

---

# Selecting the independent coordinates of manifolds with large aspect ratios

---

**Yu-Chia Chen**

Department of Electrical & Computer Engineering  
University of Washington  
Seattle, WA 98195  
yuchaz@uw.edu

**Marina Meilă**

Department of Statistics  
University of Washington  
Seattle, WA 98195  
mmp2@uw.edu

## Abstract

Many manifold embedding algorithms fail apparently when the data manifold has a large aspect ratio (such as a long, thin strip). Here, we formulate success and failure in terms of finding a smooth embedding, showing also that the problem is pervasive and more complex than previously recognized. Mathematically, success is possible under very broad conditions, provided that embedding is done by carefully selected eigenfunctions of the Laplace-Beltrami operator  $\Delta$ . Hence, we propose a bicriterial *Independent Eigencoordinate Selection (IES)* algorithm that selects smooth embeddings with few eigenvectors. The algorithm is grounded in theory, has low computational overhead, and is successful on synthetic and large real data.

## 1 Motivation

We study a well-documented deficiency of manifold learning algorithms. Namely, as shown in [GZKR08], algorithms that find their output  $\mathbf{Y} = \varphi(\mathbf{X})$  by minimizing a quadratic form under some normalization constraints, fail spectacularly when the data manifold has a large aspect ratio, that is, it extends much more in one direction than in others, as the long, thin strip illustrated in Figure 1. This class, called *output normalized (ON)* algorithms, includes Locally Linear Embedding (LLE), Laplacian Eigenmap, Local Tangent Space Alignment (LTSA), Hessian Eigenmaps (HLE), and Diffusion maps. The problem, often observed in practice, was formalized in [GZKR08] as failure to find an  $\mathbf{Y}$  affinely equivalent to  $\mathbf{X}$ . They give sufficient conditions for failure, using a linear algebraic perspective. The conditions show that, especially when noise is present, the problem is pervasive.

In the present paper, we revisit the problem from a *differential geometric* perspective. First, we define failure not as distortion, but as drop in the *rank* of the mapping  $\phi$  represented by the embedding algorithm. In other words, the algorithm fails when the map  $\phi$  is not invertible, or, equivalently, when the dimension  $\dim \phi(\mathcal{M}) < \dim \mathcal{M} = d$ , where  $\mathcal{M}$  represents the idealized data manifold, and  $\dim$  denotes the intrinsic dimension. Figure 1 demonstrates that the problem is fixed by choosing the eigenvectors with care. In fact, as we show in Section 6, it is known [Bat14] that for the DM and LE algorithms, under mild geometric conditions, one can *always* find a finite set of  $m$  eigenfunctions that provide a smooth  $d$ -dimensional map. We call this problem the *Independent Eigencoordinate Selection (IES)* problem, formulate it and explain its challenges in Section 3.

Our second main contribution (Section 4) is to design a bicriterial method that will select from a set of *coordinate functions*  $\phi_1, \dots, \phi_m$ , a subset  $S$  of small size that provides a smooth full-dimensional embedding of the data. The IES problem requires searching over a combinatorial number of sets. We show (Section 4) how to drastically reduce the computational burden per set for our algorithm. Third, we analyze the proposed criterion under asymptotic limit (Section 6). Finally (Section 7), we show

examples of successful selection on real and synthetic data. The experiments also demonstrate that users of manifold learning for other than toy data *must* be aware of the IES problem and have tools for handling it. Notations table, proofs, a library of hard examples, extra experiments and analyses are in Supplements A–G; Figure/Table/Equation references with prefix S are in the Supplement.

## 2 Background on manifold learning

**Manifold learning (ML) and intrinsic geometry** Suppose we observe data  $\mathbf{X} \in \mathbb{R}^{n \times D}$ , with data points denoted by  $\mathbf{x}_i \in \mathbb{R}^D \forall i \in [n]$ , that are sampled from a *smooth*<sup>1</sup>  $d$ -dimensional submanifold  $\mathcal{M} \subset \mathbb{R}^D$ . Manifold Learning algorithms map  $\mathbf{x}_i, i \in [n]$  to  $\mathbf{y}_i = \phi(\mathbf{x}_i) \in \mathbb{R}^s$ , where  $d \leq s \ll D$ , thus reducing the dimension of the data  $\mathbf{X}$  while preserving (some of) its properties. Here we present the LE/DM algorithm, but our results can be applied to other ML methods with slight modification.

The first two steps of LE/DM [CL06, NLCK06] algorithms are generic; they are performed by most ML algorithms. First we encode the neighborhood relations in a *neighborhood graph*, which is an undirected graph  $G(V, E)$  with vertex set  $V$  be the collection of all points  $i \in [n]$  and edge set  $E$  be the collections of tuple  $(i, j) \in V^2$  such that  $i$  is  $j$ 's neighbor (and vice versa). Common methods for building neighborhood graphs include  $r$ -radius graph and  $k$ -nearest neighbor graph. Readers are encouraged to refer to [HAvL07, THJ10] for details. In this paper,  $r$ -radius graph are considered, for which principled methods to select the neighborhood size and dimension exist. For such method, the edge set  $E$  of the neighborhood graph is constructed as follow:  $E = \{(i, j) \in V^2 : \|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq r\}$ . Closely related to the neighborhood graph is the *kernel matrix*  $\mathbf{K} \in \mathbb{R}^{n \times n}$  whose elements are  $K_{ij} = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|}{\varepsilon}\right)$  if  $(i, j) \in E$  and 0 otherwise. Typically, the radius  $r$  and the *bandwidth* parameter  $\varepsilon$  are related by  $r = c\varepsilon$  with  $c$  a small constant greater than 1, e.g.,  $c \in [3, 10]$ . This ensures that  $\mathbf{K}$  is close to its limit when  $r \rightarrow \infty$  while remaining sparse, with sparsity structure induced by the neighborhood graph. Having obtained the kernel matrix  $\mathbf{K}$ , we then construct the *renormalized graph Laplacian* matrix  $\mathbf{L}$  [CL06], also called the *sample Laplacian*, or *Diffusion Maps Laplacian*, by the following:  $\mathbf{L} = \mathbf{I}_n - \tilde{\mathbf{W}}^{-1} \mathbf{W}^{-1} \mathbf{K} \mathbf{W}^{-1}$ , where  $\mathbf{W} = \text{diag}(\mathbf{K} \mathbf{1}_n)$  with  $\mathbf{1}_n$  be all one vectors and  $\tilde{\mathbf{W}} = \text{diag}(\mathbf{W}^{-1} \mathbf{K} \mathbf{W}^{-1} \mathbf{1}_n)$ . The method of constructing  $\mathbf{L}$  as described above guarantees that if the data are sampled from a manifold  $\mathcal{M}$ ,  $\mathbf{L}$  converges to  $\Delta_{\mathcal{M}}$  [HAvL05, THJ10]. A summary of the construction of  $\mathbf{L}$  can be found in Algorithm S2 LAPLACIAN. The last step of LE/DM algorithm embeds the data by solving the minimum eigenproblem of  $\mathbf{L}$ . The desired  $m$  dimensional embedding coordinates are obtained from the second to  $m + 1$ -th principal eigenvectors of graph Laplacian  $\mathbf{L}$ , with  $0 = \lambda_0 < \lambda_1 \leq \dots \leq \lambda_m$ , i.e.,  $\mathbf{y}_i = (\phi_1(\mathbf{x}_i), \dots, \phi_m(\mathbf{x}_i))$  (see also Supplement B).

To analyze ML algorithms, it is useful to consider the limit of the mapping  $\phi$  when the data is the entire manifold  $\mathcal{M}$ . We denote this limit also by  $\phi$ , and its image by  $\phi(\mathcal{M}) \in \mathbb{R}^m$ . For standard algorithms such as LE/DM, it is known that this limit exists [CL06, BN07, HAvL05, HAvL07, THJ10]. One of the fundamental requirements of ML is to preserve the neighborhood relations in the original data. In mathematical terms, we require that  $\phi : \mathcal{M} \rightarrow \phi(\mathcal{M})$  is a *smooth embedding*, i.e., that  $\phi$  is a smooth function (i.e. does not break existing neighborhood relations) whose Jacobian  $\mathbf{D}\phi(\mathbf{x})$  is full rank  $d$  at each  $\mathbf{x} \in \mathcal{M}$  (i.e. does not create new neighborhood relations).

**The pushforward Riemannian metric** A smooth  $\phi$  does not typically preserve geometric quantities such as distances along curves in  $\mathcal{M}$ . These concepts are captured by *Riemannian geometry*, and we additionally assume that  $(\mathcal{M}, g)$  is a *Riemannian manifold*, with the metric  $g$  induced from  $\mathbb{R}^D$ . One can always associate with  $\phi(\mathcal{M})$  a Riemannian metric  $g_{*\phi}$ , called the *pushforward Riemannian metric* [Lee03], which preserves the geometry of  $(\mathcal{M}, g)$ ;  $g_{*\phi}$  is defined by

$$\langle \mathbf{u}, \mathbf{v} \rangle_{g_{*\phi}(\mathbf{x})} = \langle \mathbf{D}\phi^{-1}(\mathbf{x})\mathbf{u}, \mathbf{D}\phi^{-1}(\mathbf{x})\mathbf{v} \rangle_{g(\mathbf{x})} \text{ for all } \mathbf{u}, \mathbf{v} \in \mathcal{T}_{\phi(\mathbf{x})}\phi(\mathcal{M}) \quad (1)$$

In the above,  $\mathcal{T}_{\mathbf{x}}\mathcal{M}$ ,  $\mathcal{T}_{\phi(\mathbf{x})}\phi(\mathcal{M})$  are tangent subspaces,  $\mathbf{D}\phi^{-1}(\mathbf{x})$  maps vectors from  $\mathcal{T}_{\phi(\mathbf{x})}\phi(\mathcal{M})$  to  $\mathcal{T}_{\mathbf{x}}\mathcal{M}$ , and  $\langle \cdot, \cdot \rangle$  is the Euclidean scalar product. For each  $\phi(\mathbf{x}_i)$ , the associated push-forward Riemannian metric expressed in the coordinates of  $\mathbb{R}^m$ , is a symmetric, semi-positive definite  $m \times m$  matrix  $\mathbf{G}(i)$  of rank  $d$ . The scalar product  $\langle \mathbf{u}, \mathbf{v} \rangle_{g_{*\phi}(\mathbf{x}_i)}$  takes the form  $\mathbf{u}^\top \mathbf{G}(i) \mathbf{v}$ . Given an embedding  $\mathbf{Y} = \phi(\mathbf{X})$ ,  $\mathbf{G}(i)$  can be estimated by Algorithm 1 (RMETRIC) of [PM13]. The RMETRIC algorithm also returns the *co-metric*  $\mathbf{H}(i)$ , which is the pseudo-inverse of the metric  $\mathbf{G}(i)$ , and

<sup>1</sup>In this paper, a smooth function or manifold will be assumed to be of class at least  $C^3$ .

its Singular Value Decomposition  $\Sigma(i), \mathbf{U}(i) \in \mathbb{R}^{m \times d}$ . The latter represents an orthogonal basis of  $\mathcal{T}_{\phi(\mathbf{x})}(\phi(\mathcal{M}))$ .

---

**Algorithm 1:** RMETRIC

---

**Input :** Embedding  $\mathbf{Y} \in \mathbb{R}^{n \times m}$ , Laplacian  $\mathbf{L}$ , intrinsic dimension  $d$

- 1 **for** all  $\mathbf{y}_i \in \mathbf{Y}, k = 1 \rightarrow m, l = 1 \rightarrow m$  **do**
- 2   |  $[\tilde{\mathbf{H}}(i)]_{kl} = \sum_{j \neq i} L_{ij}(y_{jl} - y_{il})(y_{jk} - y_{ik})$
- 3 **end**
- 4 **for**  $i = 1 \rightarrow n$  **do**
- 5   |  $\mathbf{U}(i), \Sigma(i) \leftarrow \text{REDUCEDRANKSVD}(\tilde{\mathbf{H}}(i), d)$
- 6   |  $\mathbf{H}(i) = \mathbf{U}(i)\Sigma(i)\mathbf{U}(i)^\top$
- 7   |  $\mathbf{G}(i) = \mathbf{U}(i)\Sigma^{-1}(i)\mathbf{U}(i)^\top$
- 8 **end**

**Return:**  $\mathbf{G}(i), \mathbf{H}(i) \in \mathbb{R}^{m \times m}, \mathbf{U}(i) \in \mathbb{R}^{m \times d}, \Sigma(i) \in \mathbb{R}^{d \times d}$ , for  $i \in [n]$

---

### 3 IES problem, related work, and challenges

**An example** Consider a continuous two dimensional strip with width  $W$ , height  $H$ , and *aspect ratio*  $W/H \geq 1$ , parametrized by coordinates  $w \in [0, W], h \in [0, H]$ . The eigenvalues and eigenfunctions of the Laplace-Beltrami operator  $\Delta$  with von Neumann boundary conditions [Str07] are  $\lambda_{k_1, k_2} = \left(\frac{k_1 \pi}{W}\right)^2 + \left(\frac{k_2 \pi}{H}\right)^2$ , respectively  $\phi_{k_1, k_2}(w, h) = \cos\left(\frac{k_1 \pi w}{W}\right) \cos\left(\frac{k_2 \pi h}{H}\right)$ . Eigenfunctions  $\phi_{1,0}, \phi_{0,1}$  are in bijection with the  $w, h$  coordinates (and give a full rank embedding), while the mapping by  $\phi_{1,0}, \phi_{2,0}$  provides no extra information regarding the second dimension  $h$  in the underlying manifold (and is rank 1). Theoretically, one can choose as coordinates eigenfunctions indexed by  $(k_1, 0), (0, k_2)$ , but, in practice,  $k_1$ , and  $k_2$  are usually unknown, as the eigenvalues are index by their rank  $0 = \lambda_0 < \lambda_1 \leq \lambda_2 \leq \dots$ . For a two dimensional strip, it is known [Str07] that  $\lambda_{1,0}$  always corresponds to  $\lambda_1$  and  $\lambda_{0,1}$  corresponds to  $\lambda_{\lceil W/H \rceil}$ . Therefore, when  $W/H > 2$ , the mapping of the strip to  $\mathbb{R}^2$  by  $\phi_1, \phi_2$  is low rank, while the mapping by  $\phi_1, \phi_{\lceil W/H \rceil}$  is full rank. Note that other mappings of rank 2 exist, e.g.,  $\phi_1, \phi_{\lceil W/H \rceil + 2}$  ( $k_1 = k_2 = 1$  in Figure 1b). These embeddings reflect progressively higher frequencies, as the corresponding eigenvalues grow larger.

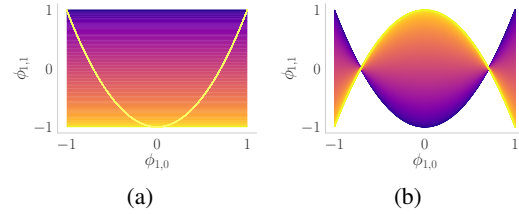


Figure 1: (a) Eigenfunction  $\phi_{1,0}$  versus  $\phi_{2,0}$  (curve) or  $\phi_{0,1}$  (two dimensional manifold). (b) Eigenfunction  $\phi_{1,0}$  versus  $\phi_{1,1}$ . All three manifolds are colored by the parameterization  $h$ .

