. Each t_i can be estimated separately:

$$\hat{t}_i = \arg \min_t \|x_i - g(t; p_{z_i})\|^2. \quad (1)$$

This is equivalent to solving $0 = (x_i - g(t; p_{z_i}))^\top (\dot{g}(t; p_{z_i}))$. Setting $c^{(s)} = \sum_{r=0}^s (-1)^{s-r} \binom{R}{r} p_{z_i}^{(r)}$ for $s \neq 0$ and $c^{(0)} = p_{z_i}^{(0)} - x_i$, let $c = \begin{bmatrix} c^{(0)} & \dots & c^{(R)} \end{bmatrix}^\top$. Then solving Eq. 1 is equivalent to finding the real roots of a polynomial with coefficients that are the sums of the reverse diagonals of CD^\top , where $C_{ij} = c_{ij}(-1)^i \binom{R}{i}$ and $D_{ij} = c_{i-1,j}(-1)^{i-1} \binom{R-1}{i-1}$.

Theorem 5. *Let each $g(\cdot; p_k)$ be a nonintersecting Bezier polynomial of order R , and a GRDPG is drawn from vectors that lie on the curves. Suppose we observe the true labels of m_k vertices from each community, and each $m_k > R + 1$. Suppose further that latent vectors $x_j = g(t_i; p_{z_j})$ that correspond to vertices with observed labels are such that Then as $n \rightarrow \infty$, the proportion of misclustered vertices from K -curves clustering approaches 0 with probability 1.*

4 Examples

Example 2. Here, $K = 2$ with $g_1(t) = \begin{bmatrix} t^2 & 2t(1-t) \end{bmatrix}^\top$ and $g_2(t) = \begin{bmatrix} 2t(1-t) & (1-t)^2 \end{bmatrix}^\top$. We draw $n_1 = n_2 = 2^8$ points uniformly from each curve.

Algorithm 3: Semi-supervised K -curves clustering.

Data: Adjacency matrix A , number of communities K , embedding dimensions p , q , stopping criterion ϵ , $m_k \leq n_k$ known community assignments for each community

Result: Community assignments $1, \dots, K$, curves g_1, \dots, g_K

- 1 Compute X , the ASE of A using the p most positive and q most negative eigenvalues and their corresponding eigenvectors.
 - 2 Fit curves g_1, \dots, g_K using each of the m_1, \dots, m_K points with known community labels by minimizing $\sum_{j=1}^{m_i} \|x_j - g_k(t_j)\|^2$.
 - 3 Assign labels z_1, \dots, z_n to each x_1, \dots, x_n by minimizing $\|x_i - g_k(t_i)\|^2$ for k , holding the initial known labels constant.
 - 4 **repeat**
 - 5 **for** $k = 1, \dots, K$ **do**
 - 6 Define X_k as the rows of X for which $z_i = k$.
 - 7 Fit curve g_k and positions t_{k_i} to X_k by minimizing $\sum_{k_i} \|x_{k_i} - g_k(t_{k_i})\|^2$.
 - 8 **end**
 - 9 **for** $k = 1, \dots, K$ **do**
 - 10 Assign $z_i \leftarrow \arg \min_{\ell} \|x_i - g_{\ell}(t_i)\|^2$, holding the known initial labels constant.
 - 11 **end**
 - 12 **until** the change in $\sum_k \sum_{i \in C_k} \|x_i - g_k(t_i)\|^2$ is less than ϵ
-

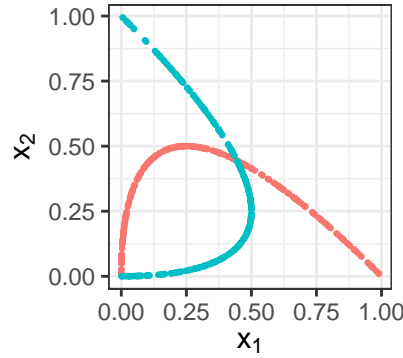


Figure 2: Latent positions, labeled by curve/community.

