

Introduction to Machine Learning

Lecture 7 Representation and Clustering - Nonlinear Dimensionality Reduction

Hongteng Xu



中國人民大學
RENMIN UNIVERSITY OF CHINA

高瓴人工智能學院
Gaoling School of Artificial Intelligence

Outline

Review

- ▶ **Curse of dimensionality**
- ▶ **Principal component analysis:** principle, implementation, and statistical ML viewpoint
- ▶ **Other linear dimensionality reduction method:** RPCA, NMF, Subspace clustering, compressive sensing

Outline

Review

- ▶ **Curse of dimensionality**
- ▶ **Principal component analysis:** principle, implementation, and statistical ML viewpoint
- ▶ **Other linear dimensionality reduction method:** RPCA, NMF, Subspace clustering, compressive sensing

Today

- ▶ Manifold learning (MDS, ISOMAP, LLE, Diffusion Map, ...)
- ▶ Large-scale manifold learning (t-SNE)
- ▶ Kernel method (Kernel PCA)
- ▶ Autoencoders

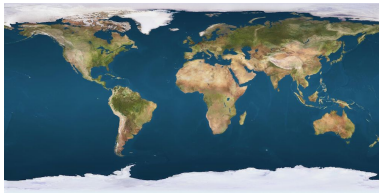
Linear and Nonlinear Dimensionality Reduction



3D data



Linear map to 2D space



Nonlinear map to 2D space.

Recall The Desired Properties of Dimensionality Reduction

- ▶ **Minimizing reconstruction error**

$$\exists g : \mathcal{Z} \mapsto \mathcal{X}, \mathbf{x} \approx g(f(\mathbf{x})). \quad (1)$$

- ▶ (Equivalently,) **Maximizing mutual information** (or minimizing information loss)

Recall The Desired Properties of Dimensionality Reduction

- ▶ **Minimizing reconstruction error**

$$\exists g : \mathcal{Z} \mapsto \mathcal{X}, \mathbf{x} \approx g(f(\mathbf{x})). \quad (1)$$

- ▶ (Equivalently,) **Maximizing mutual information** (or minimizing information loss)
- ▶ **(Nearly) Isometry:**

$$d_{\mathcal{Z}}(f(\mathbf{x}), f(\mathbf{x}')) \approx d_{\mathcal{X}}(\mathbf{x}, \mathbf{x}'). \quad (2)$$

The Properties of Various Linear DR Methods

Method	(R)PCA	NMF	Compressive Sensing
--------	--------	-----	---------------------

The Properties of Various Linear DR Methods

Method	(R)PCA	NMF	Compressive Sensing
Reconstruction power	Yes	Yes	Conditionally, Yes

The Properties of Various Linear DR Methods

Method	(R)PCA	NMF	Compressive Sensing
Reconstruction power	Yes	Yes	Conditionally, Yes
(Nearly) Isometry	No	No	Conditionally, Yes

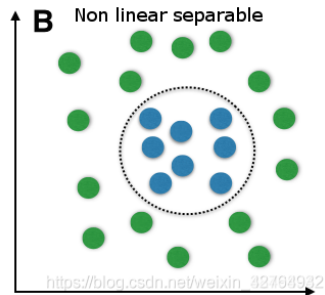
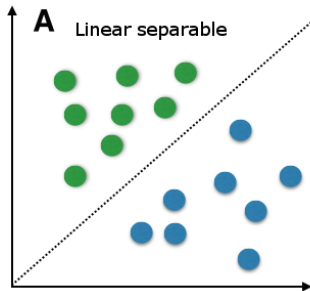
The Properties of Various Linear DR Methods

Method	(R)PCA	NMF	Compressive Sensing
Reconstruction power	Yes	Yes	Conditionally, Yes
(Nearly) Isometry	No	No	Conditionally, Yes

- The conditions of compressive sensing is based on the “sparse representation” assumption of data and the restricted isometric property (RIP) of random projection.

A Typical Case Breaking Isometry Seriously

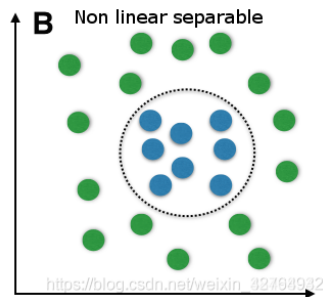
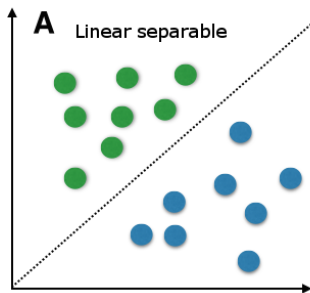
Linear inseparable data



- ▶ Classic linear DR method does not work well for linear inseparable data in general.