**Prior work** [GZKR08] is the first work to give the IES problem a rigorous analysis. Their paper focuses on rectangles, and the failure illustrated in Figure 1a is defined as obtaining a mapping  $\mathbf{Y} = \phi(\mathbf{X})$  that is not *affinely equivalent* with the original data. They call this the *Price of Normalization* and explain it in terms of the variances along  $w$  and  $h$ . [DTCK18] is the first to frame the failure in terms of the rank of  $\phi_S = \{\phi_k : k \in S \subseteq [m]\}$ , calling it the *repeated eigendirection problem*. They propose a heuristic, LLRCOORDSEARCH, based on the observation that if  $\phi_k$  is a repeated eigendirection of  $\phi_1, \dots, \phi_{k-1}$ , one can fit  $\phi_k$  with *local linear regression* on predictors  $\phi_{[k-1]}$  with low leave-one-out errors  $r_k$ .

**Existence of solution** Before trying to find an algorithmic solution to the IES problem, we ask the question whether this is even possible, in the smooth manifold setting. Positive answers are given in [Por16], which proves that isometric embeddings by DM with finite  $m$  are possible, and more recently in [Bat14], which proves that any closed, connected Riemannian manifold  $\mathcal{M}$  can be smoothly embedded by its Laplacian eigenfunctions  $\phi_{[m]}$  into  $\mathbb{R}^m$  for some  $m$ , which depends only on the intrinsic dimension  $d$  of  $\mathcal{M}$ , the volume of  $\mathcal{M}$ , and lower bounds for *injectivity radius* and *Ricci curvature*. The example in Figure 1a demonstrates that, typically, not all  $m$  eigenfunctions

are needed. I.e., there exists a set  $S \subset [m]$ , so that  $\phi_S$  is also a smooth embedding. We follow [DTCK18] in calling such a set  $S$  *independent*. It is not known how to find an independent  $S$  analytically for a given  $\mathcal{M}$ , except in special cases such as the strip. In this paper, we propose a *finite sample* and algorithmic solution, and we support it with asymptotic theoretical analysis.

**The IES Problem** We are given data  $\mathbf{X}$ , and the output of an embedding algorithm (DM for simplicity)  $\mathbf{Y} = \phi(\mathbf{X}) = [\phi_1, \dots, \phi_m] \in \mathbb{R}^{n \times m}$ . We assume that  $\mathbf{X}$  is sampled from a  $d$ -dimensional manifold  $\mathcal{M}$ , with known  $d$ , and that  $m$  is sufficiently large so that  $\phi(\mathcal{M})$  is a smooth embedding. Further, we assume that there is a set  $S \subseteq [m]$ , with  $|S| = s \leq m$ , so that  $\phi_S$  is also a smooth embedding of  $\mathcal{M}$ . We propose to find such set  $S$  so that the rank of  $\phi_S$  is  $d$  on  $\mathcal{M}$  and  $\phi_S$  varies as slowly as possible.

**Challenges** (1) Numerically, and on a finite sample, distinguishing between a full rank mapping and a rank-defective one is imprecise. Therefore, we substitute for rank the volume of a unit parallelogram in  $\mathcal{T}_{\phi(\mathbf{x}_i)}\phi(\mathcal{M})$ . (2) Since  $\phi$  is *not* an isometry, we must separate the local distortions introduced by  $\phi$  from the estimated rank of  $\phi$  at  $\mathbf{x}$ . (3) Finding the optimal balance between the above desired properties. (4) In [Bat14] it is strongly suggested that  $s$  the number of eigenfunctions needed may exceed the *Whitney embedding dimension* ( $\leq 2d$ ) [Lee03], and that this number may depend on injectivity radius, aspect ratio, and so on. Section 5 shows an example of a flat 2-manifold, the *strip with cavity*, for which  $s > 2$ . In this paper, we assume that  $s$  and  $m$  are given and focus on selecting  $S$  with  $|S| = s$  unless otherwise stated; for completeness, in Section 5 we present a heuristic to select  $s$ .

**(Global) functional dependencies, knots and crossings** Before we proceed, we describe three different ways a mapping  $\phi(\mathcal{M})$  can fail to be invertible. The first, (*global*) *functional dependency* is the case when  $\text{rank } \mathbf{D}\phi < d$  on an open subset of  $\mathcal{M}$ , or on all of  $\mathcal{M}$  (yellow curve in Figure 1a); this is the case most widely recognized in the literature (e.g., [GZKR08, DTCK18]). The *knot* is the case when  $\text{rank } \mathbf{D}\phi < d$  at an isolated point (Figure 1b). Third, the *crossing* (Figure S8 in Supplement G) is the case when  $\phi : \mathcal{M} \rightarrow \phi(\mathcal{M})$  is not invertible at  $\mathbf{x}$ , but  $\mathcal{M}$  can be covered with open sets  $U$  such that the restriction  $\phi : U \rightarrow \phi(U)$  has full rank  $d$ . Combinations of these three exemplary cases can occur. The criteria and approach we define are based on the (surrogate) rank of  $\phi$ , therefore they will not rule out all crossings. We leave the problem of crossings in manifold embeddings to future work, as we believe that it requires an entirely separate approach (based, e.g., on the injectivity radius or density in the co-tangent bundle rather than differential structure).

## 4 Criteria and algorithm

### 4.1 A geometric criterion

We start with the main idea in evaluating the quality of a subset  $S$  of coordinate functions. At each data point  $i$ , we consider the orthogonal basis  $\mathbf{U}(i) \in \mathbb{R}^{m \times d}$  of the  $d$  dimensional tangent subspace  $\mathcal{T}_{\phi(\mathbf{x}_i)}\phi(\mathcal{M})$ . The projection of the columns of  $\mathbf{U}(i)$  onto the subspace  $\mathcal{T}_{\phi(\mathbf{x}_i)}\phi_S(\mathcal{M})$  is  $\mathbf{U}(i)[S, :] \equiv \mathbf{U}_S(i)$ . The following Lemma connects  $\mathbf{U}_S(i)$  and the co-metric  $\mathbf{H}_S(i)$  defined by  $\phi_S$ , with the *full*  $\mathbf{H}(i)$ .

**Lemma 1.** *Let  $\mathbf{H}(i) = \mathbf{U}(i)\mathbf{\Sigma}(i)\mathbf{U}(i)^\top$  be the co-metric defined by embedding  $\phi$ ,  $S \subseteq [m]$ ,  $\mathbf{H}_S(i)$  and  $\mathbf{U}_S(i)$  defined above. Then  $\mathbf{H}_S(i) = \mathbf{U}_S(i)\mathbf{\Sigma}(i)\mathbf{U}_S(i)^\top = \mathbf{H}(i)[S, S]$ .*

The proof is straightforward and left to the reader. Note that Lemma 1 is responsible for the efficiency of the search over sets  $S$ , given that the push-forward co-metric  $\mathbf{H}_S$  can be readily obtained as a submatrix of  $\mathbf{H}$ . Denote by  $\mathbf{u}_k^S(i)$  the  $k$ -th column of  $\mathbf{U}_S(i)$ . We further normalize each  $\mathbf{u}_k^S$  to length 1 and define the *normalized projected volume*  $\text{Vol}_{\text{norm}}(S, i) = \frac{\sqrt{\det(\mathbf{U}_S(i)^\top \mathbf{U}_S(i))}}{\prod_{k=1}^d \|\mathbf{u}_k^S(i)\|_2}$ . Conceptually,  $\text{Vol}_{\text{norm}}(S, i)$  is the volume spanned by a (non-orthonormal) “basis” of unit vectors in  $\mathcal{T}_{\phi_S(\mathbf{x}_i)}\phi_S(\mathcal{M})$ ;  $\text{Vol}_{\text{norm}}(S, i) = 1$  when  $\mathbf{U}_S(i)$  is orthogonal, and it is 0 when  $\text{rank } \mathbf{H}_S(i) < d$ . In Figure 1a, the  $\text{Vol}_{\text{norm}}(\{1, 2\})$  with  $\phi_{\{1, 2\}} = \{\phi_{1,0}, \phi_{2,0}\}$  is close to zero, since the projection of the two tangent vectors is parallel to the yellow curve; however  $\text{Vol}_{\text{norm}}(\{1, \lceil w/h \rceil\}, i)$  is almost 1, because the projections of the tangent vectors  $\mathbf{U}(i)$  will be (approximately) orthogo-

nal. Hence,  $\text{Vol}_{\text{norm}}(S, i)$  away from 0 indicates a non-singular  $\phi_S$  at  $i$ , and we use the average log  $\text{Vol}_{\text{norm}}(S, i)$ , which penalizes values near 0 highly, as the *rank quality*  $\mathfrak{R}(S)$  of  $S$ .

Higher frequency  $\phi_S$  maps with high  $\mathfrak{R}(S)$  may exist, being either smooth, such as the embeddings of the strip mentioned previously, or containing knots involving only small fraction of points, such as  $\phi_{\phi_{1,0}, \phi_{1,1}}$  in Figure 1a. To choose the lowest frequency, slowest varying smooth map, a regularization term consisting of the eigenvalues  $\lambda_k$ ,  $k \in S$ , of the graph Laplacian  $\mathbf{L}$  is added, obtaining the criterion

$$\mathfrak{L}(S; \zeta) = \underbrace{\frac{1}{n} \sum_{i=1}^n \log \sqrt{\det(\mathbf{U}_S(i)^\top \mathbf{U}_S(i))}}_{\mathfrak{R}_1(S) = \frac{1}{n} \sum_{i=1}^n \mathfrak{R}_1(S; i)} - \underbrace{\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^d \log \|\mathbf{u}_k^S(i)\|_2}_{\mathfrak{R}_2(S) = \frac{1}{n} \sum_{i=1}^n \mathfrak{R}_2(S; i)} - \zeta \sum_{k \in S} \lambda_k \quad (2)$$

## 4.2 Search algorithm

With this criterion, the IES problem turns into a subset selection problem parametrized by  $\zeta$

$$S_*(\zeta) = \underset{S \subseteq [m]; |S|=s; 1 \in S}{\text{argmax}} \mathfrak{L}(S; \zeta) \quad (3)$$

Note that we force the first coordinate  $\phi_1$  to always be chosen, since this coordinate cannot be functionally dependent on previous ones, and, in the case of DM, it also has lowest frequency. Note also that  $\mathfrak{R}_1$  and  $\mathfrak{R}_2$  are both submodular set function (proof in Supplement C.1). For large  $s$  and  $d$ , algorithms for optimizing over the difference of submodular functions can be used (e.g., see [IB12]). For the experiments in this paper, we have  $m = 20$  and  $d, s = 2 \sim 4$ , which enables us to use exhaustive search to handle (3). The exact search algorithm is summarized in Algorithm 2 INDEIGENSEARCH. A greedy variant is also proposed and analyzed in Supplement D.

---

### Algorithm 2: INDEIGENSEARCH

---

**Input** : Data  $\mathbf{X}$ , bandwidth  $\varepsilon$ , intrinsic dimension

$d$ , embedding dimension  $s$ , regularizer  $\zeta$

1  $\mathbf{Y} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{L}, \boldsymbol{\lambda} \in \mathbb{R}^m \leftarrow \text{DIFFMAP}(\mathbf{X}, \varepsilon)$

2  $\mathbf{U}(1), \dots, \mathbf{U}(n) \leftarrow \text{RMETRIC}(\mathbf{Y}, \mathbf{L}, d)$

3 **for**  $S \in \{S' \subseteq [m] : |S'| = s, 1 \in S'\}$  **do**

4      $\mathfrak{R}_1(S) \leftarrow 0; \mathfrak{R}_2(S) \leftarrow 0$

5     **for**  $i = 1, \dots, n$  **do**

6          $\mathbf{U}_S(i) \leftarrow \mathbf{U}(i)[S, :]$

7          $\mathfrak{R}_1(S) += \frac{1}{2n} \cdot \log \det(\mathbf{U}_S(i)^\top \mathbf{U}_S(i))$

8          $\mathfrak{R}_2(S) += \frac{1}{n} \cdot \sum_{k=1}^d \log \|\mathbf{u}_k^S(i)\|_2$

9     **end**

10      $\mathfrak{L}(S; \zeta) = \mathfrak{R}_1(S) - \mathfrak{R}_2(S) - \zeta \sum_{k \in S} \lambda_k$

11 **end**

12  $S_* = \text{argmax}_S \mathfrak{L}(S; \zeta)$

**Return**: Independent eigencoordinates set  $S_*$

---

## 4.3 Regularization path and choosing $\zeta$

According to (2), the optimal subset  $S_*$  depends on the parameter  $\zeta$ . The regularization path  $\ell(\zeta) = \max_{S \subseteq [m]; |S|=s; 1 \in S} \mathfrak{L}(S; \zeta)$  is the upper envelope of multiple lines (each correspond to a set  $S$ ) with slopes  $-\sum_{k \in S} \lambda_k$  and intercepts  $\mathfrak{R}(S)$ . The larger  $\zeta$  is, the more the lower frequency subset penalty prevails, and for sufficiently large  $\zeta$  the algorithm will output  $[s]$ . In the supervised learning framework, the regularization parameters are often chosen by cross validation. Here we propose a second criterion, that effectively limits how much  $\mathfrak{R}(S)$  may be ignored, or alternatively, bounds  $\zeta$  by a data dependent quantity. Define the *leave-one-out regret* of point  $i$  as follows

$$\mathfrak{D}(S, i) = \mathfrak{R}(S_*^i; [n] \setminus \{i\}) - \mathfrak{R}(S; [n] \setminus \{i\}) \text{ with } S_*^i = \text{argmax}_{S \subseteq [m]; |S|=s; 1 \in S} \mathfrak{R}(S; i) \quad (4)$$

In the above, we denote  $\mathfrak{R}(S; T) = \frac{1}{|T|} \sum_{i \in T} \mathfrak{R}_1(S; i) - \mathfrak{R}_2(S; i)$  for some subset  $T \subseteq [n]$ . The quantity  $\mathfrak{D}(S, i)$  in (4) measures the gain in  $\mathfrak{R}$  if all the other points  $[n] \setminus \{i\}$  choose the optimal subset  $S_*^i$ . If the regret  $\mathfrak{D}(S, i)$  is larger than zero, it indicates that the alternative choice might be better compared to original choice  $S$ . Note that the mean value for all  $i$ , i.e.,  $\frac{1}{n} \sum_i \mathfrak{D}(S, i)$  depends also on the variability of the optimal choice of points  $i$ ,  $S_*^i$ . Therefore, it might not favor an  $S$ , if  $S$  is optimal for every  $i \in [n]$ . Instead, we propose to inspect the distribution of  $\mathfrak{D}(S, i)$ , and remove the sets  $S$  for which  $\alpha$ 's percentile are larger than zero, e.g.,  $\alpha = 75\%$ , recursively from  $\zeta = \infty$  in decreasing order. Namely, the chosen set is  $S_* = S_*(\zeta')$  with  $\zeta' = \max_{\zeta \geq 0} \text{PERCENTILE}(\{\mathfrak{D}(S_*(\zeta), i)\}_{i=1}^n, \alpha) \leq 0$ . The optimal  $\zeta_*$  value is simply chosen to be the midpoint of all the  $\zeta$ 's that outputs set  $S_*$  i.e.,  $\zeta_* = \frac{1}{2}(\zeta' + \zeta'')$ , where  $\zeta'' = \min_{\zeta \geq 0} S_*(\zeta) = S_*(\zeta')$ . The procedure REGUPARAMSEARCH is summarized in Algorithm 3.