We draw $A \sim \text{RDPG}(X)$ and obtain the following ASE:

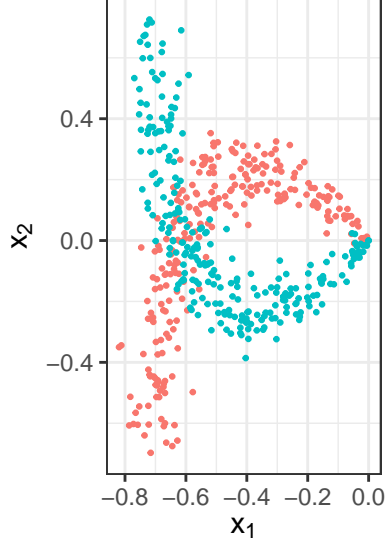


Figure 3: ASE of an RDPG drawn from the latent positions, labeled by curve/community.

We then try applying K -curves clustering to this graph. The first three are with random initial labels, forcing the intercept to be zero. The fourth initializes the labels randomly but allows the intercept to be nonzero. The fifth initializes the labels by spectral clustering with the normalized Laplacian, again forcing the intercept to be zero. The sixth also initializes via spectral clustering but allows the intercept to be nonzero.

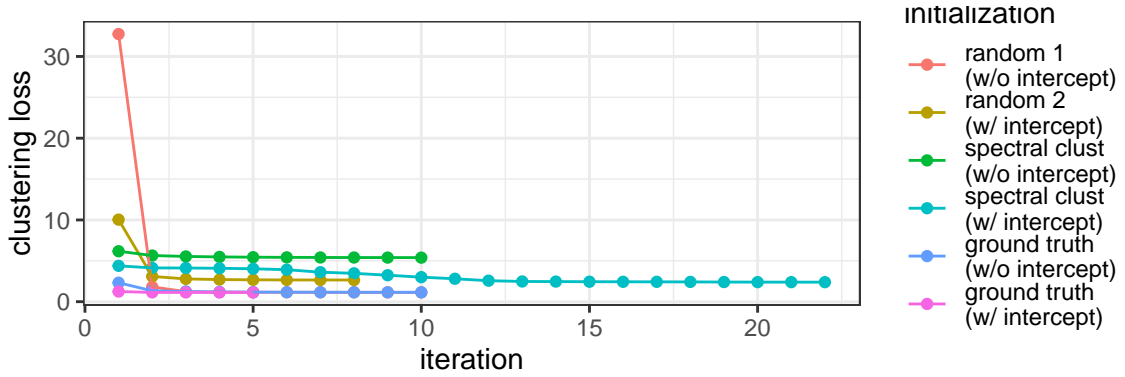


Figure 4: Clustering loss vs. iteration for each run of K -curve clustering.