A Typical Case Breaking Isometry Seriously

Linear inseparable data



- ▶ Classic linear DR method does not work well for linear inseparable data in general.
- ▶ **Q: Can we obtain better isometry? If yes, how to do it?**

Typical Nonlinear Dimensionality Reduction Methods

(Classic) Manifold Learning (1995 - 2010)

- ▶ Multi-dimensional Scaling (MDS)
- ▶ ISOMAP
- ▶ Locally Linear Embedding (LLE)
- ▶ Eigenmap (Next lecture)
- ▶ Local Tangent Space Alignment (LTSA)

Typical Nonlinear Dimensionality Reduction Methods

(Classic) Manifold Learning (1995 - 2010)

- ▶ Multi-dimensional Scaling (MDS)
- ▶ ISOMAP
- ▶ Locally Linear Embedding (LLE)
- ▶ Eigenmap (Next lecture)
- ▶ Local Tangent Space Alignment (LTSA)

Kernel Methods (also can be viewed as manifold learning)

- ▶ Diffusion Map
- ▶ Kernel PCA
- ▶ t-Distributed Stochastic Neighbor Embedding (t-SNE)

Typical Nonlinear Dimensionality Reduction Methods

(Classic) Manifold Learning (1995 - 2010)

- ▶ Multi-dimensional Scaling (MDS)
- ▶ ISOMAP
- ▶ Locally Linear Embedding (LLE)
- ▶ Eigenmap (Next lecture)
- ▶ Local Tangent Space Alignment (LTSA)

Kernel Methods (also can be viewed as manifold learning)

- ▶ Diffusion Map
- ▶ Kernel PCA
- ▶ t-Distributed Stochastic Neighbor Embedding (t-SNE)

Autoencoding (also can be viewed as manifold learning)

- ▶ Various autoencoders

(Metric) Multi-dimensional Scaling (MDS)

Motivation:

- ▶ Keep isometry as much as possible.

(Metric) Multi-dimensional Scaling (MDS)

Motivation:

- Keep isometry as much as possible.

Principle:

- Given a set of data $\{\mathbf{x}_n\}_{n=1}^N$, we can compute a distance matrix:

$$\mathbf{D} = [d_{ij}] \in \mathbb{R}^{N \times N}, \quad d_{ij} = d(\mathbf{x}_i, \mathbf{x}_j). \quad (3)$$

where $d(\cdot, \cdot)$ can be any valid metric of the sample space.

- Metric MDS aims at finding low-dimensional latent representations $\{\mathbf{z}_n\}_{n=1}^N$ via

$$\min_{\{\mathbf{z}_n\}_{n=1}^N \in \Omega} \underbrace{\left(\sum_{i \neq j} (d_{ij} - \|\mathbf{z}_i - \mathbf{z}_j\|_p)^2 \right)^{1/2}}_{\text{Stress}_d(\{\mathbf{z}_n\}_{n=1}^N)}, \quad (4)$$

where $p = 1$ or 2 in general.

(Classic) Multi-dimensional Scaling (MDS)

Classic MDS is a special case of metric MDS:

- ▶ Use Euclidean distance $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$.

(Classic) Multi-dimensional Scaling (MDS)

Classic MDS is a special case of metric MDS:

- ▶ Use Euclidean distance $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$.
- ▶ Replace the stress loss with a **strain** loss:

$$\min \text{Strain}_d(\{\mathbf{z}_n\}_{n=1}^N) = \left(\frac{\sum_{i,j=1}^N (k_{ij} - \mathbf{z}_i^T \mathbf{z}_j)^2}{\sum_{i,j=1}^N k_{ij}^2} \right)^{1/2},$$

(Classic) Multi-dimensional Scaling (MDS)

Classic MDS is a special case of metric MDS:

- ▶ Use Euclidean distance $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$.
- ▶ Replace the stress loss with a **strain** loss:

$$\min \text{Strain}_d(\{\mathbf{z}_n\}_{n=1}^N) = \left(\frac{\sum_{i,j=1}^N (k_{ij} - \mathbf{z}_i^T \mathbf{z}_j)^2}{\sum_{i,j=1}^N k_{ij}^2} \right)^{1/2}, \quad (5)$$

where

$$\mathbf{K} = [k_{ij}] = -\frac{1}{2}\mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C}, \quad \mathbf{C} = \mathbf{I}_N - \frac{1}{N}\mathbf{1}_{N \times N} \text{ (Centering matrix)}$$

(Classic) Multi-dimensional Scaling (MDS)

Classic MDS is a special case of metric MDS:

- ▶ Use Euclidean distance $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$.
- ▶ Replace the stress loss with a **strain** loss:

$$\min \text{Strain}_d(\{\mathbf{z}_n\}_{n=1}^N) = \left(\frac{\sum_{i,j=1}^N (k_{ij} - \mathbf{z}_i^T \mathbf{z}_j)^2}{\sum_{i,j=1}^N k_{ij}^2} \right)^{1/2}, \quad (5)$$

where

$$\begin{aligned} \mathbf{K} = [k_{ij}] &= -\frac{1}{2} \mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C}, \quad \mathbf{C} = \mathbf{I}_N - \frac{1}{N} \mathbf{1}_{N \times N} \text{ (Centering matrix)} \\ \Rightarrow \mathbf{K} &= \tilde{\mathbf{X}}\tilde{\mathbf{X}}^T = \mathbf{C}\mathbf{X}\mathbf{X}^T\mathbf{C} \text{ (Derive It)} \end{aligned}$$

(Classic) Multi-dimensional Scaling (MDS)

Classic MDS is a special case of metric MDS:

