

Lecture 10: Dimension Reduction Techniques

Radu Balan

Department of Mathematics, AMSC, CSCAMM and NWC
University of Maryland, College Park, MD

April 17, 2018

Input Data

It is assumed that there is a set of points $\{x_1, \dots, x_n\} \subset \mathbb{R}^N$, however either partial, or different information is available:

- 1 Geometric Graph: For a threshold $\tau \geq 0$, $\mathcal{G}_\tau = (\mathcal{V}, \mathcal{E}, \mu)$ where \mathcal{V} is the set of n vertices (nodes), \mathcal{E} is the set of edges between nodes i and j if $\|x_i - x_j\| \leq \tau$ and $\mu : \mathcal{E} \rightarrow \mathbb{R}$ the set of distances $\|x_i - x_j\|$ between nodes connected by an edge.
- 2 Weighted graph: $\mathcal{G} = (\mathcal{V}, W)$ a undirected weighted graph with n nodes and weight matrix W , where $W_{i,j}$ is inverse monotonically dependent to distances $\|x_i - x_j\|$; the smaller the distance $\|x_i - x_j\|$ the larger the weight $W_{i,j}$.
- 3 Unweighted graph: For a threshold $\tau \geq 0$, $\mathcal{U}_\tau = (\mathcal{V}, \mathcal{E})$ where \mathcal{V} is the set of n nodes, and \mathcal{E} is the set of edges connected node i to node j if $\|x_i - x_j\| \leq \tau$. Note the distance information is not available.

Thus we look for a dimension $d > 0$ and a set of points

$\{y_1, y_2, \dots, y_n\} \subset \mathbb{R}^d$ so that all $d_{i,j} = \|y_i - y_j\|$'s are compatible with raw data as defined above.

Approaches

Popular Approaches:

- 1 Principal Component Analysis
- 2 Independent Component Analysis
- 3 Laplacian Eigenmaps
- 4 Local Linear Embeddings (LLE)
- 5 Isomaps

If points were supposed to belong to a lower dimensional manifold, the problem is known under the term *manifold learning*. If the manifold is linear (affine), then the Principal Component Analysis (PCA) or Independent Component Analysis (ICA) would suffice. However, if the manifold is not linear, then nonlinear methods are needed. In this respect, Laplacian Eigenmaps, LLE and ISOMAP can be thought of as *nonlinear PCA*. Also known as *nonlinear embeddings*.

Principal Component Analysis

Approach

Data: We are given a set $\{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^N$ of n points in \mathbb{R}^N .

Goal: We want to find a linear (or affine) subspace V of dimension d that best approximates this set. Specifically, if $P = P_V$ denotes the orthogonal projection onto V , then the goal is to minimize

$$J(V) = \sum_{k=1}^n \|x_k - P_V x_k\|_2^2.$$

If V is linear space (i.e. passes through the origin) then P is $N \times N$ linear operator (i.e. matrix) that satisfies $P = P^T$, $P^2 = P$, and $\text{Ran}(P) = V$. If V is an affine space (i.e. a linear space shifted by a constant vector), then the projection onto the affine space is $T(x) = Px + b$ where b is a constant vector (the "shift").

The affine space case can be easily reduced to the linear space: just append 1 to the bottom of each vector x_k : $\tilde{x}_k = [x_k; 1]$. Now b becomes a column of the extended matrix $\tilde{P} = [P \ b]$.

Principal Component Analysis

Algorithm

Algorithm (Principal Component Analysis)

Input: Data vectors $\{x_1, \dots, x_n\} \in \mathbb{R}^N$; dimension d .