---

**Algorithm 3: REGUPARAMSEARCH**

---

**Input** : Threshold parameter  $\alpha$ 

```
1 for  $\zeta = \zeta_{\max} \rightarrow 0$  do
   $\triangleright \zeta_{\max}$  should be sufficiently large such that  $S_*(\zeta_{\max}) = [s]$ 
2    $S \leftarrow S_*(\zeta)$ ;  $S_* \leftarrow \text{NULL}$ ;  $\zeta'' \leftarrow \text{NULL}$ 
3   for  $i \in [n]$  do
4      $\mathcal{D}(S, i) \leftarrow \mathfrak{R}(S_*^i; [n] \setminus \{i\}) - \mathfrak{R}(S; [n] \setminus \{i\})$  from equation (4)
5   end
6   if  $\text{PERCENTILE}(\{\mathcal{D}(S, i)\}_{i=1}^n, \alpha) \leq 0$  and  $S_* = \text{NULL}$  then
7     Optimal set  $S_* \leftarrow S$ 
8      $\zeta' \leftarrow \zeta \triangleright$  First found a set that satisfies the criterion.
9   else if  $S_* \neq \text{NULL}$  and  $S_* = S_*(\zeta)$  then
10     $\zeta'' \leftarrow \zeta \triangleright$  Searching for  $\zeta''$ 
11  else if  $S_* \neq \text{NULL}$  and  $\zeta'' \neq \text{NULL}$  and  $S_* \neq S_*(\zeta)$  then
12     $\zeta_* \leftarrow \frac{1}{2}(\zeta' + \zeta'')$ 
13    break  $\triangleright$  Leave the loop when found  $\zeta'' = \min_{\zeta \geq 0} S_*(\zeta') = S_*(\zeta)$ 
14  else
15    continue
16  end
17 end
```

**Return:** Optimal set  $S_*$ , optimal regularization parameter  $\zeta_*$ 

---

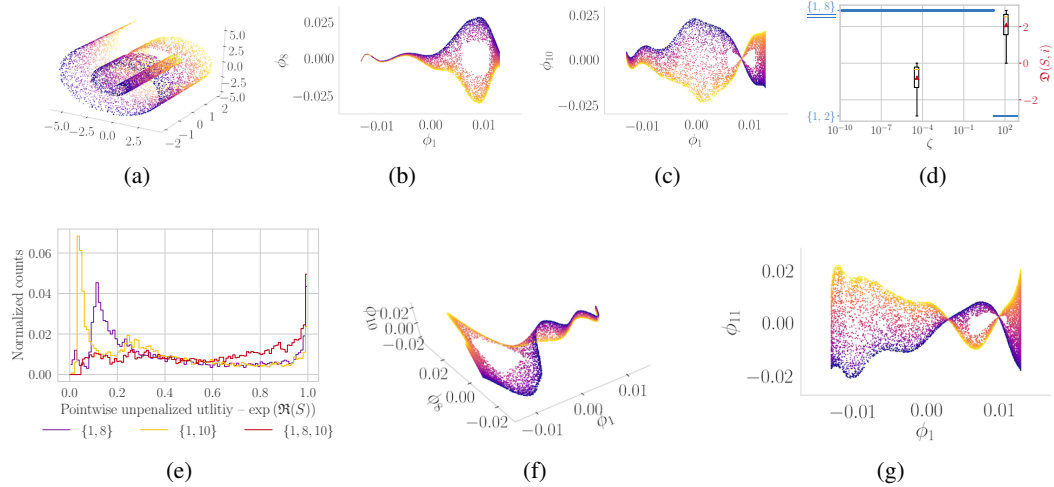


Figure 2: (a) Original data of  $\mathcal{D}_4$ , *swiss roll with hole* dataset. Embeddings with coordinate subset to be (b)  $S = \{1, 8\}$ , (c)  $S = \{1, 10\}$ , (f)  $S = \{1, 8, 10\}$  and (g)  $S = \{1, 11\}$  on  $\mathcal{D}_4$ . (e) Histogram of point-wise normalized projected volume on  $\mathcal{D}_4$  for top two ranking of subsets (purple and yellow) and the union of two sets (red) obtain from INDEIGENSEARCH algorithm.

## 5 A heuristic to determine whether $s$ is sufficiently large

In this section, we propose a heuristic method to determine whether the given  $s$  is large enough. Our method is based on the histogram of  $\text{Vol}_{\text{norm}}(S, i) = \exp(\mathfrak{R}(S, i))$ , the *normalized projected volume* of each point  $i$ . Recall that this volume is bounded between 0 and 1. Ideally, a perfect choice of cardinality  $|S|$  will result in a concentration of mass in larger  $\text{Vol}_{\text{norm}}$  region. The heuristic works as follow: at first we check the histogram of unpenalized  $\text{Vol}_{\text{norm}}$  on the top few ranked subsets in terms of  $\mathcal{L}$ . If spikes in the small  $\text{Vol}_{\text{norm}}$  regions are witnessed in the histogram, taking the union

of the subsets and inspecting the histogram of unpenalized  $\text{Vol}_{\text{norm}}$  on the combined set again. If spikes in small  $\text{Vol}_{\text{norm}}$  region diminished, one can conclude that a larger cardinality size  $|S|$  is needed for such manifold.

We illustrate the idea on *swiss roll with hole* dataset in Figure 2a. Figure 2b is the optimal subset of coordinates  $S_* = \{1, 8\}$  selected by the proposed algorithm that best parameterize the underlying manifold. Figure 2d suggests one should eliminate  $S_0 = \{1, 2\}$  because  $\mathfrak{D}(S_0, i) \geq 0$  for all the points by REGUPARAMSEARCH. However, as shown in Figure 2b, though it has low frequency and having rank 2 for most of the places, set  $\{1, 8\}$  might not be suitable for data analysis for the very thin arms in left side of the embedding. Figure 2e is the histograms of the point-wise unpenalized  $\text{Vol}_{\text{norm}}$  on different subsets. Purple and yellow curves correspond to the histogram of top two ranked subsets  $S$  from INDEIGENSEARCH. Both curves show a concentration of masses in small  $\text{Vol}_{\text{norm}}$  region. The histogram of point-wise unpenalized  $\text{Vol}_{\text{norm}}$  on  $\{1, 8, 10\}$  (red curve), which is the union of the aforementioned two subsets, shows less concentration in the small  $\text{Vol}_{\text{norm}}$  region and implies that  $|S| = 3$  might be a better choice for data analysis. Figure 2f shows the embedding with coordinate  $S = \{1, 8, 10\}$ , which represents a two dimensional strip embedded in three dimensional space. The thin arc in Figure 2b turns out to be a collapsed two dimensional manifold via projection, as shown in the upper right part of Figure 2f and left part of Figure 2c. Here we have to restate that the embedding in Figure 2b, although is a *degenerated* embedding, is still the best set one can choose for  $s = 2$  such that the embedding varies slowest and has rank 2. However, choosing  $s = 3$  might be better for data analysis.

## 6 $\mathfrak{R}$ as Kullbach-Leibler divergence

In this section we analyze  $\mathfrak{R}$  in its population version, and show that it is reminiscent of a Kullbach-Leibler divergence between *unnormalized* measures on  $\phi_S(\mathcal{M})$ . The population version of the regularization term takes the form of a well-known *smoothness* penalty on the embedding coordinates  $\phi_S$ .

**Volume element and the Riemannian metric** Consider a Riemannian manifold  $(\mathcal{M}, g)$  mapped by a smooth embedding  $\phi_S$  into  $(\phi_S(\mathcal{M}), g_{*\phi_S})$ ,  $\phi_S : \mathcal{M} \rightarrow \mathbb{R}^s$ , where  $g_{*\phi_S}$  is the *push-forward* metric defined in (1). A Riemannian metric  $g$  induces a *Riemannian measure* on  $\mathcal{M}$ , with volume element  $\sqrt{\det g}$ . Denote now by  $\mu_{\mathcal{M}}$ , respectively  $\mu_{\phi_S(\mathcal{M})}$  the Riemannian measures corresponding to the metrics induced on  $\mathcal{M}$ ,  $\phi_S(\mathcal{M})$  by the ambient spaces  $\mathbb{R}^D, \mathbb{R}^s$ ; let  $g$  be the former metric.

**Lemma 2.** *Let  $S, \phi, \phi_S, \mathbf{H}_S(\mathbf{x}), \mathbf{U}_S(\mathbf{x}), \Sigma(\mathbf{x})$  be defined as in Section 4 and Lemma 1. For simplicity, we denote by  $\mathbf{H}_S(\mathbf{y}) \equiv \mathbf{H}_S(\phi_S^{-1}(\mathbf{y}))$ , and similarly for  $\mathbf{U}_S(\mathbf{y}), \Sigma(\mathbf{y})$ . Assume that  $\phi_S$  is a smooth embedding. Then, for any measurable function  $f : \mathcal{M} \rightarrow \mathbb{R}$ ,*

$$\int_{\mathcal{M}} f(\mathbf{x}) d\mu_{\mathcal{M}}(\mathbf{x}) = \int_{\phi_S(\mathcal{M})} f(\phi_S^{-1}(\mathbf{y})) j_S(\mathbf{y}) d\mu_{\phi_S(\mathcal{M})}(\mathbf{y}), \quad (5)$$

with

$$j_S(\mathbf{y}) = 1/\text{Vol}(\mathbf{U}_S(\mathbf{y})\Sigma_S^{1/2}(\mathbf{y})). \quad (6)$$

**Proof.** Let  $\mu_{\phi_S(\mathcal{M})}^*$  denote the Riemannian measure induced by  $g_{*\phi_S}$ . Since  $(\mathcal{M}, g)$  and  $(\phi_S(\mathcal{M}), g_{*\phi_S})$  are isometric by definition,  $\int_{\mathcal{M}} f(\mathbf{x}) d\mu_{\mathcal{M}}(\mathbf{x}) = \int_{\phi_S(\mathcal{M})} f(\phi_S^{-1}(\mathbf{y})) d\mu_{\phi_S(\mathcal{M})}^*(\mathbf{y}) = \int_{\phi_S(\mathcal{M})} f(\phi_S^{-1}(\mathbf{y})) \sqrt{\det g_{*\phi_S}(\mathbf{y})} d\mu_{\phi_S(\mathcal{M})}(\mathbf{y})$  follows from the change of variable formula. It remains to find the expression of  $j_S(\mathbf{y}) = \sqrt{\det g_{*\phi_S}(\mathbf{y})}$ . The matrix  $\mathbf{U}_S(\mathbf{y})$  (note that  $\mathbf{U}_S(\mathbf{y})$  is *not orthogonal*) can be written as

$$\mathbf{U}_S(\mathbf{y}) = \mathbf{V}\mathbf{Q}_S(\mathbf{y}) \quad (7)$$

where  $\mathbf{V} \in \mathbb{R}^{s \times d}$  is an orthogonal matrix and  $\mathbf{Q}_S(\mathbf{y}) \in \mathbb{R}^{d \times d}$  is upper triangular. Then,

$$\mathbf{H}_S(\mathbf{y}) = \mathbf{U}_S(\mathbf{y})\Sigma(\mathbf{y})\mathbf{U}_S(\mathbf{y})^\top = \mathbf{V}_S(\mathbf{y}) \underbrace{(\mathbf{Q}_S(\mathbf{y})\Sigma(\mathbf{y})\mathbf{Q}_S(\mathbf{y})^\top)}_{\hat{\mathbf{H}}_S(\mathbf{y})} \mathbf{V}_S(\mathbf{y})^\top. \quad (8)$$

In the above  $\tilde{\mathbf{H}}_S(y)$  is the co-metric expressed in the new coordinate system induced by  $\mathbf{V}_S(\mathbf{y})$ . Hence, in the same basis,  $g_{*\phi_S}$  is expressed by

$$\tilde{\mathbf{G}}_S(y) = \tilde{\mathbf{H}}_S(y)^{-1} = (\mathbf{Q}_S(\mathbf{y})\boldsymbol{\Sigma}(\mathbf{y})\mathbf{Q}_S(\mathbf{y})^\top)^{-1}. \quad (9)$$

The volume element, which is invariant to the chosen coordinate system, is

$$\det(\mathbf{Q}_S(\mathbf{y})\boldsymbol{\Sigma}(\mathbf{y})\mathbf{Q}_S(\mathbf{y})^\top)^{-1/2} = \prod_{k=1}^d \sigma_k(\mathbf{y})^{-1/2} q_{S,kk}(\mathbf{y})^{-1}. \quad (10)$$

From (7), it follows also that

$$\det(\mathbf{Q}_S(\mathbf{y})\boldsymbol{\Sigma}(\mathbf{y})\mathbf{Q}_S(\mathbf{y})^\top)^{-1/2} = 1/\text{Vol}(\mathbf{U}_S(\mathbf{y})\boldsymbol{\Sigma}(\mathbf{y})^{1/2}) \quad (11)$$

■

**Asymptotic limit of  $\mathfrak{R}$**  We now study the first term of our criterion in the limit of infinite sample size. We make the following assumptions.

**Assumption 1.** *The manifold  $\mathcal{M}$  is compact of class  $\mathcal{C}^3$ , and there exists a set  $S$ , with  $|S| = s$  so that  $\phi_S$  is a smooth embedding of  $\mathcal{M}$  in  $\mathbb{R}^s$ .*

**Assumption 2.** *The data are sampled from a distribution on  $\mathcal{M}$  continuous with respect to  $\mu_{\mathcal{M}}$ , whose density is denoted by  $p$ .*

**Assumption 3.** *The estimate of  $\mathbf{H}_S$  in Algorithm 1 computed w.r.t. the embedding  $\phi_S$  is consistent.*

We know from [Bat14] that Assumption 1 is satisfied for the DM/LE embedding. The remaining assumptions are minimal requirements ensuring that limits of our quantities exist. Now consider the setting in Sections 3, in which we have a larger set of eigenfunctions,  $\phi_{[m]}$  so that  $[m]$  contains the set  $S$  of Assumption 1. Denote by  $\tilde{j}_S(\mathbf{y}) = \prod_{k=1}^d (|u_k^S(\mathbf{y})| |\sigma_k(\mathbf{y})|^{1/2})^{-1}$  a new volume element.