- ▶ Use Euclidean distance $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$.
- ▶ Replace the stress loss with a **strain** loss:

$$\min \text{Strain}_d(\{\mathbf{z}_n\}_{n=1}^N) = \left(\frac{\sum_{i,j=1}^N (k_{ij} - \mathbf{z}_i^T \mathbf{z}_j)^2}{\sum_{i,j=1}^N k_{ij}^2} \right)^{1/2}, \quad (5)$$

where

$$\begin{aligned} \mathbf{K} = [k_{ij}] &= -\frac{1}{2} \mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C}, \quad \mathbf{C} = \mathbf{I}_N - \frac{1}{N} \mathbf{1}_{N \times N} \text{ (Centering matrix)} \\ \Rightarrow \mathbf{K} &= \tilde{\mathbf{X}}\tilde{\mathbf{X}}^T = \mathbf{C}\mathbf{X}\mathbf{X}^T\mathbf{C} \text{ (Derive It)} \end{aligned} \quad (6)$$

- ▶ **Solution:** If we require L -dimensional \mathbf{z} 's, we have:

$$\mathbf{K} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T, \quad \text{and} \quad \mathbf{Z} = \mathbf{U}_L\mathbf{\Lambda}_L^{\frac{1}{2}}.$$

(Classic) Multi-dimensional Scaling (MDS)

Classic MDS is a special case of metric MDS:

- ▶ Use Euclidean distance $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$.
- ▶ Replace the stress loss with a **strain** loss:

$$\min \text{Strain}_d(\{\mathbf{z}_n\}_{n=1}^N) = \left(\frac{\sum_{i,j=1}^N (k_{ij} - \mathbf{z}_i^T \mathbf{z}_j)^2}{\sum_{i,j=1}^N k_{ij}^2} \right)^{1/2}, \quad (5)$$

where

$$\begin{aligned} \mathbf{K} = [k_{ij}] &= -\frac{1}{2} \mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C}, \quad \mathbf{C} = \mathbf{I}_N - \frac{1}{N} \mathbf{1}_{N \times N} \text{ (Centering matrix)} \\ \Rightarrow \mathbf{K} &= \tilde{\mathbf{X}}\tilde{\mathbf{X}}^T = \mathbf{C}\mathbf{X}\mathbf{X}^T\mathbf{C} \text{ (Derive It)} \end{aligned} \quad (6)$$

- ▶ **Solution:** If we require L -dimensional \mathbf{z} 's, we have:

$$\mathbf{K} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T, \quad \text{and} \quad \mathbf{Z} = \mathbf{U}_L\mathbf{\Lambda}_L^{\frac{1}{2}}. \quad (7)$$

Could you connect it with PCA or Least-Square Data Denoising?

ISOMAP

Motivation:

- ▶ Keep isometry under **geodesic distance** as much as possible.

ISOMAP

Motivation:

- Keep isometry under **geodesic distance** as much as possible.

Principle: Given a set of data $\{\mathbf{x}_n\}_{n=1}^N$

- 1 Determine the neighbors of each data point and construct a K -nearest neighbor (KNN) graph of the data.

ISOMAP

Motivation:

- ▶ Keep isometry under **geodesic distance** as much as possible.

Principle: Given a set of data $\{\mathbf{x}_n\}_{n=1}^N$

- 1 Determine the neighbors of each data point and construct a K -nearest neighbor (KNN) graph of the data.
- 2 Compute the shortest path between arbitrary two nodes and obtain an approximation of the geodesic distance matrix $\mathbf{D} = [d_{ij}]$
 - ▶ Dijkstra's algorithm
 - ▶ Floyd-Warshall algorithm

ISOMAP

Motivation:

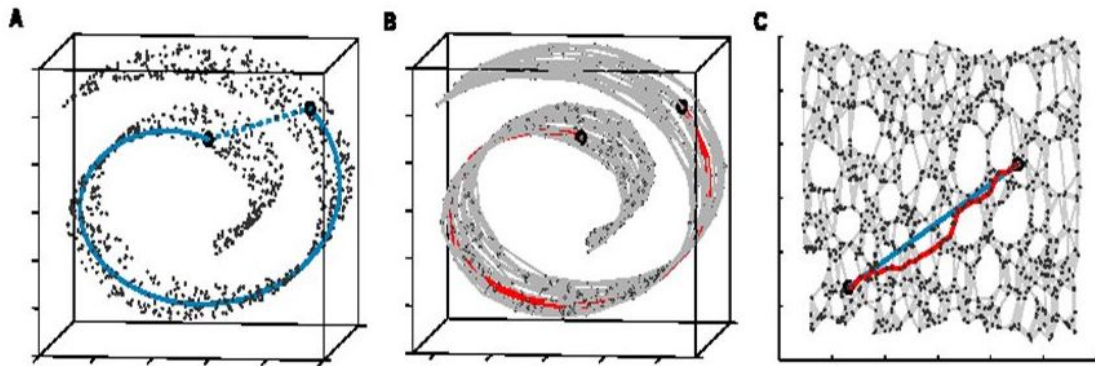
- ▶ Keep isometry under **geodesic distance** as much as possible.