- ① *If affine subspace is the goal, append '1' at the end of each data vector.*
- ① *Compute the sample covariance matrix*

$$R = \sum_{k=1}^n x_k x_k^T$$

- ② *Solve the eigenproblems $R e_k = \sigma_k^2 e_k$, $1 \leq k \leq N$, order eigenvalues $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_N^2$ and normalize the eigenvectors $\|e_k\|_2 = 1$.*

Principal Component Analysis

Algorithm - cont'ed

Algorithm (Principal Component Analysis)

- ③ *Construct the co-isometry*

$$U = \begin{bmatrix} e_1^T \\ \vdots \\ e_d^T \end{bmatrix}.$$

- ④ *Project the input data*

$$y_1 = Ux_1, \quad y_2 = Ux_2, \quad \dots, \quad y_n = Ux_n.$$

Output: Lower dimensional data vectors $\{y_1, \dots, y_n\} \in \mathbb{R}^d$.

The orthogonal projection is given by $P = \sum_{k=1}^d e_k e_k^T$ and the optimal subspace is $V = \text{Ran}(P)$.

Principal Component Analysis

Derivation

Here is the derivation in the case of linear space. The reduced dimensional data is given by Px_k . Expand the criterion $J(V)$:

$$J(V) = \sum_{k=1}^n \|x_k\|^2 - \sum_{k=1}^n \langle Px_k, x_k \rangle = \sum_{k=1}^n \|x_k\|^2 - \text{trace}(PR)$$

where $R = \sum_{k=1}^n x_k x_k^T$. It follows the minimizer of $J(V)$ maximizes $\text{trace}(PR)$ subject to $P = P^T$, $P^2 = P$ and $\text{trace}(P) = d$. It follows the optimal P is given by the orthogonal projection onto the top d eigenvectors, hence the algorithm.

Independent Component Analysis

Approach

Model (Setup): $x = As$, where A is an unknown invertible $N \times N$ matrix, and $s \in \mathbb{R}^N$ is a random vector of *independent* components.

Data: We are given a set of measurement $\{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^N$ of n points in \mathbb{R}^N of the model $x_k = As_k$, where each $\{s_1, \dots, s_n\}$ is drawn from the same distribution $p_s(s)$ of N -vectors with independent components.

Goal: We want to estimate the invertible matrix A and the (source) signals $\{s_1, \dots, s_n\}$. Specifically, we want a square matrix W such that Wx has independent components.

Principle: Perform PCA first so the decorrelated signals have unit variance. Then find an orthogonal matrix (that is guaranteed to preserve decorrelation) that creates statistical independence as much as possible.

Caveat: Two inherent ambiguities: (1) Permutation: If W is a solution to the unmixing problem so is ΠW , where Π is a permutation matrix; (2)

Scaling: If W is a solution to unmixing problem, so is DW where D is a diagonal matrix.

Independent Component Analysis

Algorithm

Algorithm (Independent Component Analysis)

Input: Data vectors $\{x_1, \dots, x_n\} \in \mathbb{R}^N$.

- ① *Compute the sample mean $b = \frac{1}{n} \sum_{k=1}^n x_k$, and sample covariance matrix $R = \frac{1}{n} \sum_{k=1}^n (x_k - b)(x_k - b)^T$.*
- ② *Solve the eigenproblem $RE = E\Lambda$, where E is the $N \times N$ orthogonal matrix whose columns are eigenvectors, and Λ is the diagonal matrix of eigenvalues.*
- ③ *Compute $F = R^{-1/2} := E\Lambda^{-1/2}E^T$ and apply it on data, $z_k = F(x_k - b)$, $1 \leq k \leq n$.*
- ④ *Compute the orthogonal matrix Q using the JADE algorithm below.*
- ⑤ *Apply Q on whitened data, $\hat{s}_k = Qz_k$, $1 \leq k \leq n$. Compute $W = QF$.*

Output: Matrix W and independent vectors $\{\hat{s}_1, \hat{s}_2, \dots, \hat{s}_n\}$.

Independent Component Analysis – Cont.

Joint Approximate Diagonalization of Eigenmatrices (JADE)

Algorithm (Cardoso's 4th Order Cumulants Algorithm'92)

Input: Whitened data vectors $\{z_1, \dots, z_n\} \in \mathbb{R}^N$.