**Theorem 3** (Limit of  $\mathfrak{R}$ ). *Under Assumptions 1–3,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_i \ln \mathfrak{R}(S, \mathbf{x}_i) = \mathfrak{R}(S, \mathcal{M}), \quad (12)$$

and

$$\mathfrak{R}(S, \mathcal{M}) = - \int_{\phi_S(\mathcal{M})} \ln \frac{j_S(\mathbf{y})}{\tilde{j}_S(\mathbf{y})} p(\phi_S^{-1}(\mathbf{y})) j_S(\mathbf{y}) d\mu_{\phi_S(\mathcal{M})}(\mathbf{y}) \stackrel{\text{def}}{=} -D(pj_S \| p\tilde{j}_S) \quad (13)$$

**Proof.** Because  $\phi_S$  is a smooth embedding,  $j_S(\mathbf{y}) > 0$  on  $\phi_S(\mathcal{M})$ , and because  $\mathcal{M}$  is compact,  $\min_{\phi_S(\mathcal{M})} j_S(\mathbf{y}) > 0$ . Similarly, noting that  $\tilde{j}_S(\mathbf{y}) \geq \prod_{k=1}^d \sigma_k^{-1/2}(\mathbf{y})$ , we conclude that  $\tilde{j}_S(\mathbf{y})$  is also bounded away from 0 on  $\mathcal{M}$ . Therefore  $\ln j_S(\mathbf{y})$  and  $\ln \tilde{j}_S(\mathbf{y})$  are bounded, and the integral in the r.h.s. of (13) exists and has a finite value. Now,

$$\frac{1}{n} \sum_i \ln \mathfrak{R}(S, \mathbf{x}_i) \rightarrow \int_{\mathcal{M}} \ln \mathfrak{R}(S, \mathbf{x}) p(\mathbf{x}) d\mu_{\mathcal{M}}(\mathbf{x}) = \mathfrak{R}(S, \mathcal{M}). \quad (14)$$

$$\begin{aligned} & \int_{\mathcal{M}} \ln \mathfrak{R}(S, \mathbf{x}) p(\mathbf{x}) d\mu_{\mathcal{M}}(\mathbf{x}) \\ &= \int_{\phi_S(\mathcal{M})} \ln \mathfrak{R}(\phi_S^{-1}(\mathbf{y})) p(\phi_S^{-1}(\mathbf{y})) j_S(\mathbf{y}) d\mu_{\phi_S(\mathcal{M})}(\mathbf{y}) \\ &= \int_{\phi_S(\mathcal{M})} \left[ \frac{1}{2} \ln \frac{\text{Vol}(\mathbf{U}_S^\top(\mathbf{y})\mathbf{U}_S(\mathbf{y}))}{\tilde{j}_S(\mathbf{y})} - \frac{p(\phi_S^{-1}(\mathbf{y})) \prod_{k=1}^d \sigma_k^{1/2}(\mathbf{y})}{p(\phi_S^{-1}(\mathbf{y})) \prod_{k=1}^d \sigma_k^{1/2}(\mathbf{y})} \right] p(\phi_S^{-1}(\mathbf{y})) j_S(\mathbf{y}) d\mu_{\phi_S(\mathcal{M})}(\mathbf{y}) \\ &= \int_{\phi_S(\mathcal{M})} \ln \frac{j_S(\mathbf{y}) p(\phi_S^{-1}(\mathbf{y}))}{\tilde{j}_S(\mathbf{y}) p(\phi_S^{-1}(\mathbf{y}))} p(\phi_S^{-1}(\mathbf{y})) j_S(\mathbf{y}) d\mu_{\phi_S(\mathcal{M})}(\mathbf{y}) = -D(pj_S \| p\tilde{j}_S) \end{aligned} \quad (15)$$

■



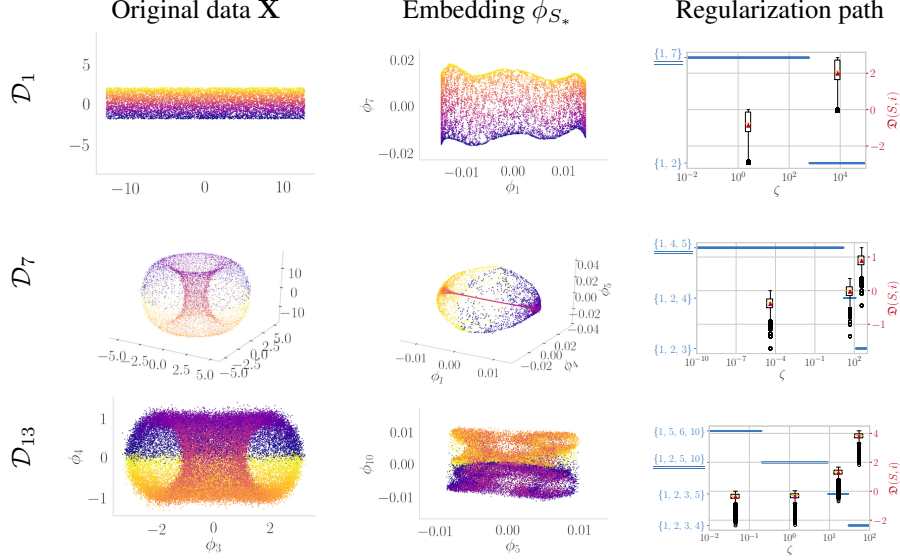


Figure 3: Experimental result for synthetic datasets. Rows correspond to different synthetic datasets (please refer to Table S2). Optimal subset  $S_*$  is selected by INDEIGENSEARCH.

Note that  $D(\cdot\|\cdot)$  is a Kullback-Leibler divergence, where the measures defined by  $p_{j_S}, p_{\tilde{j}_S}$  normalize to different values; because  $j_S \geq \tilde{j}_S$  the divergence  $D$  is always positive.

It is known that  $\lambda_k$ , the  $k$ -th eigenvalue of the Laplacian, converges under certain technical conditions [BN07] to an eigenvalue of the Laplace-Beltrami operator  $\Delta_{\mathcal{M}}$  and that

$$\lambda_k(\Delta_{\mathcal{M}}) = \langle \phi_k, \Delta_{\mathcal{M}} \phi_k \rangle = \int_{\mathcal{M}} \|\text{grad } \phi_k(\mathbf{x})\|_2^2 d\mu(\mathcal{M}). \quad (17)$$

Hence, a smaller value for the regularization term encourages the use of slow varying coordinate functions, as measured by the squared norm of their gradients, as in equation (17). Hence, under Assumptions 1, 2, 3,  $\mathcal{L}$  converges to

$$\mathcal{L}(S, \mathcal{M}) = -D(p_{j_S}\|p_{\tilde{j}_S}) - \frac{\zeta}{\lambda_1(\mathcal{M})} \sum_{k \in S} \lambda_k(\mathcal{M}). \quad (18)$$

The rescaling of  $\zeta$  in comparison with equation (2) aims to make  $\zeta$  adimensional, whereas the eigenvalues scale with the volume of  $\mathcal{M}$ .

## 7 Experiments

We demonstrate the proposed algorithm on three synthetic datasets, one where the minimum embedding dimension  $s$  equals  $d$  ( $\mathcal{D}_1$  long strip), and two ( $\mathcal{D}_7$  high torus and  $\mathcal{D}_{13}$  three torus) where  $s > d$ . The complete list of synthetic manifolds (transformations of 2 dimensional strips, 3 dimensional cubes, two and three tori, etc.) investigated can be found in Supplement G and Table S2. The examples have (i) aspect ratio of at least 4 (ii) points sampled *non-uniformly* from the underlying manifold  $\mathcal{M}$ , and (iii) Gaussian noise added. The sample size of the synthetic datasets is  $n = 10,000$  unless otherwise stated. Additionally, we analyze several real datasets from chemistry and astronomy. All embeddings are computed with the DM algorithm, which outputs  $m = 20$  eigenvectors. Hence, we examine 171 sets for  $s = 3$  and 969 sets for  $s = 4$ . No more than 2 to 5 of these sets appear on the regularization path. Detailed experimental results are in Table S3. In this section, we show the original dataset  $\mathbf{X}$ , the embedding  $\phi_{S_*}$ , with  $S_*$  selected by INDEIGENSEARCH and  $\zeta_*$  from REGUPARAMSEARCH, and the maximizer sets on the regularization path with box plots of  $\mathcal{D}(S, i)$  as discussed in Section 4. The  $\alpha$  threshold for REGUPARAMSEARCH is set to 75%. All the experiments are replicated for more than 5 times, and the outputs are similar because of the large sample size  $n$ .

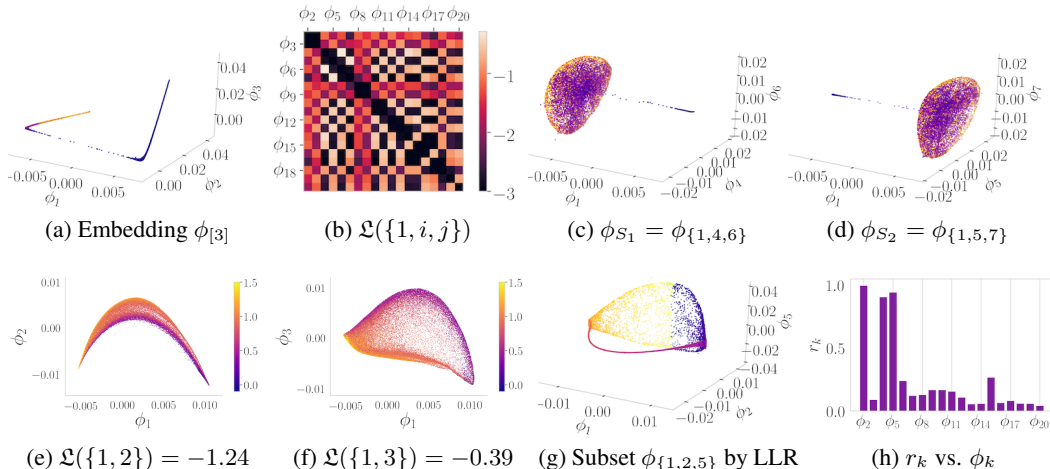


Figure 4: First row: Chloromethane dataset; second row: SDSS dataset in (e), (f) and (g), (h) show the example when LLR failed. (c) and (d) are embeddings with top two ranked subsets  $S_1$  and  $S_2$ , colored by the distances between C and two different  $\text{Cl}^-$ , respectively. (e) and (f) are embeddings of  $\phi_{\{1,2\}}$  (suboptimal set) and  $\phi_{\{1,3\}}$  (maximizer of  $\mathcal{L}$ ), respectively (values shown in caption).

**Synthetic manifolds** The results of synthetic manifolds are in Figure 3. (i) Manifold with  $s = d$ . The first synthetic dataset we considered,  $\mathcal{D}_1$ , is a two dimensional strip with aspect ratio  $W/H = 2\pi$ . Left panel of the top row shows the scatter plot of such dataset. From the theoretical analysis in Section 3, the coordinate set that corresponds to slowest varying unique eigendirection is  $S = \{1, \lceil W/H \rceil\} = \{1, 7\}$ . Middle panel, with  $S_* = \{1, 7\}$  selected by INDEIGENSEARCH with  $\zeta$  chosen by REGUPARAMSEARCH, confirms this. The right panel shows the box plot of  $\{\mathfrak{D}(S, i)\}_{i=1}^n$ . According to the proposed procedure, we eliminate  $S_0 = \{1, 2\}$  since  $\mathfrak{D}(S_0, i) \geq 0$  for almost all the points. (ii) Manifold with  $s > d$ . The second data  $\mathcal{D}_7$  is displayed in the left panel of the second row. Due to the mechanism we used to generate the data, the resultant torus is non-uniformly distributed along the z axis. Middle panel is the embedding of the optimal coordinate set  $S_* = \{1, 4, 5\}$  selected by INDEIGENSEARCH. Note that the middle region (in red) is indeed a two dimensional narrow tube when zoomed in. The right panel indicates that both  $\{1, 2, 3\}$  and  $\{1, 2, 4\}$  (median is around zero) should be removed. The optimal regularization parameter is  $\zeta_* \approx 7$ . The result of the third dataset  $\mathcal{D}_{13}$ , *three torus*, is in the third row of the figure. We displayed only projections of the penultimate and the last coordinate of original data  $\mathbf{X}$  and embedding  $\phi_{S_*}$  (which is  $\{5, 10\}$ ) colored by  $\alpha_1$  of (S7) in the left and middle panel to conserve space. A full combinations of coordinates can be found in Figure S5. The right panel implies one should eliminate the set  $\{1, 2, 3, 4\}$  and  $\{1, 2, 3, 5\}$  since both of them have more than 75% of the points such that  $\mathfrak{D}(S, i) \geq 0$ . The first remaining subset is  $\{1, 2, 5, 10\}$ , which yields an optimal regularization parameter  $\zeta_* \approx 5$ .

**Molecular dynamics dataset [FTP16]** In SN2 reaction molecular dynamics of chloromethane [FTP16] dataset, two chloride atoms substitute with each other in different configurations/points  $\mathbf{x}_i$  as described in the following chemical equation  $\text{CH}_3\text{Cl} + \text{Cl}^- \longleftrightarrow \text{CH}_3\text{Cl} + \text{Cl}^-$ . The dataset exhibits some kind of clustering structure with a sparse connection between two clusters which represents the time when the substitution happened. The dataset has size  $n \approx 30,000$  and ambient dimension  $D = 40$ , with the intrinsic dimension estimate be  $\hat{d} = 2$  The embedding with coordinate set  $S = [3]$  is shown in Figure 4a. The first three eigenvectors parameterize the same directions, which yields a one dimensional manifold in the figure. Top view ( $S = [2]$ ) of the figure is a u-shaped structure similar to the yellow curve in Figure 1a. The heat map of  $\mathcal{L}(\{1, i, j\})$  for different combinations of coordinates in Figure 4b confirms that  $\mathcal{L}$  for  $S = [3]$  is low and that  $\phi_1, \phi_2$  and  $\phi_3$  give a low rank mapping. The heat map also shows high  $\mathcal{L}$  values for  $S_1 = \{1, 4, 6\}$  or  $S_2 = \{1, 5, 7\}$ , which correspond to the top two ranked subsets. The embeddings with  $S_1, S_2$  are in Figures 4c and 4d, respectively. In this case, we obtain two optimal  $S$  sets due to the data symmetry.

**Galaxy spectra from the Sloan Digital Sky survey (SDSS)**<sup>2</sup> [AAMA<sup>+</sup>09], preprocessed as in [MMVZ16]. We display a sample of  $n = 50,000$  points from the first 0.3 million points which correspond to closer galaxies. Figures 4e and 4f show that the first two coordinates are almost dependent; the embedding with  $S_* = \{1, 3\}$  is selected by INDEIGENSEARCH with  $d = 2$ . Both plots are colored by the blue spectrum magnitude, which is correlated to the number of young stars in the galaxy, showing that this galaxy property varies smoothly and non-linearly with  $\phi_1, \phi_3$ , but is not smooth w.r.t.  $\phi_1, \phi_2$ .

**Comparison with [DTCK18]** The LLRCOORDSEARCH method outputs similar candidate coordinates as our proposed algorithm most of the time (see Table S3). However, the results differ for *high torus* as in Figure 4. Figure 4h is the leave one out (LOO) error  $r_k$  versus coordinates. The coordinates chosen by LLRCOORDSEARCH was  $S = \{1, 2, 5\}$ , as in Figure 4g. The embedding is clearly shown to be suboptimal, for it failed to capture the cavity within the torus. This is because the algorithm searches in a sequential fashion; the noise eigenvector  $\phi_2$  in this example appears before the signal eigenvectors e.g.,  $\phi_4$  and  $\phi_5$ .