Principle: Given a set of data $\{\mathbf{x}_n\}_{n=1}^N$

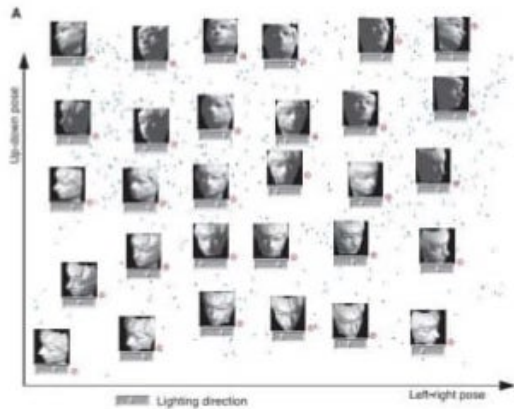
- 1 Determine the neighbors of each data point and construct a K -nearest neighbor (KNN) graph of the data.
- 2 Compute the shortest path between arbitrary two nodes and obtain an approximation of the geodesic distance matrix $\mathbf{D} = [d_{ij}]$
 - ▶ Dijkstra's algorithm
 - ▶ Floyd-Warshall algorithm
- 3 Compute low-dimensional embedding by MDS:

$$\mathbf{K} = -\frac{1}{2}\mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T \Rightarrow \mathbf{Z} = \mathbf{U}_L\mathbf{\Lambda}_L^{\frac{1}{2}} \quad (8)$$

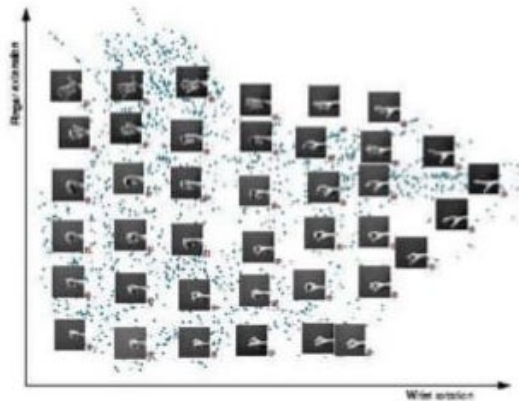
ISOMAP v.s. Metric MDS



ISOMAP



Face varying in pose and illumination



Hand varying in finger extension & wrist rotation

Locally Linear Embedding (LLE)

Motivation:

- ▶ Keep isometry indirectly through inheriting local linear self-representation power.

Locally Linear Embedding (LLE)

Motivation:

- ▶ Keep isometry indirectly through inheriting local linear self-representation power.

Principle:

- ▶ **Local linear self-representation:** Given a sample \mathbf{x} and its K neighbors $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_K]$, where $d(\mathbf{x}, \mathbf{x}_k) \leq \tau$ for $k = 1, \dots, K$.

$$\exists \mathbf{w} \in \mathbb{R}^K, \text{ such that } \mathbf{x} \approx \mathbf{X}\mathbf{w}.$$

Locally Linear Embedding (LLE)

Motivation:

- ▶ Keep isometry indirectly through inheriting local linear self-representation power.

Principle:

- ▶ **Local linear self-representation:** Given a sample \mathbf{x} and its K neighbors $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_K]$, where $d(\mathbf{x}, \mathbf{x}_k) \leq \tau$ for $k = 1, \dots, K$.

$$\exists \mathbf{w} \in \mathbb{R}^K, \text{ such that } \mathbf{x} \approx \mathbf{X}\mathbf{w}. \quad (9)$$

- ▶ **LLE (Locally Linear Embedding):** Given $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$, find low-dimensional representations $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]$ that inherit the local linear self-representation relations.

Locally Linear Embedding (LLE)

Solution:

- 1 Compute the linear coefficients

$$\min_{\mathbf{W}} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{X}_i \mathbf{w}_i\|_2^2, \quad s.t. \quad \sum_{k=1}^K w_{ik_i} = 1, \text{ and } \{k_i\}_{k=1}^K = \text{neighbors' index} \quad (10)$$

where $\mathbf{W} = [w_{ij}] \in \mathbb{R}^{N \times N}$, and $\forall i \in \{1, \dots, N\}$

$$w_{ij} = \begin{cases} w_{ik_i} & j \in \{k_i\}_{k=1}^K \\ 0 & \text{Otherwise.} \end{cases}$$

Locally Linear Embedding (LLE)

Solution:

- 1 Compute the linear coefficients

$$\min_{\mathbf{W}} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{X}_i \mathbf{w}_i\|_2^2, \quad s.t. \sum_{k=1}^K w_{ik_i} = 1, \text{ and } \{k_i\}_{i=1}^N = \text{neighbors' index} \quad (10)$$

where $\mathbf{W} = [\mathbf{w}_{ij}] \in \mathbb{R}^{N \times N}$, and $\forall i \in \{1, \dots, N\}$

$$w_{ij} = \begin{cases} w_{ik_i} & j \in \{k_i\}_{k=1}^K \\ 0 & \text{Otherwise.} \end{cases} \quad (11)$$

- 2 Compute the embeddings:

$$\min_{\mathbf{Z}} \sum_{i=1}^N \|\mathbf{z}_i - \mathbf{Z}_i \mathbf{w}_i\|_2^2 = \min_{\mathbf{Z}} \text{tr}(\mathbf{Z} \Phi \mathbf{Z}^T), \quad s.t. \mathbf{Z} \mathbf{Z}^T = \mathbf{I}_L, \mathbf{Z} \mathbf{1}_L = \mathbf{0}. \quad (12)$$

where $\Phi = [\phi_{ij}]$ and $\phi_{ij} = \delta_{ij} - w_{ij} - w_{ji} + \mathbf{w}_i^T \mathbf{w}_j$.

