



浙江大学

ZheJiang University



人工智能研究所

Institute of Artificial Intelligence

Artificial Intelligence

Continuous Latent Variables

Donghui Wang
AI Institute@ZJU
2015.4



Contents

- Principal Component Analysis (PCA)
- Probabilistic PCA
- Kernel PCA
- Nonlinear Latent Variable Models

References:

1. Bishop. *“Pattern Recognition and Machine Learning”, Chapter 12.* 2006.



浙江大学

ZheJiang University



人工智能研究所

Institute of Artificial Intelligence

Principal Component Analysis (PCA)



What's PCA?

- PCA - Principal Component Analysis
 - *Karhunen-Loève* transform
- Two commonly used definitions:
 - The orthogonal projection of the data onto a lower dimensional linear space, known as the *principal subspace*, such that the variance of the projected data is maximized.
 - The linear projection that minimizes the average projection cost, defined as the mean squared distance between the data points and their projections.

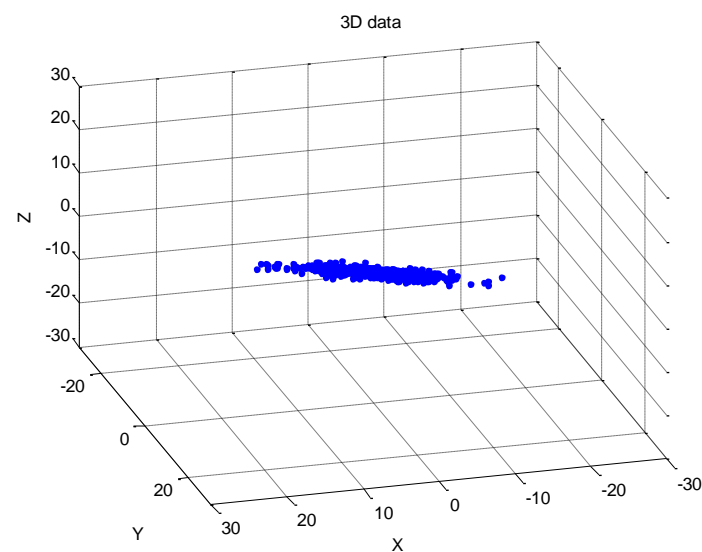
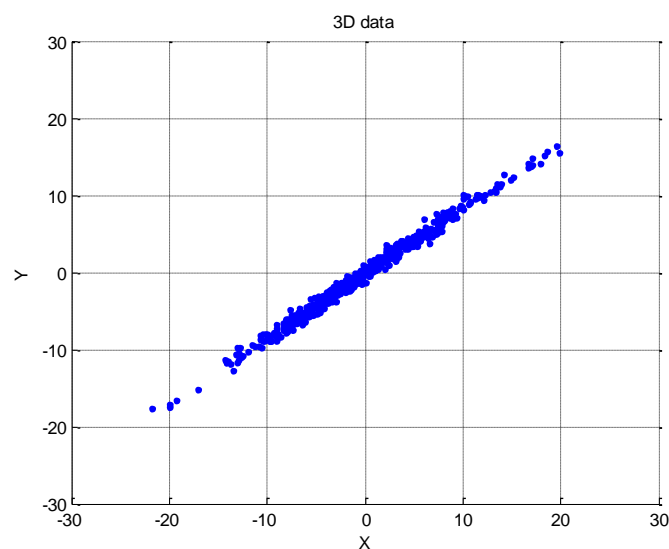
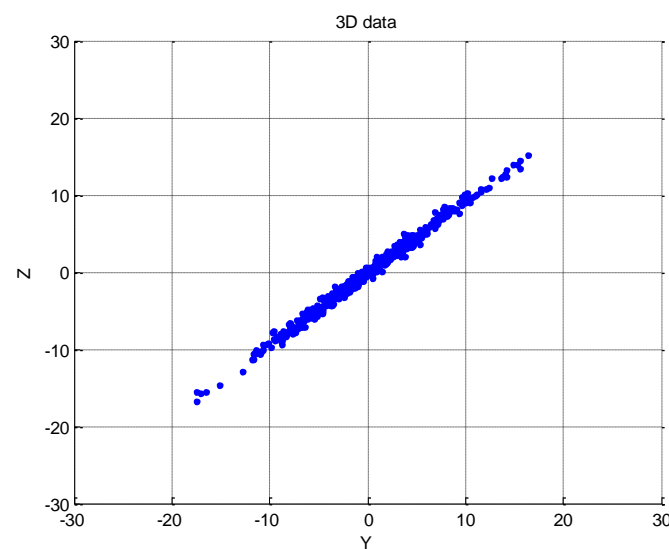
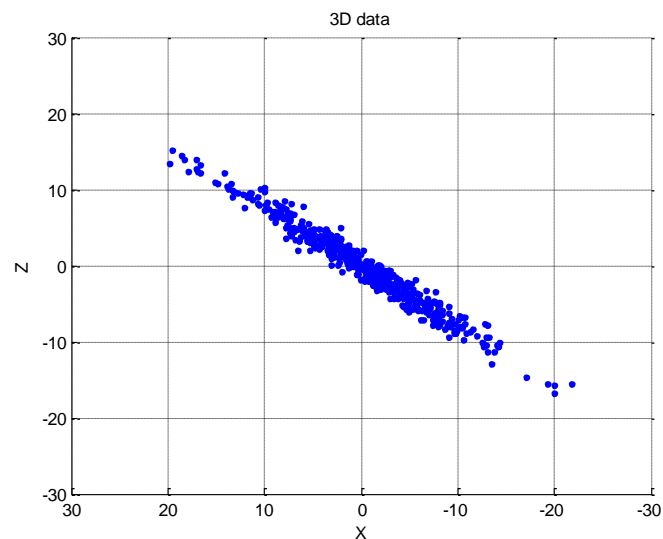


Why PCA?