**Additional experiments with real data** are shown in Table 1. Not surprisingly, for most real data sets we examined, the independent coordinates are not the first  $s$ . They also show that the algorithm scales well and is robust to the noise present in real data.

Table 1: Results for other real datasets. Columns from left to right are sample size  $n$ , ambient dimension of data  $D$ , average degree of neighbor graph  $\text{deg}_{\text{avg}}$ ,  $(s, d)$  and runtime for IES, and the chosen set  $S^*$ , respectively. Last three datasets are from [CTS<sup>+</sup>17].

	$n$	$D$	$\text{deg}_{\text{avg}}$	$(s, d)$	$t$ (sec)	$S^*$
SDSS (full)	298,511	3750	144.91	(2, 2)	106.05	(1, 3)
Aspirin	211,762	244	101.03	(4, 3)	85.11	(1, 2, 3, 7)
Ethanol	555,092	102	107.27	(3, 2)	233.16	(1, 2, 4)
Malondialdehyde	993,237	96	106.51	(3, 2)	459.53	(1, 2, 3)

The asymptotic runtime of LLRCOORDSEARCH has quadratic dependency on  $n$ , while for our algorithm is linear in  $n$ . Details of runtime analysis are Supplement E. LLRCOORDSEARCH was too slow to be tested on the four larger datasets (see also Figure S1).

## 8 Conclusion

Algorithms that use eigenvectors, such as DM, are among the most promising and well studied in ML. It is known since [GZKR08] that when the aspect ratio of a low dimensional manifold exceeds a threshold, the choice of eigenvectors becomes non-trivial, and that this threshold can be as low as 2. Our experimental results confirm the need to augment ML algorithms with IES methods in order to successfully apply ML to real world problems. Surprisingly, the IES problem has received little attention in the ML literature, to the extent that the difficulty and complexity of the problem have not been recognized. Our paper advances the state of the art by (i) introducing for the first time a differential geometric definition of the problem, (ii) highlighting geometric factors such as injectivity radius that, in addition to aspect ratio, influence the number of eigenfunctions needed for a smooth embedding, (iii) constructing selection criteria based on *intrinsic manifold quantities*, (iv) which have analyzable asymptotic limits, (v) can be computed efficiently, and (vi) are also robust to the noise present in real scientific data. The library of hard synthetic examples we constructed will be made available along with the python software implementation of our algorithms.

## Acknowledgements

The authors acknowledge partial support from the U.S. Department of Energy’s Office of Energy Efficiency and Renewable Energy (EERE) under the Solar Energy Technologies Office Award Number DE-EE0008563 and from the NSF DMS PD 08-1269 and NSF IIS-0313339 awards. They are

<sup>2</sup>The Sloan Digital Sky Survey data can be downloaded from <https://www.sdss.org>

grateful to the Tkatchenko and Pfaendtner labs and in particular to Stefan Chmiela and Chris Fu for providing the molecular dynamics data and for many hours of brainstorming and advice.

## Disclaimer

The views expressed herein do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

## References

- [AAMA<sup>+</sup>09] Kevork N Abazajian, Jennifer K Adelman-McCarthy, Marcel A Agüeros, Sahar S Allam, Carlos Allende Prieto, Deokkeun An, Kurt SJ Anderson, Scott F Anderson, James Annis, Neta A Bahcall, et al. The seventh data release of the sloan digital sky survey. *The Astrophysical Journal Supplement Series*, 182(2):543, 2009.
- [Bat14] Jonathan Bates. The embedding dimension of laplacian eigenfunction maps. *Applied and Computational Harmonic Analysis*, 37(3):516–530, 2014.
- [BN07] Mikhail Belkin and Partha Niyogi. Convergence of laplacian eigenmaps. In B. Schölkopf, J. C. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19*, pages 129–136. MIT Press, 2007.
- [CL06] R. R. Coifman and S. Lafon. Diffusion maps. *Applied and Computational Harmonic Analysis*, 30(1):5–30, 2006.
- [CTS<sup>+</sup>17] Stefan Chmiela, Alexandre Tkatchenko, Huziel E Saucedo, Igor Poltavsky, Kristof T Schütt, and Klaus-Robert Müller. Machine learning of accurate energy-conserving molecular force fields. *Science advances*, 3(5):e1603015, 2017.
- [Dry16] I. L. (Ian L.) Dryden. *Statistical shape analysis : with applications in R*. Wiley series in probability and statistics. Wiley, Chichester, West Sussex, England, 2nd ed. edition, 2016.
- [DS13] Sanjoy Dasgupta and Kaushik Sinha. Randomized partition trees for exact nearest neighbor search. In *Conference on Learning Theory*, pages 317–337, 2013.
- [DTCK18] Carmeline J Dsilva, Ronen Talmon, Ronald R Coifman, and Ioannis G Kevrekidis. Parsimonious representation of nonlinear dynamical systems through manifold learning: A chemotaxis case study. *Applied and Computational Harmonic Analysis*, 44(3):759–773, 2018.
- [FTP16] Kelly L. Fleming, Pratyush Tiwary, and Jim Pfaendtner. New approach for investigating reaction dynamics and rates with ab initio calculations. *Journal of Physical Chemistry A*, 120(2):299–305, 2016.
- [GZKR08] Yair Goldberg, Alon Zakai, Dan Kushnir, and Ya’acov Ritov. Manifold learning: The price of normalization. *Journal of Machine Learning Research*, 9(Aug):1909–1939, 2008.
- [Har98] David A Harville. *Matrix algebra from a statistician’s perspective*, 1998.
- [HAvL05] Matthias Hein, Jean-Yves Audibert, and Ulrike von Luxburg. From graphs to manifolds - weak and strong pointwise consistency of graph laplacians. In *Learning Theory, 18th Annual Conference on Learning Theory, COLT 2005, Bertinoro, Italy, June 27-30, 2005, Proceedings*, pages 470–485, 2005.
- [HAvL07] Matthias Hein, Jean-Yves Audibert, and Ulrike von Luxburg. Graph laplacians and their convergence on random neighborhood graphs. *Journal of Machine Learning Research*, 8:1325–1368, 2007.
- [HHJ90] Roger A Horn, Roger A Horn, and Charles R Johnson. *Matrix analysis*. Cambridge university press, 1990.
- [IB12] Rishabh Iyer and Jeff Bilmes. Algorithms for approximate minimization of the difference between submodular functions, with applications. In *Proceedings of the Twenty-Eighth Conference on Uncertainty in Artificial Intelligence, UAI’12*, pages 407–417, Arlington, Virginia, United States, 2012. AUAI Press.

- [LB05] Elizaveta Levina and Peter J Bickel. Maximum likelihood estimation of intrinsic dimension. In *Advances in neural information processing systems*, pages 777–784, 2005.
- [Lee03] John M. Lee. Introduction to smooth manifolds, 2003.
- [MHM18] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- [MMVZ16] James McQueen, Marina Meilă, Jacob VanderPlas, and Zhongyue Zhang. Megaman: Scalable manifold learning in python. *Journal of Machine Learning Research*, 17(148):1–5, 2016.
- [NLCK06] Boaz Nadler, Stephane Lafon, Ronald Coifman, and Ioannis Kevrekidis. Diffusion maps, spectral clustering and eigenfunctions of Fokker-Planck operators. In Y. Weiss, B. Schölkopf, and J. Platt, editors, *Advances in Neural Information Processing Systems 18*, pages 955–962, Cambridge, MA, 2006. MIT Press.
- [NWF78] George L Nemhauser, Laurence A Wolsey, and Marshall L Fisher. An analysis of approximations for maximizing submodular set functions—i. *Mathematical programming*, 14(1):265–294, 1978.
- [PM13] D. Perraul-Joncas and M. Meila. Non-linear dimensionality reduction: Riemannian metric estimation and the problem of geometric discovery. *ArXiv e-prints*, May 2013.
- [Por16] Jacobus W Portegies. Embeddings of riemannian manifolds with heat kernels and eigenfunctions. *Communications on Pure and Applied Mathematics*, 69(3):478–518, 2016.
- [Str07] Walter A Strauss. *Partial differential equations: An introduction*. Wiley, 2007.
- [THJ10] Daniel Ting, Ling Huang, and Michael I. Jordan. An analysis of the convergence of graph laplacians. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 1079–1086, 2010.

## Supplement to Selecting the independent coordinates of manifolds with large aspect ratios

### A Notational table

Table S1: Notational table	
Matrix operation	
$\mathbf{M}$	Matrix
$\mathbf{m}_i$	Vector represents the $i$ -th row of $\mathbf{M}$
$\mathbf{m}_{:,j}^T$	Vector represents the $j$ -th column of $\mathbf{M}$
$m_{ij}$	Scalar represents $ij$ -th element of $\mathbf{M}$
$[\mathbf{M}]_{ij}$	Scalar, alternative notation for $m_{ij}$
$\mathbf{M}[\alpha, \beta]$	Submatrix of $\mathbf{M}$ of index sets $\alpha, \beta$
$\mathbf{v}$	Column vector
$v_i$	Scalar represents $i$ -th element of vector $\mathbf{v}$
$[\mathbf{v}]_i$	Scalar, alternative notation for $v_i$
Scalars	
$n$	Number of samples
$D$	Ambient dimension
$m$	Dimension of diffusion embedding
$s$	(Minimum) embedding dimension
$d$	Intrinsic dimension
Vectors & Matrices	
$\mathbf{X}$	Data matrix
$\mathbf{x}_i$	Point $i$ in ambient space
$\mathbf{Y}$	Diffusion coordinates
$\mathbf{y}_i$	Point $i$ in diffusion coordinates
$\phi_i$	The $i$ -th diffusion coordinate of all points
$\mathbf{K}$	Kernel (similarity) matrix
$\mathbf{L}$	Graph Laplacian
$\mathbf{H}(i)$	Dual metric at point $i$
$\mathbf{I}_k$	Identity matrix in $k$ dimension space
$\mathbf{1}_n$	All one vector $\in \mathbb{R}^n$
$\mathbf{1}_S$	$[\mathbf{1}_S]_i = 1$ if $i \in S$ 0 otherwise
Miscellaneous	
$G(V, E)$	Graph with vertex set $V$ and edge set $E$
$\mathcal{M}$	Data manifold
$\phi(\cdot)$	Embedding mapping
$\mathfrak{L}(S; \zeta)$	Utilities
$\mathfrak{R}$	Unpenalized utilities
$[s]$	Set $\{1, \dots, s\}$
$D(\cdot    \cdot)$	KL divergence
$\mathbf{D}$	Jacobian
$\mathfrak{D}(S, i)$	Leave-one-out regret of point $i$

### B Pseudocodes

---

#### Algorithm S1: DIFFMAP

---

**Input** : Data matrix  $\mathbf{X} \in \mathbb{R}^{n \times D}$ ,  
bandwidth  $\varepsilon$ , embedding  
dimension  $m$

- 1 Compute similarity matrix  $\mathbf{K}$  with  $K_{ij} = \begin{cases} \exp\left[-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\varepsilon^2}\right] & \text{if } \|x - y\| \leq 3\varepsilon \\ 0 & \text{otherwise} \end{cases}$
- 2  $\mathbf{L} \leftarrow \text{LAPLACIAN}(\mathbf{K}) \in \mathbb{R}^{n \times n}$  (Algorithm S2)
- 3 Compute eigenvectors of  $\mathbf{L}$  for smallest  $m + 1$  eigenvalues  
 $[\phi_0 \ \phi_1 \ \dots \ \phi_m] \in \mathbb{R}^{n \times (m+1)}$

**Return**:  $\Phi = [\phi_1 \ \dots \ \phi_m] \in \mathbb{R}^{n \times m}$  The  
*embedding coordinates* of  $\mathbf{x}_i$  are  
 $(\Phi_{i1}, \dots, \Phi_{im}) \in \mathbb{R}^m$

---



---

#### Algorithm S2: LAPLACIAN

---

**Input** : Symmetric similarity matrix  $\mathbf{K}$

- 1 Calculate the *degree* of node  $i$ ,  
 $[\mathbf{w}]_i = \sum_{j=1}^n K_{ij} \triangleright \text{Set } \mathbf{W} = \text{diag}(\mathbf{w})$
- 2  $\tilde{\mathbf{L}} = \mathbf{W}^{-1} \mathbf{K} \mathbf{W}^{-1}$
- 3  $[\tilde{\mathbf{w}}]_i \leftarrow \sum_{j=1}^n \tilde{L}_{ij} \triangleright \text{Set } \tilde{\mathbf{W}} = \text{diag}(\tilde{\mathbf{w}})$
- 4  $\mathbf{L} = \mathbf{I}_n - \tilde{\mathbf{W}}^{-1} \tilde{\mathbf{L}}$

**Return**: Renormalized graph Laplacian  $\mathbf{L}$

---

---

**Algorithm S3: LLRCOORDSEARCH**

---

**Input** : Embedding

$$\mathbf{Y} = [\phi_1, \dots, \phi_m] \in \mathbb{R}^{n \times m}$$

1 Set the leave-one-out validation error

$$\mathbf{r} = [1, \dots, 1] \in \mathbb{R}^m$$

2 **for**  $s = 2 \rightarrow m$  **do**3     Bandwidth of LLR:  $h \leftarrow$ 

$$\frac{1}{3} \cdot \text{MEDIAN}(\text{PAIRWISEDIST}(\phi_{[s-1]}))$$

4      $\hat{\phi}_s \leftarrow$ 

$$\text{LOCALLINEARREGRESSION}(\phi_s, \phi_{[s-1]}, h)$$

5      $r_s = \sqrt{\frac{\|\hat{\phi}_s - \phi_s\|^2}{\|\phi_s\|^2}}$ 6 **end**7  $S_* \leftarrow \text{ARGSORT}(\mathbf{r})$      $\triangleright$  Sort in descending order.**Return:** Sorted independent coordinates  $S_*$ 

---

## C Extra theorems

### C.1 Submodularity of the objective functions

**Theorem S1.** For a rank  $d$  tangent space matrix  $\mathbf{U} \in \mathbb{R}^{m \times d}$ , if any submatrix  $\mathbf{U}_S$ , with index set  $S \subseteq [m]$  and  $|S| = s \geq d$ , is rank  $d$ , we have  $\mathfrak{R}_1$  be a submodular set function.