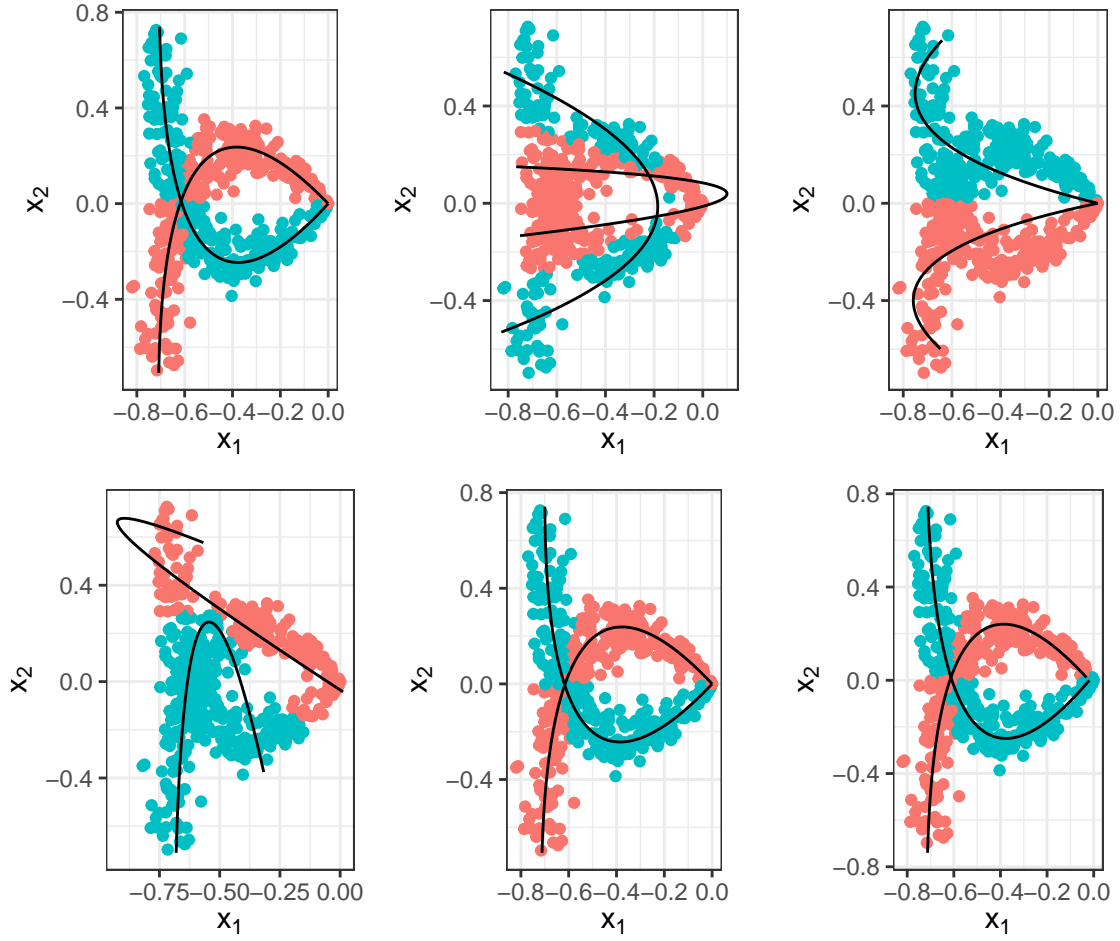
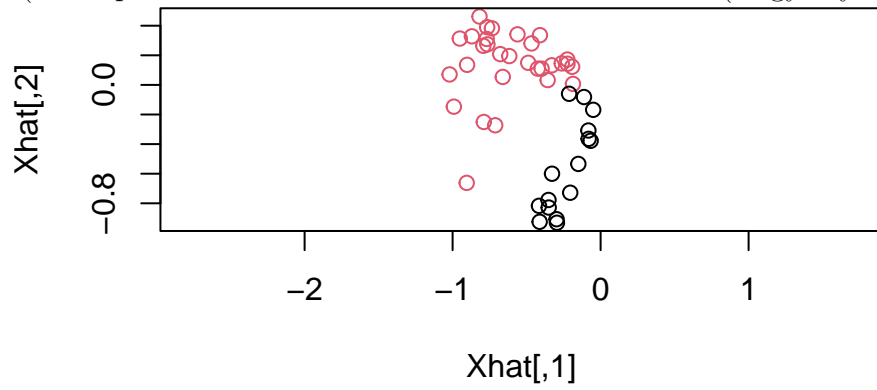
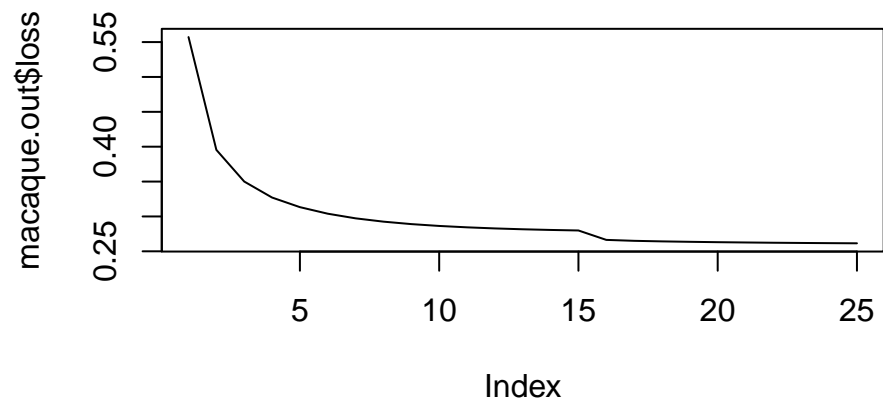
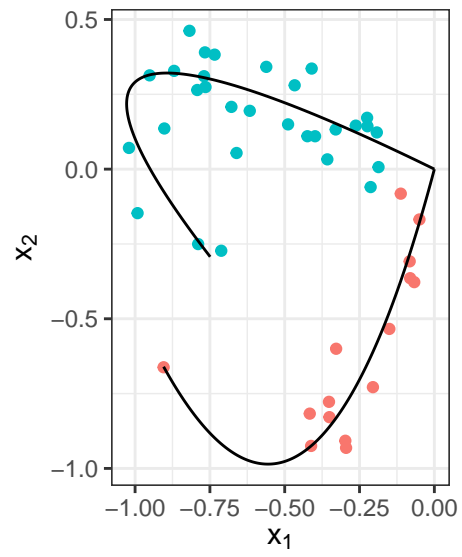