- 1 Compute the sample 4th order symmetric cumulant tensor

$$F_{ijkl} = \frac{1}{N} \sum_{t=1}^N z_t(i)z_t(j)z_t(k)z_t(l) - \delta_{i,j}\delta_{k,l} - \delta_{i,k}\delta_{j,l} - \delta_{i,l}\delta_{j,k}.$$

- 2 Compute N eigenmatrices $M_{i,j}$, so that $F(M_{i,j}) = \lambda_{i,j}M_{i,j}$.
- 3 Maximize the criterion

$$J_{\text{JADE}}(Q) = \sum_{i,j} |\lambda_{i,j}|^2 \|\text{diag}(QM_{i,j}Q^T)\|_2^2$$

over orthogonal matrices Q by performing successive rotations marching through all pairs (a, b) of distinct indices in $\{1, \dots, N\}$.

Independent Component Analysis – Cont.

Wang-Amari Natural Stochastic Gradient Algorithm of Bell-Sejnowski MaxEntropy

Algorithm (Wang-Amari'97; Bell-Sejnowski'95)

Input: Sphered data vectors $\{z_1, \dots, z_n\} \in \mathbb{R}^N$; Cumulative distribution functions g_k of each component of s ; Learning rate η .

- 1 Initialize $W^{(0)} = F$.
- 2 Repeat until convergence, or until maximum number of steps reached:
 - 1 Draw a data vector z randomly from data vectors, and compute

$$W^{(t+1)} = W^{(t)} + \eta(I + (1 - 2g(z))z^T)W^{(t)}.$$
 - 2 increment $t \leftarrow t + 1$.

Output: Unmixing $N \times N$ matrix $W = W^{(T)}$.

Independent Component Analysis

Derivation

Dimension Reduction using Laplacian Eigenmaps

Idea

First, convert any relevant data into an undirected weighted graph, hence a symmetric weight matrix W .

The Laplacian eigenmaps solve the following optimization problem:

$$(LE) : \begin{array}{ll} \text{minimize} & \text{trace} \left\{ Y \Delta Y^T \right\} \\ \text{subject to} & Y D Y^T = I_d \end{array}$$

where $\Delta = D - W$ with D the diagonal matrix $D_{ii} = \sum_{k=1}^n W_{i,k}$

The $d \times n$ matrix $Y = [y_1 | \cdots | y_n]$ contains the embedding.

Dimension Reduction using Laplacian Eigenmaps

Algorithm

Algorithm (Dimension Reduction using Laplacian Eigenmaps)

Input: A geometric graph $\{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^N$. Parameters: threshold τ , weight coefficient α , and dimension d .

- 1 *Compute the set of pairwise distances $\|x_i - x_j\|$ and convert them into a set of weights:*

$$W_{i,j} = \begin{cases} \exp(-\alpha\|x_i - x_j\|^2) & \text{if } \|x_i - x_j\| \leq \tau \\ 0 & \text{if otherwise} \end{cases}$$

- 2 *Compute the $d + 1$ bottom eigenvectors of the normalized Laplacian matrix $\tilde{\Delta} = I - D^{-1/2}WD^{-1/2}$, $\tilde{\Delta}e_k = \lambda_k e_k$, $1 \leq k \leq d + 1$, $0 = \lambda_0 \leq \dots \leq \lambda_{d+1}$, where $D = \text{diag}(\sum_{k=1}^n W_{i,k})_{1 \leq i \leq n}$.*

Dimension Reduction using Laplacian Eigenmaps

Algorithm - cont'd

Algorithm (Dimension Reduction using Laplacian Eigenmaps-cont'd)

- ③ Construct the $d \times n$ matrix Y ,

$$Y = \begin{bmatrix} e_2^T \\ \vdots \\ e_{d+1}^T \end{bmatrix} D^{-1/2}$$

- ④ The new geometric graph is obtained by converting the columns of Y into n d -dimensional vectors:

$$\begin{bmatrix} y_1 & | & \cdots & | & y_n \end{bmatrix} = Y$$

Output: $\{y_1, \dots, y_n\} \subset \mathbb{R}^d$.