**Proof.** W.L.O.G, set  $n = 1$ , with slightly abuse of notation, let  $\mathbf{U} = \mathbf{U}_{T \cup \{i\}} \in \mathbb{R}^{(|T|+1) \times d}$ . The matrix can be written in the following form

$$\mathbf{U} = \begin{bmatrix} \mathbf{T} \\ \mathbf{a} \end{bmatrix} = \begin{bmatrix} \mathbf{S} \\ \mathbf{V} \\ \mathbf{a} \end{bmatrix} \in \mathbb{R}^{(|T|+1) \times d}$$

With  $\mathbf{U}_S = \mathbf{S}$ ,  $\mathbf{U}_T = \mathbf{T}$  and  $\mathbf{U}_{\{i\}} = \mathbf{a}$  for set  $S \subseteq T \subseteq [m]$  and  $i \in [m] \setminus T$ . Here  $\mathbf{a} \in \mathbb{R}^{1 \times d}$ . By the definition of  $\mathfrak{R}_1$  in (2), one has (ignoring the constants)

$$\begin{aligned} \mathfrak{R}_1(S) &= \log \det(\mathbf{S}^\top \mathbf{S}) \\ \mathfrak{R}_1(T) &= \log \det(\mathbf{T}^\top \mathbf{T}) \\ \mathfrak{R}_1(S \cap \{i\}) &= \log \det \left( \begin{bmatrix} \mathbf{S} \\ \mathbf{a} \end{bmatrix}^\top \begin{bmatrix} \mathbf{S} \\ \mathbf{a} \end{bmatrix} \right) \\ \mathfrak{R}_1(T \cap \{i\}) &= \log \det(\mathbf{U}^\top \mathbf{U}) \end{aligned}$$

Denote  $\partial_i f(S) = f(S \cup \{i\}) - f(S)$  for some function  $f$ , we have

$$\begin{aligned} \partial_i \mathfrak{R}_1(S) &= \log \det(\mathbf{S}^\top \mathbf{S} + \mathbf{a}^\top \mathbf{a}) - \log \det(\mathbf{S}^\top \mathbf{S}) \\ \partial_i \mathfrak{R}_1(T) &= \log \det(\mathbf{T}^\top \mathbf{T} + \mathbf{a}^\top \mathbf{a}) - \log \det(\mathbf{T}^\top \mathbf{T}) \end{aligned}$$

The full rank of any submatrices guarantees the positive definiteness of  $\mathbf{S}^\top \mathbf{S}$ ,  $\mathbf{T}^\top \mathbf{T}$ , by matrix determinant lemma [Har98], we have

$$\det(\mathbf{S}^\top \mathbf{S} + \mathbf{a}^\top \mathbf{a}) = \det(\mathbf{S}^\top \mathbf{S}) (1 + \mathbf{a}(\mathbf{S}^\top \mathbf{S})^{-1} \mathbf{a}^\top)$$

Therefore

$$\partial_i \mathfrak{R}_1(S) = 1 + \mathbf{a}(\mathbf{S}^\top \mathbf{S})^{-1} \mathbf{a}^\top$$

Similar equation holds for set  $T$ . Therefore,

$$\partial_i \mathfrak{R}_1(S) - \partial_i \mathfrak{R}_1(T) = \log \frac{1 + \mathbf{a}(\mathbf{S}^\top \mathbf{S})^{-1} \mathbf{a}^\top}{1 + \mathbf{a}(\mathbf{T}^\top \mathbf{T})^{-1} \mathbf{a}^\top}$$

Because  $\mathbf{T}^\top \mathbf{T} \succeq \mathbf{S}^\top \mathbf{S}$ , we have  $(\mathbf{S}^\top \mathbf{S})^{-1} \succeq (\mathbf{T}^\top \mathbf{T})^{-1}$  [HHJ90], which implies  $\partial_i \mathfrak{R}_1(S) - \partial_i \mathfrak{R}_1(T) \geq 0$  for all  $S \subseteq T \subseteq [m]$  and  $i \in [m] \setminus T$ . This completes the proof.  $\blacksquare$

**Theorem S2.**  $\mathfrak{R}_2$  is a submodular set function.

**Proof.** W.L.O.G, set  $n, d = 1$ . With slightly abuse of notation, let  $\mathbf{u} \leftarrow \mathbf{u}_1(i)$  and  $\mathbf{u}_S \leftarrow \mathbf{u}_1^S(i)$ . For any set  $S \subseteq T \subseteq [m]$  and  $i \in [m] \setminus T$ , we have

$$\begin{aligned} \partial_i \mathfrak{R}_2(S) &= \mathfrak{R}_2(S \cap \{i\}) - \mathfrak{R}_2(S) = \log \frac{\sum_{k \in S} u_k^2 + u_i^2}{\sum_{k \in S} u_k^2} = \log \frac{\Sigma_S + u_i^2}{\Sigma_S} \\ \partial_i \mathfrak{R}_2(T) &= \mathfrak{R}_2(T \cap \{i\}) - \mathfrak{R}_2(T) = \log \frac{\sum_{k \in T} u_k^2 + u_i^2}{\sum_{k \in T} u_k^2} = \log \frac{\Sigma_S + \Sigma_{T \setminus S} + u_i^2}{\Sigma_S + \Sigma_{T \setminus S}} \end{aligned}$$

Where  $\Sigma_S = \sum_{k \in S} u_k^2$ . By definition, we have  $\Sigma_S, \Sigma_{T \setminus S}, u_i^2 \geq 0$ . Therefore,

$$\begin{aligned} \partial_i \mathfrak{R}_2(S) - \partial_i \mathfrak{R}_2(T) &= \log \frac{(\Sigma_S + u_i^2) \cdot (\Sigma_S + \Sigma_{T \setminus S})}{\Sigma_S \cdot (\Sigma_S + \Sigma_{T \setminus S} + u_i^2)} \\ &= \log \underbrace{\left[ \frac{\Sigma_S^2 + \Sigma_S (\Sigma_{T \setminus S} + u_i^2) + u_i^2 \Sigma_{T \setminus S}}{\Sigma_S^2 + \Sigma_S (\Sigma_{T \setminus S} + u_i^2)} \right]}_{\geq 1} \geq 0 \end{aligned}$$

Which completes the proof.  $\blacksquare$



## D Greedy search

---

**Algorithm S4:** GREEDYINDEIGENSEARCH

---

**Input :** Orthogonal basis  $\{\mathbf{U}(i)\}_{i=1}^n$ , eigenvalues  $\lambda$ , intrinsic dimension  $d$ , regularization parameter  $\zeta$

- 1 Solve  $S_* \leftarrow \operatorname{argmax}_{S \subseteq [m]; |S|=d; 1 \in S} \mathcal{L}(S; \zeta)$ .
- 2 **for**  $s = d + 1 \rightarrow m$  **do**
- 3      $k_* = \operatorname{argmax}_{k \in [m] \setminus S_*} \mathcal{L}(S_* \cup \{k\}; \zeta)$
- 4      $S_* \leftarrow S_* \cup \{k_*\}$  ▷ Record order
- 5 **end**

**Return:** Independent coordinates  $S_*$

---

Inspired by the greedy version of submodular maximization [NWF78], a greedy heuristic has been proposed, as in Algorithm S4. The algorithm starts from an observation that the optimal value of the  $S' = \operatorname{argmax}_{S; d \leq |S| \leq s} \mathcal{L}(S; \zeta)$  will often time be a subset of the optimal  $S_*$  of (3). Since the appropriate cardinality of the set  $S$  is unknown, we can simply scan from  $|S| = d$  to  $m$ . The order of the returned elements indicates the significance of the corresponding coordinate.

## E Computational complexity analysis

### E.1 The proposed algorithms

For computation complexity analysis, we assume the embedding has already been obtained. Therefore, the computational complexity for building neighbor graph and solving the eigen-problem of graph Laplacian can be omitted. This is also the case for LLRCOORDSEARCH.

**Co-metrics and orthogonal basis** According to [PM13], time complexity for computing  $\mathbf{H}(i) \in \mathbb{R}^{m \times m} \forall i \in [n]$  is  $\mathcal{O}(nm^2\delta)$ , with  $\delta$  be the average degree of the neighbor graph  $G(V, E)$ . In manifold learning, the graph will be sparse therefore  $\delta \ll n$ . Time complexity for obtaining principal space  $\mathbf{U}(i)$  of point  $i$  via SVD will be  $\mathcal{O}(m^3)$ . Total time complexity will be  $\mathcal{O}(nm^2\delta + nm^3)$ .

**Exact search** Evaluating the objective function  $\mathcal{L}$  for each point  $i$  takes  $\mathcal{O}(sd^2)$  in computing  $\mathbf{U}_S(i)^\top \mathbf{U}_S(i)$ ,  $\mathcal{O}(d^3)$  in evaluating the determinant of a  $d \times d$  matrix. Normalization ( $\mathfrak{R}_2$  term) takes  $\mathcal{O}(ds)$ . Exhaustive search over all the subset with cardinality  $s$  takes  $\mathcal{O}\left(\binom{m}{s}\right)$ . The total computational complexity will therefore be  $\mathcal{O}(nm^s(d^3 + d^2s) + nm^2\delta + nm^3) = \mathcal{O}(nm^{s+3} + nm^2\delta)$ .

**Greedy algorithm** First step of greedy algorithm includes solving  $\operatorname{argmax}_{S \subseteq [m]; |S|=d} \mathcal{L}(S, d)$ , which takes  $\mathcal{O}(nm^d d^3) = \mathcal{O}(nm^{d+3})$ . Starting from  $s = d + 1 \rightarrow m$ , each step includes exhaustively search over  $m - s$  candidates, with the time complexity of evaluating  $\mathcal{L}$  be  $n(d^3 + d^2s)$ . Putting things together, one has the second part of the greedy algorithm be

$$\sum_{s=d}^m n(m-s)(d^3 + d^2s) = \mathcal{O}(nm^5) \quad (\text{S1})$$

The total computational complexity will therefore be  $\mathcal{O}(n(m^{d+3} + m^5 + m^2\delta))$ .

### E.2 Time complexity of [DTCK18] & discussion

The Algorithm LLRCOORDSEARCH is summarized in Algorithm S3. For searching over fixed coordinate  $s$ , the algorithm first build a kernel for local linear regression by constructing a neighbor graph, which takes  $\mathcal{O}(n \log(n)s)^3$  using approximate nearest neighbor search. The  $s$  dependency come from the dimension of the feature. For each point  $i$ , a ordinary least square (OLS) problem is solved, which results in  $\mathcal{O}(n^2s^2 + ns^3)$  time complexity.

<sup>3</sup>This is a simplified lower bound, see [DS13] for details.

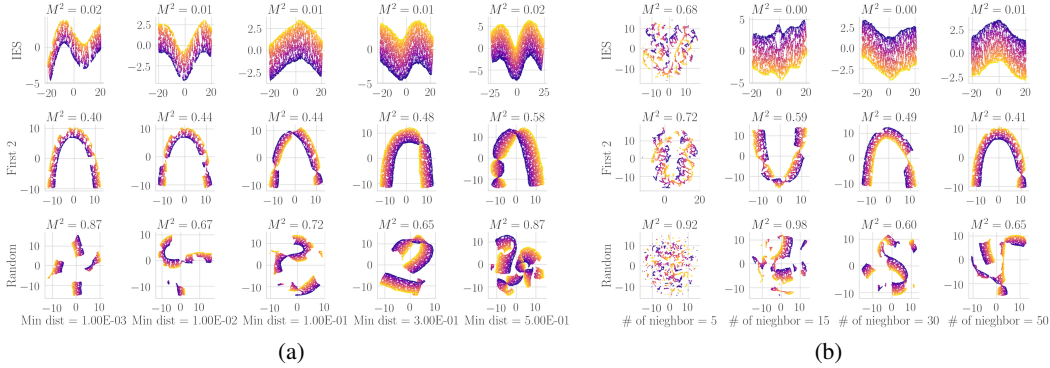


Figure S2: UMAP embeddings of 2D long stripe with different initializations and choices of hyper-parameters. Rows from top to bottom correspond to UMAP embedding initialized with DM which coordinates chosen by INDEIGENSEARCH, naïve DM and random initialization, respectively. Columns represent different choices of (a) points separation and (b) number of neighbors.

Searching from  $s = 2 \rightarrow m$  will make the total time complexity be

$$\sum_{s=2}^m n^2 s^2 + ns^3 + ns \log n = \mathcal{O}(n^2 m^3 + nm^4) \quad (\text{S2})$$

For a sparse graph, the overheads of the INDEIGENSEARCH and GREEDYINDEIGENSEARCH algorithms come from the enumeration of the subset  $S$ . Because of the linear dependency on the sample size  $n$ , the algorithm is tractable for small  $s$  and  $d$ . However, LLRCOORDSEARCH has a quadratic dependency on sample size  $n$ , which is more computationally intensive for large sample size. For large  $s$  and  $d$ , one can use the techniques in difference between submodular function optimization (e.g. [IB12]) as  $\mathfrak{R}_1, \mathfrak{R}_2$  are both submodular set function from Theorems S1 and S2. An empirical runtime plot for different algorithms can be found in Figure S1. The runtime was evaluated on two dimensional long strip with  $s = d = 2$  and was performed on a single desktop computer running Linux with 32GB RAM and a 8-Core 4.20GHz Intel® Core™ i7-7700K CPU.

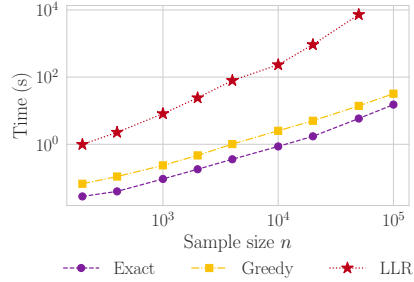


Figure S1: Runtimes of different IES algorithms on two dimensional long strip. Purple, yellow and red curves correspond to INDEIGENSEARCH, GREEDYINDEIGENSEARCH and LLRCOORDSEARCH algorithm, respectively.

## F A discussion on UMAP

UMAP [MHM18] is a commonly used data visualization alternative of t-SNE. The authors proposed to use the spectral embedding of the graph Laplacian as an initialization to the algorithm for faster convergence (compared to random initialization). In this section, we showed empirically that, (1) given reasonable computing resources, the IES problem also appears in the UMAP embedding and (2) by initializing with spectral embedding with carefully selected coordinate set chosen by INDEIGENSEARCH, one can obtain a faster convergence and a globally interpretable embedding. Figure S2 is the UMAP embedding of 2D long stripe dataset  $\mathcal{D}_1$  with different choices of hyper-parameters (points separation in S2a and number of neighbors in S2b), with total number of epochs be 500. The first row of both plots are the embedding initialized with INDEIGENSEARCH, the second row corresponds to those initialized with naïve DM. The embeddings in the third row are initialized randomly. As shown in the results, algorithmic/random artifacts can be easily seen in the embeddings with Naïve DM or random initialization (2nd and 3rd rows). More precisely, unwanted patterns which reduce the interpretability of the embeddings, e.g., the “knots” in the second row or disconnected components in the second/third rows, are generated. The sum of square procrustes

error  $M^2$  between ground truth dataset and the embedding shown on each subplots also confirm our statement. Note that it is possible to unroll the algorithmic artifact with more epochs in the sampling steps of UMAP. (3 to 5 times more iterations are needed in this example.) However, due to the efficiency of performing INDEIGENSEARCH, it is beneficial to initialize with the embedding selected by INDEIGENSEARCH.

## G Additional experiments & details of the used datasets

In this paper, a total of 13 different synthetic manifolds are considered. Table S2 summarized the synthetic manifolds constructed and its abbreviations (from  $\mathcal{D}_1$  to  $\mathcal{D}_{13}$ ). Embedding results for the synthetic manifolds are in Figures S3, S4 and S5. The ranking of the first few candidate sets  $S$  from INDEIGENSEARCH, GREEDYINDEIGENSEARCH and LLRCOORDSEARCH can be found in Table S3. The table shows the optimal subsets return by three different algorithms are often time the same, with exception for  $\mathcal{D}_7$  *high torus* as discussed in Section 7.