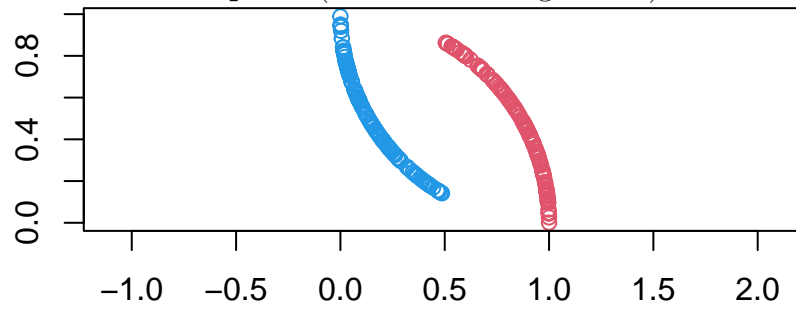
Figure 5: ASE labeled by estimated community labels for each initialization strategy.

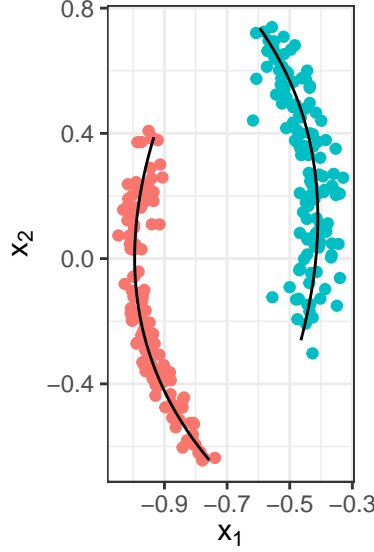
Example 3 (Macaque visuotactile brain areas and connections (Négyessy et al. 2006)).





Example 4 (Non-intersecting curves).





5 Simulation Study

6 Discussion

A Proofs of Theorems

To prove theorems 1, 2, and 3, we first need to introduce some theory for order statistics. Let $D_i^{(k)} = X_{(i+1)}^{(k)} - X_{(i)}^{(k)}$, where $X_{(i)}^{(k)}$ is the i^{th} order statistic of univariate sample $X_1^{(k)}, \dots, X_n^{(k)} \stackrel{\text{iid}}{\sim} F^{(k)}$ and $F^{(k)}$ is a continuous distribution on manifold \mathcal{M}_k . If $\max_i D_i^{(k)} < \delta$, then there is sufficient separation of points between each manifold. Then it is sufficient to quantify $P(\max_i D_i^{(k)} > \delta)$ for each k as a function of n and δ and show that this converges to zero as n grows to ∞ .