Example

see:

<http://www.math.umd.edu/~rvbalan/TEACHING/AMSC663Fall2010/PROJECTS/P5/index.html>

Dimension Reduction using LLE

The Idea

Presented in [1]. If data is sufficiently dense, we expect that each data point and its neighbors to lie on or near a (locally) linear patch. We assume we are given the set $\{x_1, \dots, x_n\}$ in the high dimensional space \mathbb{R}^N . Step 1. Find a set of local weights $w_{i,j}$ that best explain the point x_i from its local neighbors:

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^n \|x_i - \sum_j w_{i,j} x_j\|^2 \\ & \text{subject to} && \sum_j w_{i,j} = 1, \quad i = 1, \dots, n \end{aligned}$$

Step 2. Find the points $\{y_1, \dots, y_n\} \subset \mathbb{R}^d$ that minimize

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^n \|y_i - \sum_j w_{i,j} y_j\|^2 \\ & \text{subject to} && \sum_{i=1}^n y_i = 0 \\ & && \frac{1}{n} \sum_{i=1}^n y_i y_i^T = I_d \end{aligned}$$

Dimension Reduction using LLE

Derivations (1)

Step 1. The weights are obtained by solving a constrained least-squares problem. The optimization problem decouples for each i . The Lagrangian at fixed $i \in \{1, 2, \dots, n\}$ is

$$L((w_{i,j})_{j \in N}, \lambda) = \|x_i - \sum_{j \in N} w_{i,j} x_j\|^2 + \lambda (\sum_{j \in N} w_{i,j} - 1)$$

where N denotes its K -neighborhood of closest K vertices. Use $\sum_{j \in N} w_{i,j} = 1$ to regroup the terms in the first term, and then expand the square:

$$L() = \left\| \sum_{j \in N} w_{i,j} (x_i - x_j) \right\|^2 + \lambda (\sum_{j \in N} w_{i,j} - 1) = w^T C w + \lambda w^T \cdot \mathbf{1} - \lambda$$

where C is the $K \times K$ covariance matrix $C_{j,k} = \langle x_j - x_i, x_k - x_i \rangle$. Set $\nabla_w L = 0$ and solve for w . $\nabla_w L = 2C \cdot w + \lambda \mathbf{1}$.

Dimension Reduction using LLE

Derivations (2)

$$w = -\frac{\lambda}{2} C^{-1} \cdot \mathbf{1}$$

The multiplier λ is obtained from the constraint $w^T \cdot \mathbf{1} = 1$:

$\lambda = -\frac{2}{\mathbf{1}^T C^{-1} \mathbf{1}}$. Thus

$$w = \frac{C^{-1} \cdot \mathbf{1}}{\mathbf{1}^T C^{-1} \mathbf{1}}$$

Step 2. The embedding in the lower dimensional space is obtained as follows. First denote $Y = [y_1 | \dots | y_n]$ a $d \times n$ matrix. Then

$$\begin{aligned} \sum_{i=1}^n \|y_i - \sum_j w_{i,j} y_j\|^2 &= \sum_{i=1}^n \langle y_i, y_i \rangle - 2 \sum_{i=1}^n \sum_j w_{i,j} \langle y_i, y_j \rangle + \sum_{i=1}^n \sum_{j,k} w_{i,j} w_{i,k} \langle y_j, y_k \rangle \\ &= \text{trace}(YY^T) - 2\text{trace}(YWY^T) + \text{trace}(YW^T WY^T) = \end{aligned}$$