Table S2: Abbreviations for different synthetic manifolds in this paper. The abbreviation with asterisk represents such dataset is discussed in main manuscript.

Manifold with $s = d$	
$\mathcal{D}_1^*$	Two dimensional strip (aspect ratio $2\pi$ )
$\mathcal{D}_2$	2D strip with cavity (aspect ratio $2\pi$ )
$\mathcal{D}_3$	Swiss roll
$\mathcal{D}_4$	Swiss roll with cavity
$\mathcal{D}_5$	Gaussian manifold
$\mathcal{D}_6$	Three dimensional cube
Manifold with $s > d$	
$\mathcal{D}_7^*$	High torus
$\mathcal{D}_8$	Wide torus
$\mathcal{D}_9$	z-asymmetrized high torus
$\mathcal{D}_{10}$	x-asymmetrized high torus
$\mathcal{D}_{11}$	z-asymmetrized wide torus
$\mathcal{D}_{12}$	x-asymmetrized wide torus
$\mathcal{D}_{13}^*$	Three-torus

### G.1 Additional experiments on synthetic manifolds with $s = d$

Below summarized the details of generating the datasets.

Table S3: Results returned from different algorithms on different synthetic datasets.

	Exact search					Greedy rank	LLR rank
	1	2	3	4	5		
$\mathcal{D}_1$	[1, 7]	[1, 8]	[1, 9]	[1, 10]	[1, 12]	[1, 7, 6, 4, 3, 2, 5]	[1, 7, 14, 16, 11, 18, 6]
$\mathcal{D}_2$	[1, 4]	[1, 8]	[1, 9]	[1, 10]	[1, 12]	[1, 4, 8, 6, 5, 3, 2]	[1, 4, 8, 5, 17, 11, 14]
$\mathcal{D}_3$	[1, 9]	[1, 10]	[1, 11]	[1, 13]	[1, 18]	[1, 9, 5, 2, 3, 4, 6]	[1, 9, 19, 16, 12, 10, 4]
$\mathcal{D}_4$	[1, 8]	[1, 10]	[1, 11]	[1, 14]	[1, 15]	[1, 8, 3, 2, 4, 10, 5]	[1, 8, 11, 10, 19, 16, 4]
$\mathcal{D}_5$	[1, 6]	[1, 8]	[1, 10]	[1, 11]	[1, 13]	[1, 6, 2, 8, 3, 10, 4]	[1, 6, 19, 8, 18, 14, 12]
$\mathcal{D}_6$	[1, 2, 8]	[1, 2, 11]	[1, 4, 8]	[1, 2, 17]	[1, 2, 13]	[1, 2, 8, 3, 4, 6, 5]	[1, 2, 8, 10, 3, 13, 6]
$\mathcal{D}_7$	[1, 4, 5]	[1, 4, 8]	[1, 5, 7]	[1, 7, 12]	[1, 7, 8]	[1, 5, 4, 3, 6, 2, 8]	[1, 2, 5, 4, 15, 6, 10]
$\mathcal{D}_8$	[1, 2, 7]	[1, 4, 7]	[1, 3, 7]	[1, 2, 9]	[1, 5, 7]	[1, 7, 2, 4, 3, 13, 5]	[1, 2, 7, 13, 12, 15, 14]
$\mathcal{D}_9$	[1, 3, 4]	[1, 3, 7]	[1, 4, 6]	[1, 3, 10]	[1, 7, 9]	[1, 3, 4, 2, 9, 7, 6]	[1, 3, 4, 2, 19, 8, 7]
$\mathcal{D}_{10}$	[1, 2, 4]	[1, 3, 4]	[1, 4, 5]	[1, 6, 9]	[1, 6, 14]	[1, 4, 2, 3, 5, 6, 8]	[1, 4, 2, 3, 8, 5, 6]
$\mathcal{D}_{11}$	[1, 2, 5]	[1, 4, 8]	[1, 4, 5]	[1, 8, 9]	[1, 2, 8]	[1, 5, 2, 4, 8, 3, 9]	[1, 2, 5, 8, 10, 9, 11]
$\mathcal{D}_{12}$	[1, 2, 5]	[1, 4, 5]	[1, 2, 7]	[1, 3, 5]	[1, 2, 8]	[1, 5, 2, 3, 4, 6, 8]	[1, 5, 2, 6, 10, 9, 4]
$\mathcal{D}_{13}$	[1, 2, 5, 10]	[1, 3, 5, 10]	[1, 4, 5, 10]	[1, 5, 6, 10]	[1, 2, 8, 10]	[1, 5, 10, 2, 4, 3, 6]	[1, 2, 10, 5, 14, 15, 16]

1.  $\mathcal{D}_1^*$ : points from this dataset are sampled uniformly from  $\mathbf{x}_i \sim \text{UNIF}([-2, 2] \times [-4\pi, 4\pi])$ .
2.  $\mathcal{D}_2$ : points are first sampled uniformly from  $[-2, 2] \times [-4\pi, 4\pi]$ . Points  $i$  are removed if  $|X_{i1}| < 4\pi/3$  and  $|X_{i2}| < 2/3$ .
3.  $\mathcal{D}_3$ : first sampling points  $\mathbf{X}_{\text{true}} = [\mathbf{x}_0, \mathbf{y}_0]$  uniformly from a two dimensional strip. The data  $\mathbf{X}$  can be obtained by the following non-linear transformation.

$$\mathbf{X} = \left[ \frac{\mathbf{x}_0 \circ \cos \mathbf{x}_0}{2}, \mathbf{y}_0, \frac{\mathbf{x}_0 \circ \sin \mathbf{x}_0}{2} \right] \quad (\text{S3})$$

With  $\circ$  denotes Hadamard (element-wise) product.

4.  $\mathcal{D}_4$ : sampling points  $\mathbf{X}_{\text{true}} = [\mathbf{x}_0, \mathbf{y}_0]$  uniformly from 2D strip with cavity then applying the transformation (S3) to get  $\mathbf{X}$ .
5.  $\mathcal{D}_5$ : sampling points  $\mathbf{X}_{\text{true}}$  uniformly from ellipse  $\left\{ (x, y) \in \mathbb{R}^2 : \left(\frac{x}{6}\right)^2 + \left(\frac{y}{2}\right)^2 = 1 \right\}$ . The data is obtained by

$$\mathbf{X} = [\mathbf{X}_{\text{true}}, \mathbf{z}]$$

$$\text{With } z_i = \exp\left(-\left(\left(\frac{X_{i1}}{3}\right)^2 + X_{i2}^2\right)/2\right)$$

6.  $\mathcal{D}_6$ : points are sampled uniformly from  $[-1, 1] \times [-2, 2] \times [-4, 4]$ .

The experimental results are in Figure S3 ( $\mathcal{D}_4$  in Figure 2).

## G.2 Additional experiments on synthetic manifolds with $s > d$

### G.2.1 Tori and asymmetrized tori

A torus can be parametrized by

$$\begin{aligned} x &= (a + b \cos \alpha) \cos \beta \\ y &= (a + b \cos \alpha) \sin \beta \\ z &= h \sin(\beta) \end{aligned} \quad (\text{S4})$$

1.  $\mathcal{D}_7^*$ : sampling  $\alpha, \beta$  uniformly from  $[0, 2\pi)$  and generating the torus with  $(a, b, h) = (3, 2, 8)$  from (S4).
2.  $\mathcal{D}_8$ : generating the torus with  $(a, b, h) = (10, 2, 2)$ .
3.  $\mathcal{D}_9$ : generating a high torus with  $(a, b, h) = (3, 2, 8)$  and applying the following transformation

$$z \leftarrow (z - \min(z))^{\gamma/\varsigma} \quad (\text{S5})$$

$$\text{with } (\gamma, \varsigma) = (3, 1500)$$

4.  $\mathcal{D}_{10}$ : generating a high torus with  $(a, b, h) = (3, 2, 8)$  and applying the following transformation

$$x \leftarrow (x - \min(x))^{\kappa/\eta} \quad (\text{S6})$$

$$\text{with } (\kappa, \eta) = (2, 10)$$

5.  $\mathcal{D}_{11}$ : generating a wide torus with  $(a, b, h) = (10, 2, 2)$  and applying transformation (S5) with  $(\gamma, \varsigma) = (3, 50)$ .
6.  $\mathcal{D}_{12}$ : generating a wide torus with  $(a, b, h) = (10, 2, 2)$  and applying transformation (S6) with  $(\kappa, \eta) = (3, 1000)$ .

The experimental results are in Figure S4.

### G.2.2 Three-torus

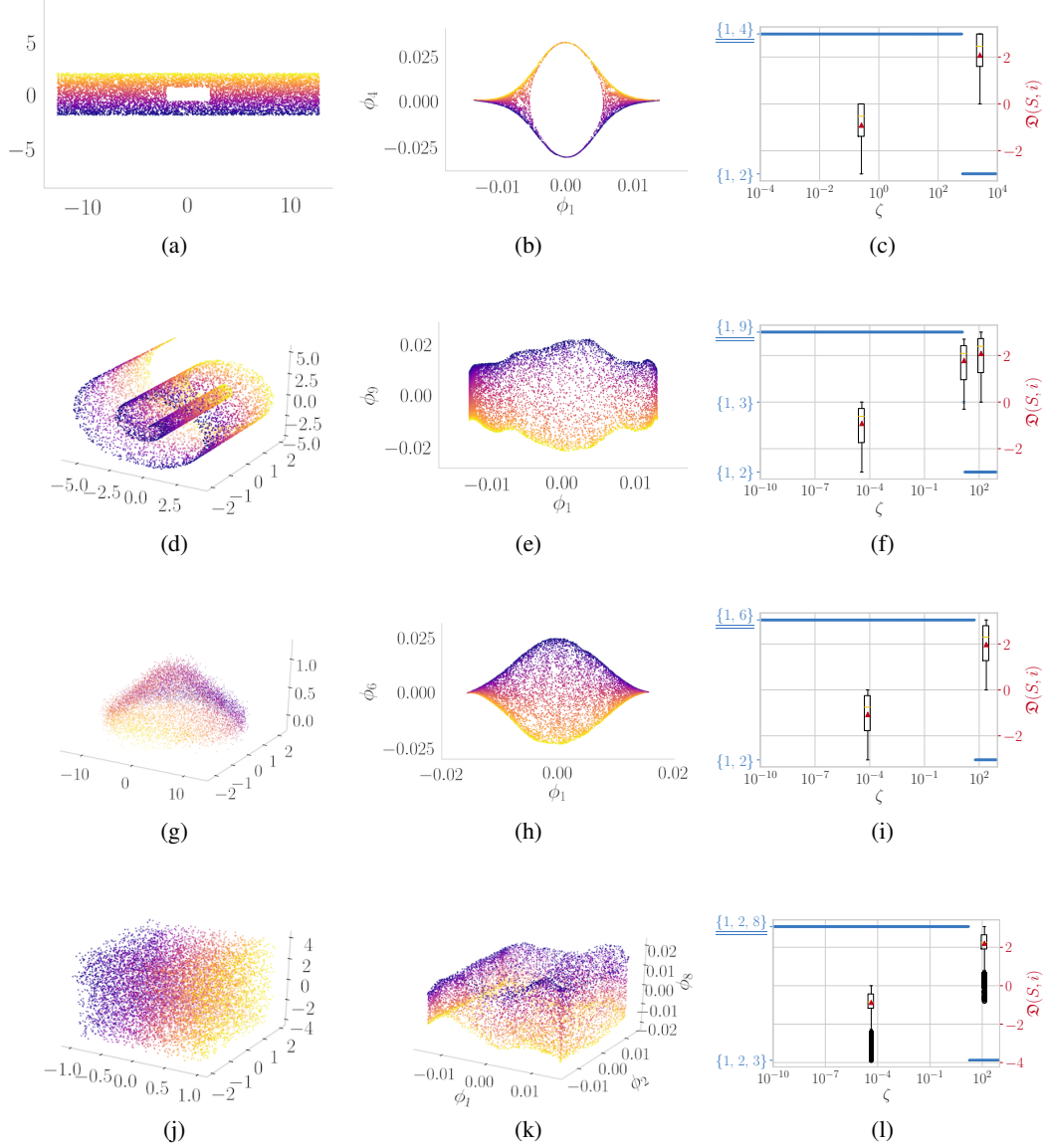


Figure S3: Synthetic manifolds with minimum embedding dimension  $s$  equals intrinsic dimension  $d$ . Rows from top to bottom represent *two dimensional strip with cavity* (aspect ratio  $W/H = 2\pi$ ), *swiss roll*, *gaussian manifold* and *three dimensional cube* dataset, respectively. Columns from left to right are the original data  $\mathbf{X}$ , embedding  $\phi_{S_*}$  with optimal coordinate sets  $S_*$  chosen by INDEIGENSEARCH and the regularization path, respectively.

The parameterization of the three torus is

$$\begin{aligned}
 x_1 &= a_1 \sin \alpha_1 \\
 x_2 &= (a_2 + a_1 \cos \alpha_1) \sin \alpha_2 \\
 x_3 &= (a_3 + (a_2 + a_1 \cos \alpha_1) \cos \alpha_2) \sin \alpha_3 \\
 x_4 &= (a_3 + (a_2 + a_1 \cos \alpha_1) \cos \alpha_2) \cos \alpha_3
 \end{aligned}
 \tag{S7}$$

To generate  $\mathcal{D}_{13}$ , we sample  $\alpha_k$  uniformly from  $[0, 2\pi)$  for  $k \in [3]$  and apply the transformation (S7) with  $(a_1, a_2, a_3) = (8, 2, 1)$ . The sample size for this dataset is  $n = 50,000$ . The experimental result of three-torus can be found in Figure S5.