For ease of notation, we drop the superscript denoting community/manifold membership until lemma 6. Denote $f(x)$ as the density of each F , $g_i(x)$ as the density of $X_{(i)}$, $g_{ij}(x, y)$ as the joint density of $X_{(i)}, X_{(j)}$, and $h_i(d)$ as the density of D_i (with corresponding capital letters for the cumulative distribution functions).

The following are taken as given¹:

1. $g_i(x) = \frac{n!}{(n-i)!(i-1)!} (F(x))^{i-1} (1 - F(x))^{n-i} f(x)$.

¹https://en.wikipedia.org/wiki/Order_statistic

2. $g_{ij}(x, y) = \frac{n!}{(i-1)!(j-i-1)!(n-j)!} (F(x))^{i-1} (F(y) - F(x))^{j-i-1} (1 - F(y))^{n-j} f(x) f(y)$.
3. By convolution, $h_i(d) = \int_0^1 g_{i,i+1}(x, x+d) dx$.

Lemma 1 (The probability density function of D_i). *Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} F$ and the support of F is the unit interval. Let $D_i = X_{(i+1)} - X_{(i)}$. Then the density of each D_i is:*

$$h_i(d) = \int_0^{1-d} \frac{n!}{(i-1)!(n-i-1)!} (F(x))^{i-1} (1 - F(x+d))^{n-i-1} f(x) f(x+d) dx \quad (2)$$

Proof. This is a direct consequence of 2 and 3. We also note that because the support of X_i is $[0, 1]$, the integral only needs to be evaluated from 0 to $1-d$ because of the $f(x+d)$ and $1 - F(x+d)$ terms. \square

Lemma 2 (The cumulative distribution function of D_i). *Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} F$ and the support of F is the unit interval. Let $D_i = X_{(i+1)} - X_{(i)}$. Then the distribution function of each D_i is:*

$$P(D_i < \delta) = H_i(\delta) = 1 - \int_0^{1-\delta} \frac{n!}{(n-i)!(i-1)!} (F(x))^{i-1} (1 - F(x+\delta))^{n-i} f(x) dx \quad (3)$$

Proof.

$$\begin{aligned} H_i(\delta) &= \int_x^{x+\delta} h_i(d) dd \\ &= \int_x^{x+\delta} \int_0^1 \frac{n!}{(i-1)!(n-i-1)!} ((F(x))^{i-1} (1 - F(x+d))^{n-i-1} f(x) f(x+d)) dx dd \\ &= \int_0^1 \frac{n!}{(i-1)!(n-i-1)!} (F(x))^{i-1} f(x) \int_x^{x+\delta} (1 - F(x+d))^{n-i-1} f(x+d) dd dx \\ &= \int_0^1 \frac{n!}{(i-1)!(n-i-1)!} (F(x))^{i-1} f(x) \int_{F(x)}^{F(x+\delta)} (1-u)^{n-i-1} du dx \\ &= \int_0^1 \frac{n!}{(i-1)!(n-i)!} (F(x))^{i-1} f(x) ((1 - F(x))^{n-i} - (1 - F(x+\delta))^{n-i}) dx \\ &= \int_0^1 g_i(x) dx - \int_0^1 \frac{n!}{(i-1)!(n-i)!} (F(x))^{i-1} (1 - F(x+\delta))^{n-i} f(x) dx \\ &= 1 - \int_0^1 \frac{n!}{(i-1)!(n-i)!} (F(x))^{i-1} (1 - F(x+\delta))^{n-i} f(x) dx \end{aligned}$$

Because of the $x + \delta$ term and the fact that the support is $[0, 1]$, the integrand is zero above $1 - \delta$, so we are left with

$$= 1 - \int_0^{1-\delta} \frac{n!}{(i-1)!(n-i)!} (F(x))^{i-1} (1 - F(x+\delta))^{n-i} f(x) dx.$$

□

We now focus on the case where points are sampled uniformly on the unit interval to build up to more general distributions.