Dimension Reduction using LLE

Derivations (3)

$$= \text{trace}(Y(I - W)^T(I - W)Y^T).$$

where W is the $n \times n$ (non-symmetric) matrix of weights. The optimization problem becomes:

$$\begin{aligned} \text{minimize} \quad & \text{trace}(Y(I - W)^T(I - W)Y^T) \\ \text{subject to} \quad & Y \cdot \mathbf{1} = 0 \\ & YY^T = I_d \end{aligned}$$

Just as the graph Laplacian, the solution is given by the eigenvectors corresponding to the smallest eigenvalues of $(I - W)^T(I - W)$. The condition $Y \cdot \mathbf{1} = 0$ rules out the lowest eigenvector (which is $\mathbf{1}$), and requires rows in Y to be orthogonal to this eigenvector. Therefore, the rows in Y are taken to be the eigenvectors associated to the smallest $d + 1$ eigenvalues, except the smallest eigenvalue.

Dimension Reduction using LLE

Algorithm

Algorithm (Dimension Reduction using Locally Linear Embedding)

Input: A geometric graph $\{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^N$. Parameters: neighborhood size K and dimension d .

- ① *Finding the weight matrix w : For each point i do the following:*
 - ① *Find its closest K neighbors, say \mathcal{V}_i ;*
 - ② *Compute the $K \times K$ local covariance matrix C ,*

$$C_{j,k} = \langle x_j - x_i, x_k - x_i \rangle.$$
 - ③ *Solve $C \cdot u = \mathbf{1}$ for u ($\mathbf{1}$ denotes the K -vector of 1's).*
 - ④ *Rescale $u = u / (u^T \cdot \mathbf{1})$.*
 - ⑤ *Set $w_{i,j} = u_j$ for $j \in \mathcal{V}_i$.*

Dimension Reduction using LLE

Algorithm - cont'd

Algorithm (Dimension Reduction using Locally Linear Embedding)

② Solving the Eigen Problem:

- ① Create the (typically sparse) matrix $L = (I - W)^T(I - W)$;
- ② Find the bottom $d + 1$ eigenvectors of L (the bottom eigenvector would be $[1, \dots, 1]^T$ associated to eigenvalue 0) $\{e_1, e_2, \dots, e_{d+1}\}$;
- ③ Discard the last vector and insert all other eigenvectors as rows into matrix Y

$$Y = \begin{bmatrix} e_2^T \\ \vdots \\ e_{d+1}^T \end{bmatrix}$$

Output: $\{y_1, \dots, y_n\} \subset \mathbb{R}^d$ as columns from

$$\begin{bmatrix} y_1 & | & \cdots & | & y_n \end{bmatrix} = Y$$

Dimension Reduction using Isomaps

The Idea

Presented in [2]. The idea is to first estimate all pairwise distances, and then use the nearly isometric embedding algorithm with full data we studied in Lecture 7.

For each node in the graph we define the distance to the nearest K neighbors using the Euclidean metric. The distance to further nodes is defined as the geodesic distance w.r.t. these local distances.

Dimension Reduction using Isomaps

Algorithm

Algorithm (Dimension Reduction using Isomap)

Input: A geometric graph $\{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^N$. Parameters: neighborhood size K and dimension d .

- ① *Construct the symmetric matrix S of squared pairwise distances:

 - ① *Construct the sparse matrix T , where for each node i find the nearest K neighbors \mathcal{V}_i and set $T_{i,j} = \|x_i - x_j\|_2$, $j \in \mathcal{V}_i$.*
 - ② *For any pair of two nodes (i, j) compute $d_{i,j}$ as the length of the shortest path, $\sum_{p=1}^L T_{k_{p-1}, k_p}$ with $k_0 = i$ and $k_L = j$, using e.g. Dijkstra's algorithm.*
 - ③ *Set $S_{i,j} = d_{i,j}^2$.**

Dimension Reduction using Isomaps

Algorithm - cont'd

Algorithm (Dimension Reduction using Isomap - cont'd)