Locally Linear Embedding (LLE)

Solution:

- 1 Compute the linear coefficients

$$\min_{\mathbf{W}} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{X}_i \mathbf{w}_i\|_2^2, \quad \text{s.t.} \quad \sum_{k=1}^K w_{ik_i} = 1, \quad \text{and} \quad \{k_i\}_{i=1}^N = \text{neighbors' index} \quad (10)$$

where $\mathbf{W} = [\mathbf{w}_{ij}] \in \mathbb{R}^{N \times N}$, and $\forall i \in \{1, \dots, N\}$

$$w_{ij} = \begin{cases} w_{ik_i} & j \in \{k_i\}_{k=1}^K \\ 0 & \text{Otherwise.} \end{cases} \quad (11)$$

- 2 Compute the embeddings:

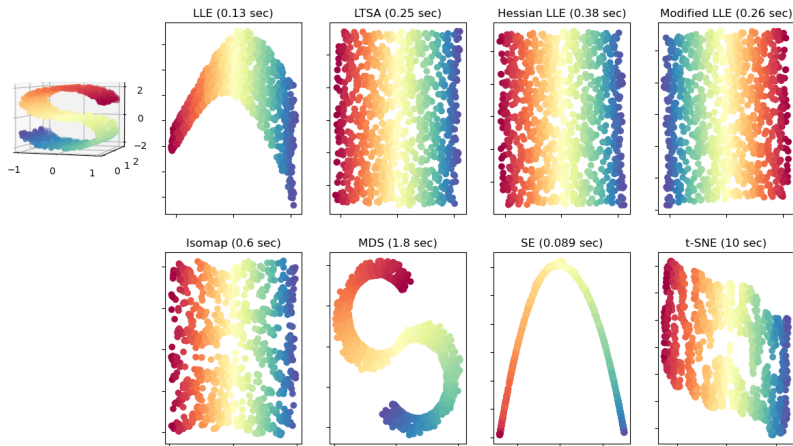
$$\min_{\mathbf{Z}} \sum_{i=1}^N \|\mathbf{z}_i - \mathbf{Z}_i \mathbf{w}_i\|_2^2 = \min_{\mathbf{Z}} \text{tr}(\mathbf{Z} \Phi \mathbf{Z}^T), \quad \text{s.t.} \quad \mathbf{Z} \mathbf{Z}^T = \mathbf{I}_L, \mathbf{Z} \mathbf{1}_L = \mathbf{0}. \quad (12)$$

where $\Phi = [\phi_{ij}]$ and $\phi_{ij} = \delta_{ij} - w_{ij} - w_{ji} + \mathbf{w}_i^T \mathbf{w}_j$.

- 3 $\Phi = \mathbf{U} \Lambda \mathbf{U}^T$, $\lambda_1 \geq \dots \geq \lambda_N$, and $\mathbf{Z} = \mathbf{U}_{N-L+1:N}$.

Classic Manifold Learning Methods

Manifold Learning with 1000 points, 10 neighbors



Most of them require construct a KNN graph first.

Nonlinear DR Methods: Kernel PCA

Motivation:

- ▶ Extend PCA to linear inseparable (equivalently, nonlinear separable) situations.

Nonlinear DR Methods: Kernel PCA

Motivation:

- ▶ Extend PCA to linear inseparable (equivalently, nonlinear separable) situations.
- ▶ Metric MDS and ISOMAP can be viewed as special cases of kernel PCA.

Nonlinear DR Methods: Kernel PCA

Motivation:

- ▶ Extend PCA to linear inseparable (equivalently, nonlinear separable) situations.
- ▶ Metric MDS and ISOMAP can be viewed as special cases of kernel PCA.

(Linear) PCA:

- ▶ Consider the covariance matrix of zero-mean samples: $\mathbf{C} = \mathbf{X}^T \mathbf{X} \in \mathbb{R}^{D \times D}$.
- ▶ Apply Eigenvalue Decomposition $\mathbf{C} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$
- ▶ Consider the projection along top- L principal components: $\mathbf{X} \mathbf{V}_L$.

Nonlinear DR Methods: Kernel PCA

Motivation:

- ▶ Extend PCA to linear inseparable (equivalently, nonlinear separable) situations.
- ▶ Metric MDS and ISOMAP can be viewed as special cases of kernel PCA.

(Linear) PCA:

- ▶ Consider the covariance matrix of zero-mean samples: $\mathbf{C} = \mathbf{X}^T \mathbf{X} \in \mathbb{R}^{D \times D}$.
- ▶ Apply Eigenvalue Decomposition $\mathbf{C} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$
- ▶ Consider the projection along top- L principal components: \mathbf{XV}_L .

SVD-based Implementation:

- ▶ Apply SVD $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$.
- ▶ Consider the projection on top- L principal components: $\mathbf{XV}_L = \mathbf{U}_L \mathbf{\Sigma}_L$.

Kernel PCA

SVD-based Implementation:

- ▶ Apply SVD $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$
- ▶ Consider the projection on top- L principal components: $\mathbf{X}\mathbf{V}_L = \mathbf{U}_L\mathbf{\Sigma}_L$.

Kernel PCA

SVD-based Implementation:

- ▶ Apply SVD $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$
- ▶ Consider the projection on top- L principal components: $\mathbf{X}\mathbf{V}_L = \mathbf{U}_L\mathbf{\Sigma}_L$.

Revisit the SVD-based Implementation from a viewpoint of linear kernel

- ▶ The Gram matrix of linear kernel $\mathbf{K} = \mathbf{X}\mathbf{X}^T \in \mathbb{R}^{N \times N}$.