Lemma 3 (Differences between order statistics of a uniform distribution). *If $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Uniform}(0, 1)$, then each $D_i \sim \text{Beta}(1, n)$.*

Proof. We begin with Eq. (2), plugging in $f(x) = 1$ and $F(x) = x$:

$$h_i(d) = \int_0^{1-d} \frac{n!}{(i-1)!(n-i-1)!} x^{i-1} (1-x-d)^{n-i-1} dx.$$

Then we proceed with integration by parts, setting $u = x^{i-1} \implies du = (i-1)x^{i-2}$ and $dv = (1-x-d)^{n-i-1} dx \implies v = -\frac{1}{n-i}(1-x-d)^{n-i}$. Note that $uv|_0^{1-d} = 0$ in this case. This yields

$$= \frac{n!}{(i-1)!(n-i-1)!} \int \frac{i-1}{n-i} x^{i-2} (1-x-d)^{n-i} dx.$$

Then applying integration by parts again until the x^p term disappears, we get:

$$\begin{aligned} &= \frac{n!}{(i-1)!(n-i-1)!} \frac{(i-1)!}{(n-i) \cdots (n-2)} \int_0^{1-d} (1-x-d)^{n-2} dx \\ &= -\frac{n(n-1)}{n-1} (1-x-d)^{n-1} \Big|_0^{1-d} \\ &= n(1-d)^{n-1}. \end{aligned}$$

This is the density function for $\text{Beta}(1, n)$, completing the proof. □

Lemma 4. *Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Uniform}(0, 1)$. Then for any ϵ and $\delta > 0$, there exists an $N = O\left(\frac{-\log \epsilon}{\delta}\right)$ such that $P(\max_i X_{(i+1)} - X_{(i)} < \delta) \geq 1 - \epsilon$ when $n > N$.*

Proof. Since $X_{(i+1)} - X_{(i)} = D_i \sim \text{Beta}(1, n)$, $P(X_{(i+1)} - X_{(i)} < \delta) = 1 - (1 - \delta)^n$. This yields

$$\begin{aligned}
P(\max_i D_i < \delta) &\geq (P(D_i < \delta))^{n-1} \\
&= (1 - (1 - \delta)^n)^{n-1} \\
&\approx e^{-n \exp(-n\delta)}.
\end{aligned}$$

In the limit $n \rightarrow \infty$, this goes to 1. □

Now we extend this lemma to general distributions on the unit interval.

Lemma 5. *Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} F$ with support $[0, 1]$, and suppose $f(x)$ is continuous and $f(x) \geq a > 0$ everywhere on the support. Let $D_i = X_{(i+1)} - X_{(i)}$. Then for any $\epsilon > 0$, there exists $N > 0$ such that $P(\max_i D_i < \delta) \geq 1 - \epsilon$ when $n > N$.*

Proof. We start with Eq. (3):

$$P(D_i \leq \delta) = 1 - \int_0^{1-\delta} \frac{n!}{(n-i)!(i-1)!} (F(x))^{i-1} (1 - F(x + \delta))^{n-i} f(x) dx.$$

Making the approximation $F(x + \delta) \approx F(x) + \delta f(x)$ and bounding $f(x) \geq a$, we get:

$$P(D_i \leq \delta) \geq 1 - \int_0^{1-\delta} \frac{n!}{(n-i)!(i-1)!} (F(x))^{i-1} (1 - F(x) - a\delta)^{n-i} f(x) dx.$$

Then making the substitution $u = F(x) \implies du = f(x) dx$, we obtain

$$1 - \int_0^{F(1-\delta)} \frac{n!}{(n-i)!(i-1)!} u^{i-1} (1 - u - a\delta)^{n-i} du$$