- ② Compute the Gram matrix G :

$$\rho = \frac{1}{2n} \mathbf{1}^T \cdot S \cdot \mathbf{1}, \quad \nu = \frac{1}{n} (S \cdot \mathbf{1} - \rho \mathbf{1})$$

$$G = \frac{1}{2} \nu \cdot \mathbf{1}^T + \frac{1}{2} \mathbf{1} \cdot \nu^T - \frac{1}{2} S$$

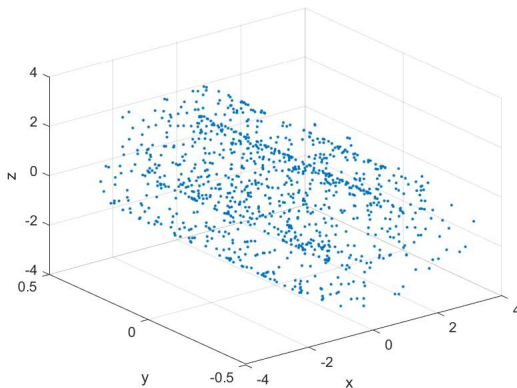
- ③ Find the top d eigenvectors of G , say e_1, \dots, e_d so that $GE = E\Lambda$, form the matrix Y and then collect the columns:

$$Y = \Lambda^{1/2} \begin{bmatrix} e_1^T \\ \vdots \\ e_d^T \end{bmatrix} = \begin{bmatrix} y_1 & | & \cdots & | & y_n \end{bmatrix}$$

Output: $\{y_1, \dots, y_n\} \subset \mathbb{R}^d$.

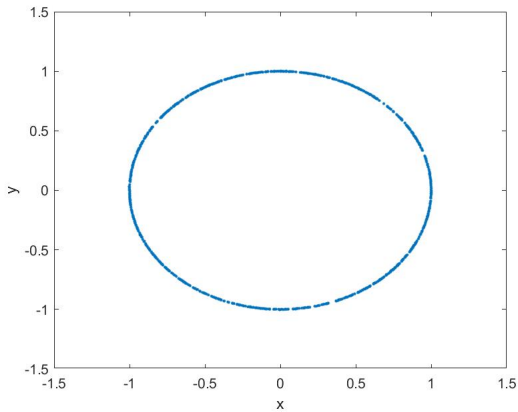
Data Sets

The Swiss Roll



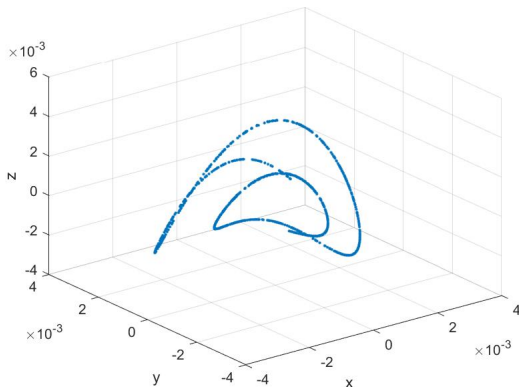
Data Sets

The Circle



Dimension Reduction for the Swiss Roll

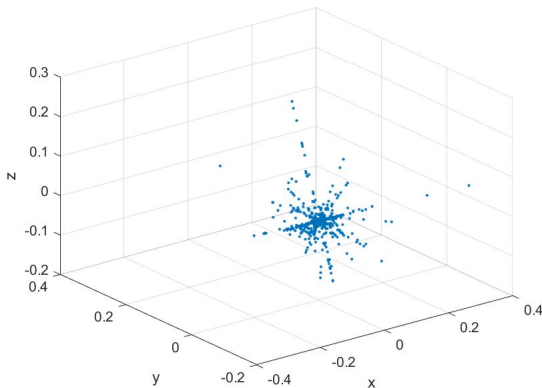
Laplacian Eigenmap



Parameters: $d = 3$, $W_{i,j} = \exp(-0.1\|x_i - x_j\|^2)$, for all i, j .