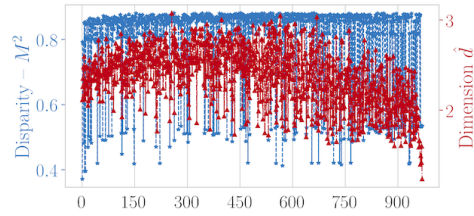


Figure S6:  $M^2$  and  $\hat{d}$  vs. ranking of  $\mathcal{D}_{13}$

### G.3 Verification of the chosen subsets on synthetic manifolds

Unlike 2D strip, the close form solution of the optimal set is oftentimes unknown in general. In this section, we verify the correctness of the chosen subset by reporting the full procrustes distance (disparity score)  $M^2$  [Dry16], which is defined to be the normalized sum of square of the point-wise difference between the procrustes transformed ground truth data  $\mathbf{X}_{\text{true}} \in \mathbb{R}^{n \times k}$  and the test data  $\mathbf{X}_{\text{test}} \in \mathbb{R}^{n \times k}$ . Namely,

$$\begin{aligned} M^2(\mathbf{X}_{\text{true}}, \mathbf{X}_{\text{test}}) &= \min_{\beta, \gamma, \Gamma} \|\mathbf{X}_{\text{true}} - \beta \mathbf{X}_{\text{test}} \Gamma - \mathbf{1}_n \gamma^\top\|_F^2 \\ \text{s.t. } &\beta > 0, \gamma \in \mathbb{R}^k, \Gamma \in SO(k) \end{aligned} \quad (\text{S8})$$

Here  $\beta$  is a scale parameter,  $\gamma$  is the centering parameter and  $\Gamma$  is a  $k \times k$  rotation matrix. We further require  $\|\mathbf{X}_{\text{true}}\|_F = 1$  so that the disparity score will be between 0 and 1. Intuitively, one can expect the optimal choice of eigencoordinates  $S_*$  will yield a small disparity score  $M^2(\mathbf{X}_{\text{true}}, \phi_{S_*})$ , with score increases as the coordinate set  $S$  contains duplicate parameterizations or  $\phi_S$  contains *knots*, *crossings*, etc. (e.g., Figure S8). Note that the score can only be calculated when the ground truth data  $\mathbf{X}_{\text{true}}$  is available. For dataset without obtainable ground truth, one cannot proposed to report

the disparity score of  $\phi_S$  and the original data  $\mathbf{X}$  as the proxy of  $\mathbf{X}_{\text{true}}$ , for  $\mathbf{X}$  might not be a affine transformation of  $\mathbf{X}_{\text{true}}$ , e.g., Swiss roll. Besides, small  $M^2$  given  $\phi_S$  does not imply  $S$  is optimal, which will be clear in the discussion of Figure 2g. Besides disparity scores, we will also report the estimated dimension  $\hat{d}$ . One can expect the estimated dimension for the optimal set  $\dim(\phi_{S_*})$  will be close to the intrinsic dimension  $d$ , while the estimated dimension for sets containing duplicate parameterizations will be smaller than the intrinsic dimension. One cannot propose to use it as a criterion to choose the optimal set, for the suboptimal sets can also have estimated dimensions closed to the intrinsic dimension, e.g., Figure 4g. Throughout the experiment, the dimension estimation method by [LB05] is used for its ability to estimate dimension among all candidate subsets fairly fast. Blue and red curves in Figure S7 and S6 show the disparity scores and estimated dimensions versus ranking of coordinate subsets for different synthetic manifolds, respectively. As expected, we have an increasing in  $M^2$  and decreasing in  $\hat{d}$  with respect to ranking. We first highlight that the set that produces the lowest disparity score is not necessarily optimal, although  $S_*$  does yield a small disparity. This can be shown in the example of  $\mathcal{D}_4$  *swiss roll with hole* dataset. Figure 2g is the embedding  $\phi_{S_3}$  of  $\mathcal{D}_4$ , with  $S_3$  is ranked third subset in terms of  $\mathcal{L}(S; \zeta)$ , that minimizes the disparity score  $M^2$  in  $\mathcal{D}_4$  as shown in Figure S7d. This is because the embedding of the subset  $S_3 = \{1, 11\}$  has larger area on the left, compared to Figure Figure 2b. This balances out the high disparity caused by the *flipped* region between two *knots* in the embedding  $\phi_{S_3}$  when matched with  $\mathbf{X}_{\text{true}}$ . Since all the ranked first subset has low disparity compared to other subsets, we have higher confidence saying that the ranked 1st subset is indeed the optimal choice for the synthetic manifolds.

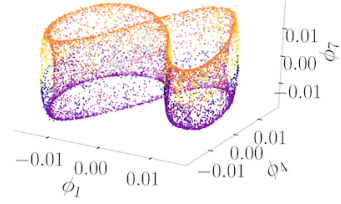


Figure S8: Embedding that has crossing.

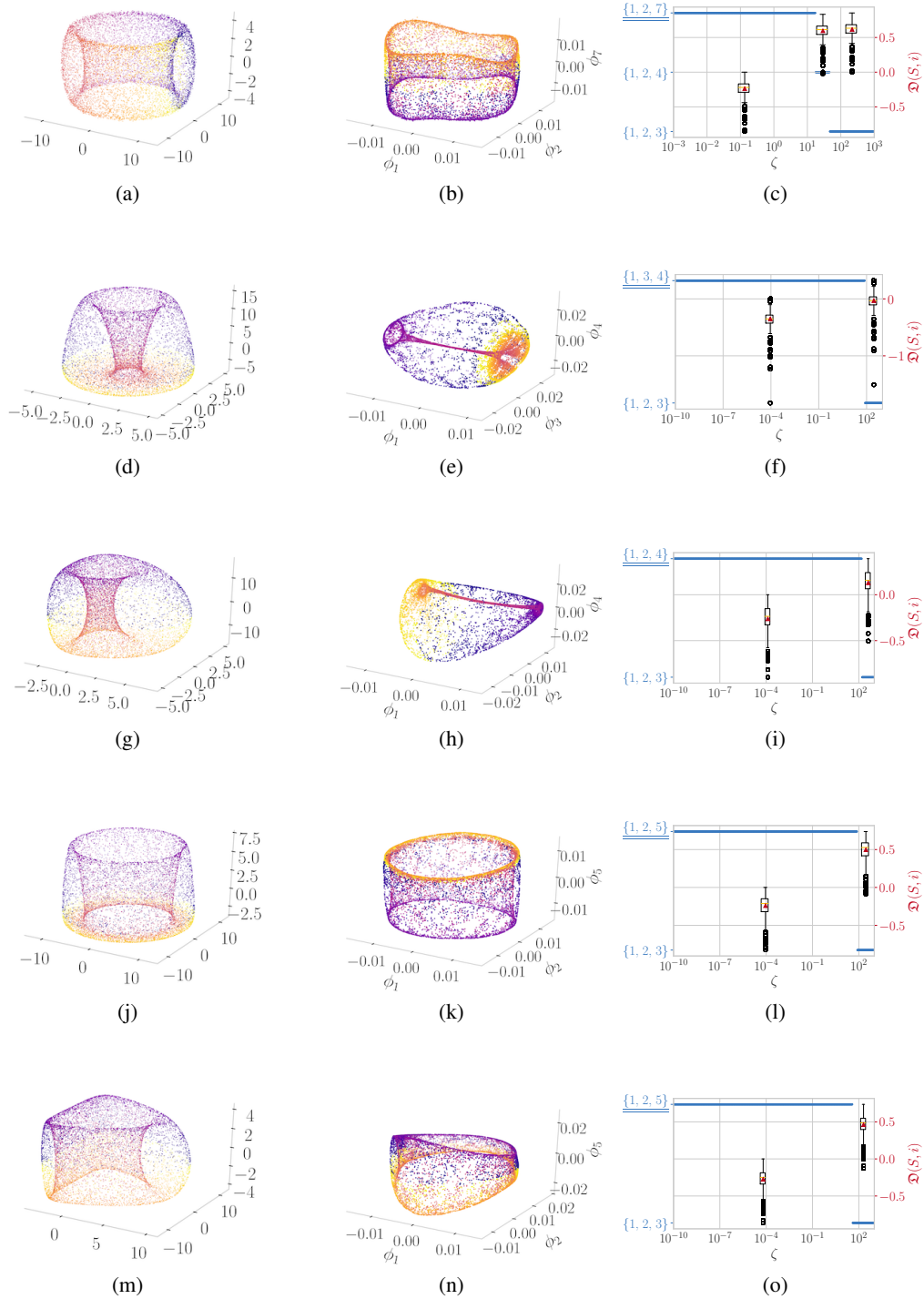
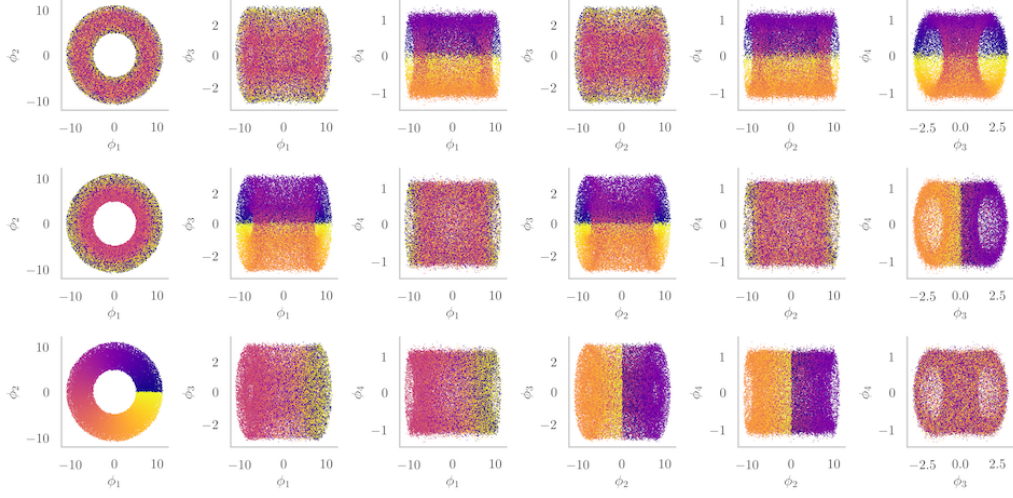
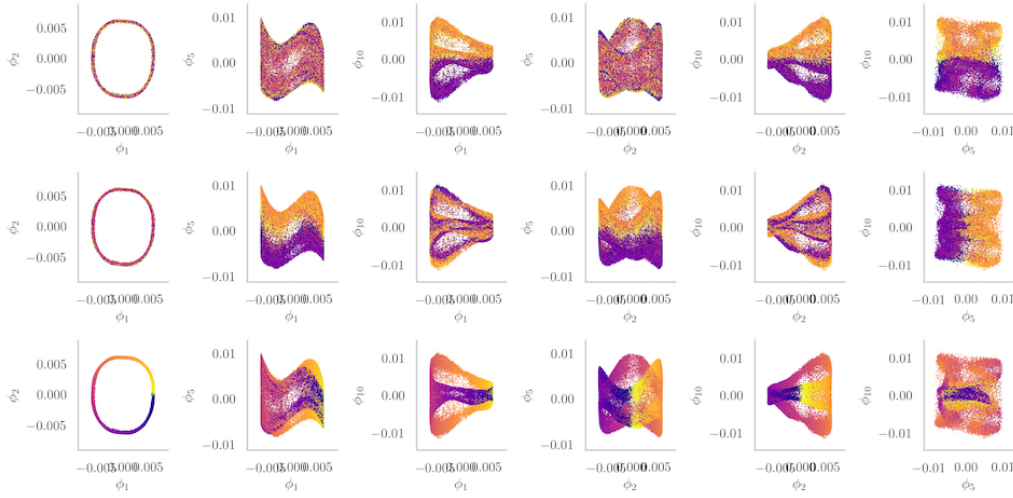


Figure S4: Synthetic manifolds with minimum embedding dimension  $s$  greater than intrinsic dimension  $d$ . Rows from top to bottom represent *wide torus*, *z-asymmetrized high torus*, *x-asymmetrized high torus*, *z-asymmetrized wide torus* and *x-asymmetrized wide torus*, respectively. Columns from left to right are the original data  $X$ , embedding  $\phi_{S_*}$  with optimal coordinate sets  $S_*$  chosen by INDEIGENSEARCH and the regularization path, respectively.



(a)



(b)

Figure S5: Experiment on *three-torus* dataset. (a) Original data  $\mathbf{X}$  of three torus. (b) Embedding  $\phi_{S_*}$  with optimal coordinate sets  $S_*$  chosen by INDEIGENSEARCH. Rows for both (a) and (b) from top to bottom are embedding colored by the parameterization  $(\alpha_1, \alpha_2, \alpha_3)$  in (S7), respectively.



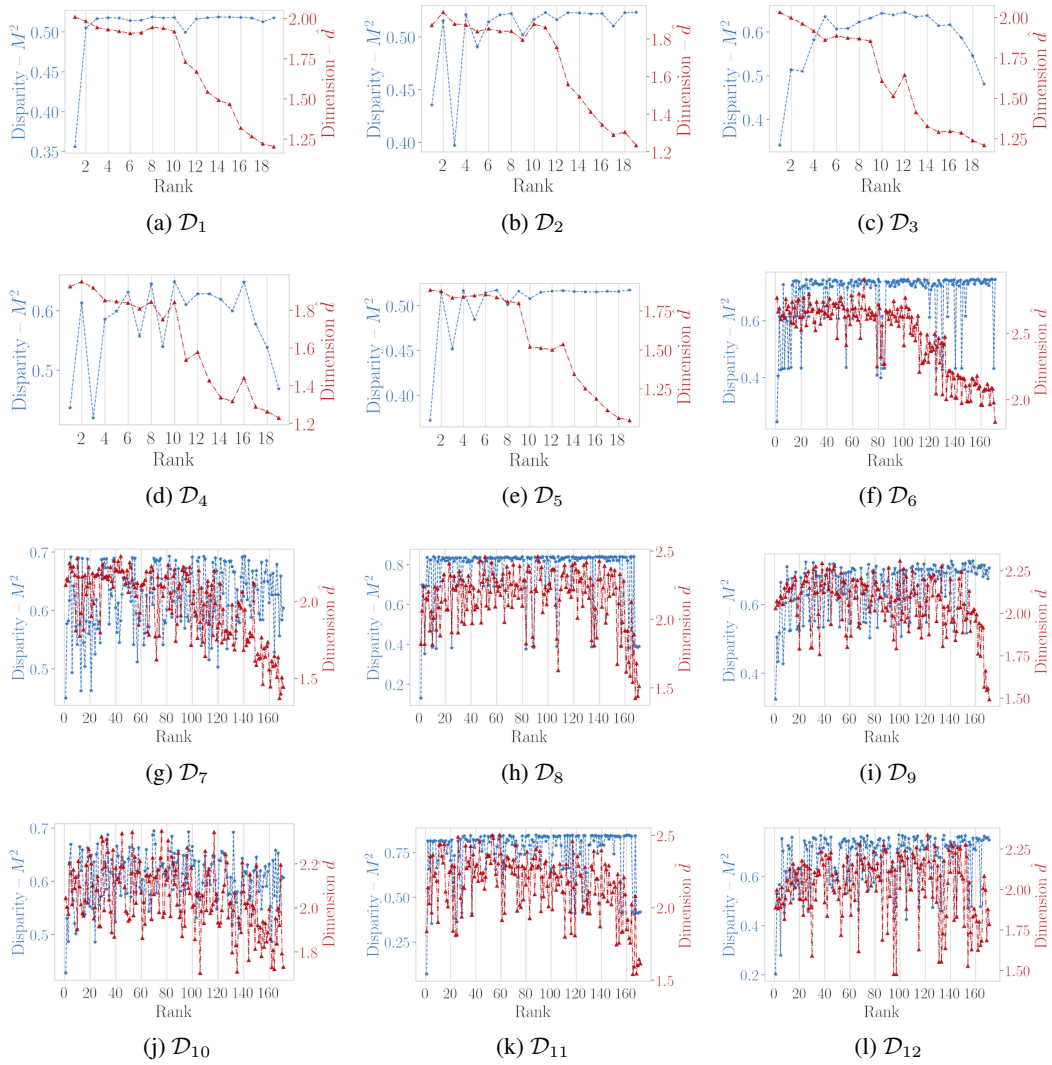


Figure S7: Verification of the correctness of the chosen sets in synthetic manifolds.