Evaluating the integral yields

$$P(D_i < \delta) = 1 - (1 - a\delta)^n + (1 - F(1 - \delta) - a\delta)^n.$$

Then as before,

$$\begin{aligned}
P(\max_i D_i < \delta) &= P(\text{all } D_i < \delta) \\
&= 1 - P(\text{some } D_i > \delta) \\
&\geq 1 - \sum_i^{n-1} P(D_i > \delta) \\
&= 1 - (n-1)(1 - a\delta)^n + (n-1)(1 - F(1 - \delta) - a\delta)^n
\end{aligned}$$

Since $(1 - a\delta) < 1$ and $(1 - F(1 - \delta) - a\delta) < 1$, this converges to 1 in the limit $n \rightarrow \infty$. Then for any $\epsilon \in (0, 1)$, we can always solve for an n such that $1 - (n - 1)(1 - a\delta)^n(n - 1)(1 - F(1 - \delta) - a\delta)^n \geq 1 - \epsilon$ holds. \square

Proof of theorem 1. This theorem requires two additional conditions: First, each $n_k \rightarrow \infty$ as $n \rightarrow \infty$, and second, the number of manifolds K remains fixed. These are satisfied by the model assumptions. Then this is a direct consequence of Lemma 5. Denoting $m = \min_k n_k$ as the sample size of the manifold with the fewest points,

$$\begin{aligned} P(\max_{k,i} \|X_{(i+1)}^{(k)} - X_{(i)}^{(k)}\| < \delta) &= P(\forall k, i \ D_i^{(k)} < \delta) \\ &\geq \prod_{k=1}^K P(D_i^{(k)} < \delta) \\ &\geq (1 - (m - 1)(1 - a\delta)^m + (m - 1)(1 - F(1 - \delta) - a\delta)^m)^K \end{aligned}$$

Since as $n \rightarrow \infty$, $m \rightarrow \infty$ as well, and K is fixed, this approaches $(1)^K = 1$. Then setting this $\geq 1 - \epsilon$, it is always possible to solve this for m . \square

Proof of theorem 2. Let $\max_i \|e_i\| \leq \xi$. \square

Proof of theorem 3. \square

Here, we extend the results on the line to the unit hypercube. The following lemma is based on lemma 2 of Trosset & Buyukbas (2020).

Lemma 6. Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} F$ with support $[0, 1]^r$, and $f(x) \geq a > 0$ everywhere on the support. Define E_n as the event that an η -neighborhood graph constructed from the sample is connected. Then for any $\epsilon > 0$, there exists $N = O\left(\frac{\log \epsilon \eta^r / r^{r/2}}{\log(1 - \frac{a\eta^r}{r^{r/2}})}\right)$ such that $P(E_n) > 1 - \epsilon$ when $n \geq N$.

Proof (sketch). Divide the hypercube $[0, 1]^r$ into a grid of sub-hypercubes of side length at most η/\sqrt{r} . E_n is satisfied if each sub-hypercube contains at least one X_i from the sample.

$$\begin{aligned}
P(E_n) &= 1 - P(\text{some cells don't contain } X_i) \\
&\geq 1 - \sum_k^{\lceil \sqrt{r}/\eta \rceil^r} \prod_i^n P(X_i \text{ is not in the } k^{\text{th}} \text{ hypercube}) \\
&\geq 1 - \lceil \sqrt{r}/\eta \rceil^r (1 - a\eta^r/r^{r/2})^n
\end{aligned}$$

Setting this $\geq 1 - \epsilon$ and solving for n yields the desired rate. \square

B Details on Fitting Bezier Curves with Noise

References

Négyessy, L., Nepusz, T., Kocsis, L. & Bazsó, F. (2006), ‘Prediction of the main cortical areas and connections involved in the tactile function of the visual cortex by network analysis’, *European Journal of Neuroscience* **23**(7), 1919–1930.

URL: <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1460-9568.2006.04678.x>

Trosset, M. W. & Buyukbas, G. (2020), ‘Rehabilitating isomap: Euclidean representation of geodesic structure’.

URL: <https://arxiv.org/abs/2006.10858>