Dimension Reduction for the Swiss Roll

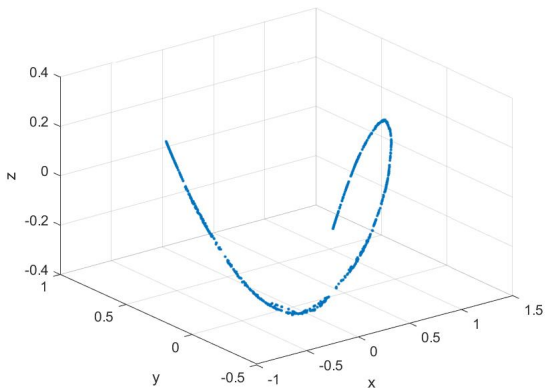
Local Linear Embedding (LLE)



Parameters: $d = 3$, $K = 2$.

Dimension Reduction for the Swiss Roll

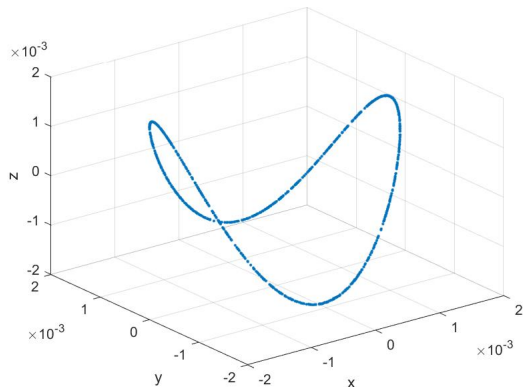
ISOMAP



Parameters: $d = 3$, $K = 10$.

Dimension Reduction for the Circle

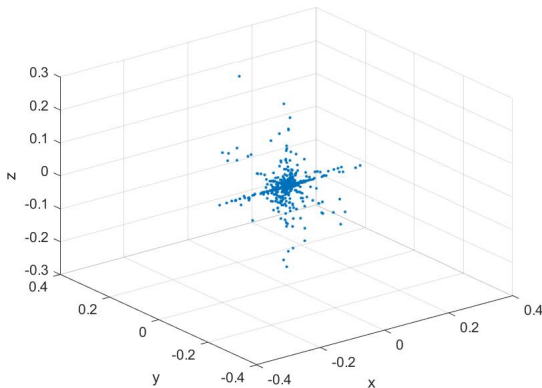
Laplacian Eigenmap



Parameters: $d = 3$, $W_{i,j} = \exp(-0.1\|x_i - x_j\|^2)$, for all i, j .

Dimension Reduction for the Circle

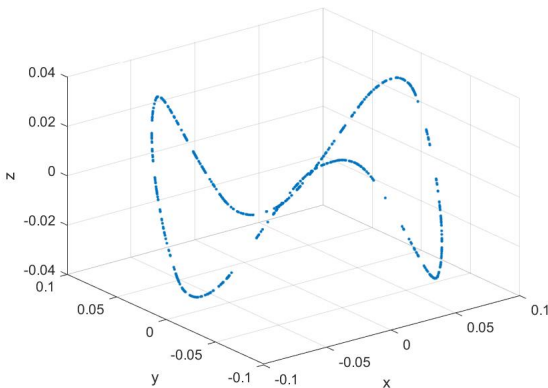
Local Linear Embedding (LLE)



Parameters: $d = 3$, $K = 2$.

Dimension Reduction for the Circle

ISOMAP



Parameters: $d = 3$, $K = 10$.

References



S.T. Roweis, L.K. Saul, Locally linear embedding, *Science* 290, 2323–2326 (2000).



J.B. Tenenbaum, V. de Silva, J.C. Langford, A global geometric framework for nonlinear dimensionality reduction, *Science* 290, 2319–2323 (2000).