Kernel PCA

SVD-based Implementation:

- ▶ Apply SVD $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$
- ▶ Consider the projection on top- L principal components: $\mathbf{X}\mathbf{V}_L = \mathbf{U}_L\mathbf{\Sigma}_L$.

Revisit the SVD-based Implementation from a viewpoint of linear kernel

- ▶ The Gram matrix of linear kernel $\mathbf{K} = \mathbf{X}\mathbf{X}^T \in \mathbb{R}^{N \times N}$.
- ▶ The Eigenvalue decomposition: $\mathbf{K} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$

Kernel PCA

SVD-based Implementation:

- ▶ Apply SVD $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$
- ▶ Consider the projection on top- L principal components: $\mathbf{X}\mathbf{V}_L = \mathbf{U}_L\mathbf{\Sigma}_L$.

Revisit the SVD-based Implementation from a viewpoint of linear kernel

- ▶ The Gram matrix of linear kernel $\mathbf{K} = \mathbf{X}\mathbf{X}^T \in \mathbb{R}^{N \times N}$.
- ▶ The Eigenvalue decomposition: $\mathbf{K} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$
- ▶ The PCA corresponds to the scaling of the top- L eigenvectors of \mathbf{K} : $\mathbf{U}_L\mathbf{\Lambda}_L^{\frac{1}{2}}$.

Kernel PCA

SVD-based Implementation:

- ▶ Apply SVD $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$
- ▶ Consider the projection on top- L principal components: $\mathbf{X}\mathbf{V}_L = \mathbf{U}_L\mathbf{\Sigma}_L$.

Revisit the SVD-based Implementation from a viewpoint of linear kernel

- ▶ The Gram matrix of linear kernel $\mathbf{K} = \mathbf{X}\mathbf{X}^T \in \mathbb{R}^{N \times N}$.
- ▶ The Eigenvalue decomposition: $\mathbf{K} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$
- ▶ The PCA corresponds to the scaling of the top- L eigenvectors of \mathbf{K} : $\mathbf{U}_L\mathbf{\Lambda}_L^{\frac{1}{2}}$.

Replacing the linear kernel to arbitrary kernel, you obtain the Kernel PCA:)

- ▶ $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^2, \exp(-\|\mathbf{x} - \mathbf{y}\|^2/h), \dots$

Kernel PCA

Principle:

- ▶ N D -dimensional samples can be linearly separated when $D \geq N$.

Kernel PCA

Principle:

- ▶ N D -dimensional samples can be linearly separated when $D \geq N$.
- ▶ Recall that a kernel function equals to the inner product of samples in an implicit, maybe infinite-dimensional, feature space (Lecture 5).

$$K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathcal{F}} = \phi(\mathbf{x})^T \phi(\mathbf{y})$$

Kernel PCA

Principle:

- ▶ N D -dimensional samples can be linearly separated when $D \geq N$.
- ▶ Recall that a kernel function equals to the inner product of samples in an implicit, maybe infinite-dimensional, feature space (Lecture 5).

$$K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathcal{F}} = \phi(\mathbf{x})^T \phi(\mathbf{y}) \quad (13)$$

- ▶ Because we are never working directly in the (implicit) feature space.
- ▶ The kernel PCA does not compute the principal components directly, but the projections of data onto the top- L components.

Kernel PCA

- Denote $\Phi(\mathbf{X}) = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)]^T \in \mathbb{R}^{N \times \dim(F)}$

Kernel PCA

- ▶ Denote $\Phi(\mathbf{X}) = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)]^T \in \mathbb{R}^{N \times \dim(F)}$
- ▶ An “functional” SVD: $\Phi(\mathbf{X}) = \mathbf{U}\Sigma\mathbf{V}(\mathbf{X})$, where $\mathbf{V}(\mathbf{X}) \in \mathbb{R}^{N \times \dim(F)}$ are principal components of \mathcal{F} .

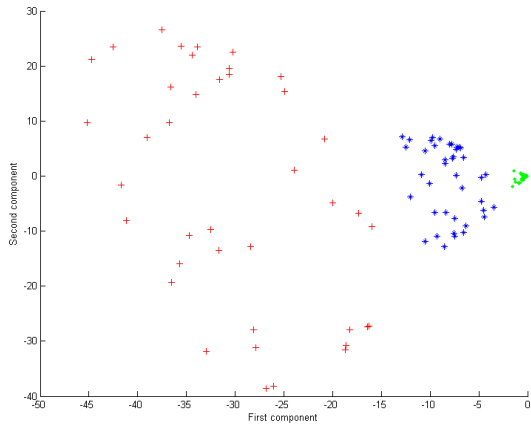
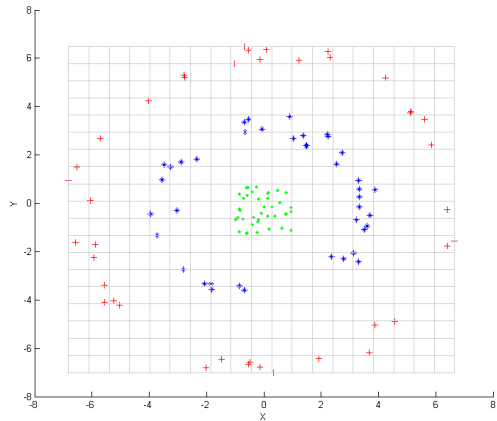
Kernel PCA

- ▶ Denote $\Phi(\mathbf{X}) = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)]^T \in \mathbb{R}^{N \times \dim(F)}$
- ▶ An “functional” SVD: $\Phi(\mathbf{X}) = \mathbf{U}\Sigma\mathbf{V}(\mathbf{X})$, where $\mathbf{V}(\mathbf{X}) \in \mathbb{R}^{N \times \dim(F)}$ are principal components of \mathcal{F} .
- ▶ **Instead of applying the functional SVD, consider the eigenvalue decomposition of $\mathbf{K} = \Phi(\mathbf{X})\Phi^T(\mathbf{X}) = \mathbf{U}\Sigma^2\mathbf{U}^T = \mathbf{U}\Lambda\mathbf{U}^T$**

Kernel PCA

- ▶ Denote $\Phi(\mathbf{X}) = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)]^T \in \mathbb{R}^{N \times \dim(F)}$
- ▶ An “functional” SVD: $\Phi(\mathbf{X}) = \mathbf{U}\Sigma\mathbf{V}(\mathbf{X})$, where $\mathbf{V}(\mathbf{X}) \in \mathbb{R}^{N \times \dim(F)}$ are principal components of \mathcal{F} .
- ▶ **Instead of applying the functional SVD, consider the eigenvalue decomposition of $\mathbf{K} = \Phi(\mathbf{X})\Phi^T(\mathbf{X}) = \mathbf{U}\Sigma^2\mathbf{U}^T = \mathbf{U}\Lambda\mathbf{U}^T$**
- ▶ Kernel PCA = Linear PCA in the feature space associated with the kernel.

Kernel PCA



Revisit MDS and ISOMAP as Kernel PCA Methods

- MDS (Linear Kernel):

$$\mathbf{K} = -\frac{1}{2}\mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C} = \mathbf{C}\mathbf{X}\mathbf{X}^T\mathbf{C} = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^T$$

Revisit MDS and ISOMAP as Kernel PCA Methods

- MDS (Linear Kernel):

$$\mathbf{K} = -\frac{1}{2}\mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C} = \mathbf{C}\mathbf{X}\mathbf{X}^T\mathbf{C} = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^T \quad (14)$$

- ISOMAP (A positive semi-definite kernel based on doubly centered geodesic distance, a.k.a., Mercer kernel):

$$\mathbf{K} = -\frac{1}{2}\mathbf{C}(\mathbf{D}_{geo} \odot \mathbf{D}_{geo})\mathbf{C} \quad (15)$$

t-Stochastic Neighborhood Embedding

Motivation:

- ▶ Large-scale dimensionality reduction.
- ▶ Visualization and clustering of high-dimensional data.

Principle:

- ▶ Given $\{\mathbf{x}_n\}_{n=1}^N$, define a probability p_{ij} that is proportional to the similarity between \mathbf{x}_i and \mathbf{x}_j :

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2N} \quad \text{and} \quad p_{ii} = 0 \quad (16)$$

$$p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / (2\sigma_i^2))}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / (2\sigma_i^2))}. \quad (17)$$

Nonlinear DR Methods: t-Stochastic Neighborhood Embedding

- ▶ t-SNE aims to learn $\{\mathbf{z}_n \in \mathbb{R}^K\}_{n=1}^N$ ($K = 2, 3$ typically) by
 - 1 Define a similarity q_{ij} between \mathbf{z}_i and \mathbf{z}_j as

$$q_{ij} = \frac{(1 + \|\mathbf{z}_i - \mathbf{z}_j\|^2)^{-1}}{\sum_k \sum_{k \neq l} (1 + \|\mathbf{z}_k - \mathbf{z}_l\|^2)^{-1}} \quad \text{and} \quad q_{ii} = 0 \quad (18)$$

where $\{q_{ij}\}$ yield a t-distribution.

Nonlinear DR Methods: t-Stochastic Neighborhood Embedding

- ▶ t-SNE aims to learn $\{\mathbf{z}_n \in \mathbb{R}^K\}_{n=1}^N$ ($K = 2, 3$ typically) by

- 1 Define a similarity q_{ij} between \mathbf{z}_i and \mathbf{z}_j as

$$q_{ij} = \frac{(1 + \|\mathbf{z}_i - \mathbf{z}_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|\mathbf{z}_k - \mathbf{z}_l\|^2)^{-1}} \quad \text{and} \quad q_{ii} = 0 \quad (18)$$

where $\{q_{ij}\}$ yield a t-distribution.

- 2 Learn the \mathbf{z} 's via

$$\min_{\mathbf{z}} \text{KL}(P\|Q) := \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}. \quad (19)$$

SGD is applied to solve this problem.

Nonlinear DR Methods: t-Stochastic Neighborhood Embedding

- ▶ t-SNE aims to learn $\{\mathbf{z}_n \in \mathbb{R}^K\}_{n=1}^N$ ($K = 2, 3$ typically) by

- 1 Define a similarity q_{ij} between \mathbf{z}_i and \mathbf{z}_j as

$$q_{ij} = \frac{(1 + \|\mathbf{z}_i - \mathbf{z}_j\|^2)^{-1}}{\sum_k \sum_{k \neq l} (1 + \|\mathbf{z}_k - \mathbf{z}_l\|^2)^{-1}} \quad \text{and} \quad q_{ii} = 0 \quad (18)$$

where $\{q_{ij}\}$ yield a t-distribution.

- 2 Learn the \mathbf{z} 's via

$$\min_{\mathbf{z}} \text{KL}(P\|Q) := \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}. \quad (19)$$

SGD is applied to solve this problem.

Rationality:

- ▶ Similar \mathbf{x} 's are modeled by nearby \mathbf{z} 's and dissimilar \mathbf{x} 's are modeled by distant \mathbf{z} 's with high probability.
- ▶ Under special settings, approximate a simple form of spectral clustering.

Nonlinear DR Methods: Autoencoders

Drawbacks of the above methods:

- ▶ From inductive paradigm to transductive paradigm
- ▶ High complexity and inscalable inference
- ▶ No reconstruction power

Nonlinear DR Methods: Autoencoders

Drawbacks of the above methods:

- ▶ From inductive paradigm to transductive paradigm
- ▶ High complexity and inscalable inference
- ▶ No reconstruction power

Motivation:

- ▶ Achieve (nearly) isometry and reconstruction power jointly.
- ▶ Achieve inductive inference
- ▶ Reduce complexity of other nonlinear method

Nonlinear DR Methods: Autoencoders

Principle:

- Revisit PCA from a viewpoint of autoencoding.

$$\mathbf{X}^* = \arg \min_{\mathbf{X}} \|\mathbf{X}_{noisy} - \mathbf{X}\|_F^2, \quad s.t. \text{rank}(\mathbf{X}) \leq L$$

Nonlinear DR Methods: Autoencoders

Principle:

- Revisit PCA from a viewpoint of autoencoding.

$$\begin{aligned}\mathbf{X}^* &= \arg \min_{\mathbf{X}} \|\mathbf{X}_{noisy} - \mathbf{X}\|_F^2, \quad s.t. \text{rank}(\mathbf{X}) \leq L \\ &= \mathbf{U}_L \mathbf{\Sigma}_L \mathbf{V}_L^T\end{aligned}$$

Nonlinear DR Methods: Autoencoders

Principle:

- Revisit PCA from a viewpoint of autoencoding.

$$\begin{aligned}\mathbf{X}^* &= \arg \min_{\mathbf{X}} \|\mathbf{X}_{noisy} - \mathbf{X}\|_F^2, \quad \text{s.t. rank}(\mathbf{X}) \leq L \\ &= \mathbf{U}_L \Sigma_L \mathbf{V}_L^T \\ \Rightarrow \text{Encoder: } \mathbf{X}_{noisy} \mathbf{V}_L &= \mathbf{U}_L \Sigma_L \\ \text{Decoder: } \mathbf{X}^* &= \mathbf{X}_{noisy} \mathbf{V}_L \mathbf{V}_L^T\end{aligned}\tag{20}$$

- \mathbf{V}_K and \mathbf{V}_K^T work as the encoder and the decoder, respectively.

Nonlinear DR Methods: Autoencoders

- In general, a typical autocoder consists of

$$\text{Encoder: } f: \mathcal{X} \mapsto \mathcal{Z}, \quad \text{Decoder: } g: \mathcal{Z} \mapsto \mathcal{X}$$

Nonlinear DR Methods: Autoencoders

- In general, a typical autocoder consists of

$$\text{Encoder: } f: \mathcal{X} \mapsto \mathcal{Z}, \quad \text{Decoder: } g: \mathcal{Z} \mapsto \mathcal{X} \quad (21)$$

- Given a set of data $\{\mathbf{x}_n\}_{n=1}^N$,

$$\min_{f,g} \underbrace{\sum_{n=1}^N \text{loss}(\mathbf{x}_n, g(f(\mathbf{x}_n)))}_{\text{loss}(\mathbf{X}, g(f(\mathbf{X})))} + \text{reg}(q_{\mathbf{Z}|\mathbf{X}}, p_{\mathbf{Z}}) \quad (22)$$

where $p_{\mathbf{Z}}$ is a predefined prior distribution of latent codes, while $q_{\mathbf{Z}|\mathbf{X}}$ is the posterior distribution of the latent codes given the corresponding data, which is determined by the encoder f .

Nonlinear DR Methods: Autoencoders

- In general, a typical autocoder consists of

$$\text{Encoder: } f: \mathcal{X} \mapsto \mathcal{Z}, \quad \text{Decoder: } g: \mathcal{Z} \mapsto \mathcal{X} \quad (21)$$

- Given a set of data $\{\mathbf{x}_n\}_{n=1}^N$,

$$\min_{f,g} \underbrace{\sum_{n=1}^N \text{loss}(\mathbf{x}_n, g(f(\mathbf{x}_n)))}_{\text{loss}(\mathbf{X}, g(f(\mathbf{X})))} + \text{reg}(q_{\mathbf{Z}|\mathbf{X}}, p_{\mathbf{Z}}) \quad (22)$$

where $p_{\mathbf{Z}}$ is a predefined prior distribution of latent codes, while $q_{\mathbf{Z}|\mathbf{X}}$ is the posterior distribution of the latent codes given the corresponding data, which is determined by the encoder f .

- More details will be given in the following lectures.

In Summary

- ▶ Introduce classic manifold learning methods
- ▶ Introduce the most widely-used manifold learning method, t-SNE
- ▶ Introduce the kernelization of PCA.
- ▶ Introduce autoencoders (briefly)

Next...

- ▶ Data representation, clustering, and unsupervised learning
- ▶ Kmeans and spectral clustering
- ▶ Evaluation of clustering method.