- Dimensionality reduction
 - Avoid the curse of dimensionality
- Lossy data compression
- Feature extraction
- Data visualization
 - How to visualize high-dimensional data?
- ...



- Dir
- Lo
- Fe
- Da
- ...





Maximum variance formulation

- Consider a data set of observations $\{\mathbf{x}_n\}$ where $n=1, \dots, N$, and \mathbf{x}_n is a D -dimensional vector.
- Goal:** project the data onto a M -dimensional subspace ($M < D$) while maximizing the variance of the projected data.

- $M=1$:

– define the direction of this space using a D -dimensional unit vector \mathbf{u}_1 and $\mathbf{u}_1^T \mathbf{u}_1 = 1$.

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad \mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T \quad \Rightarrow \quad \frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}}\}^2 = \boxed{\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1}$$

Maximizing it.

Lagrange multiplier

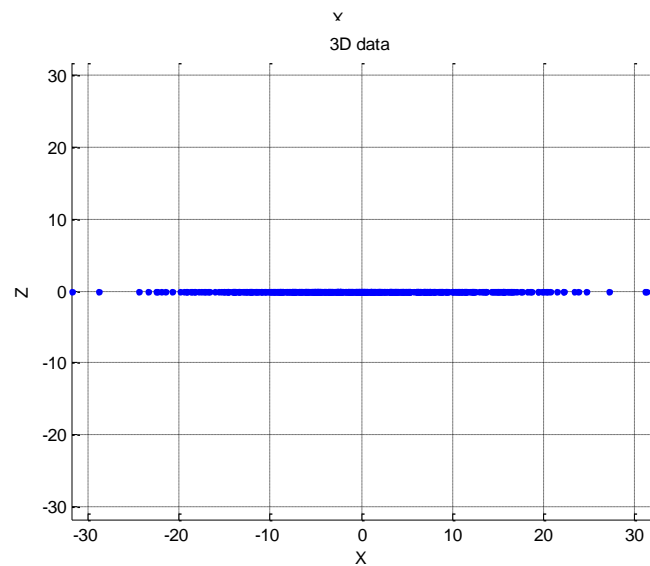
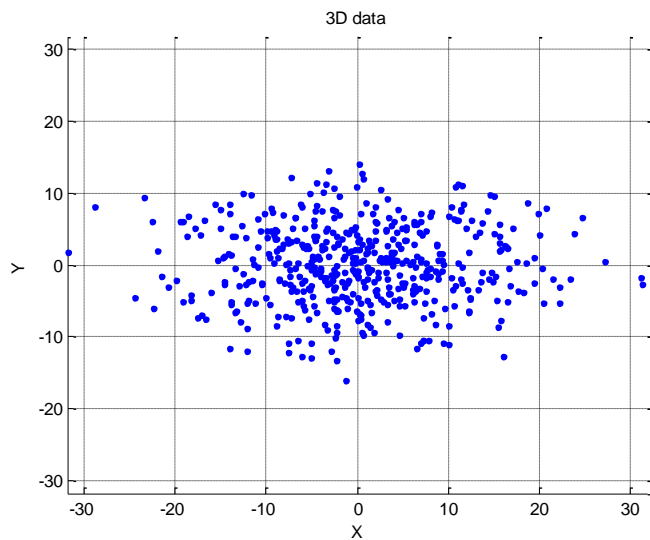
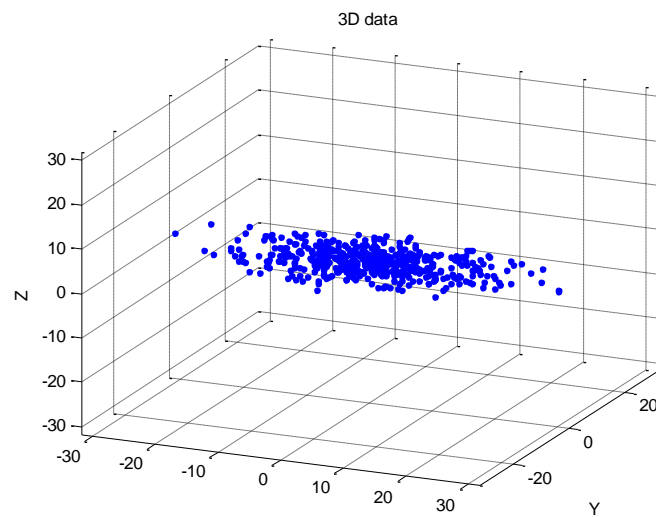
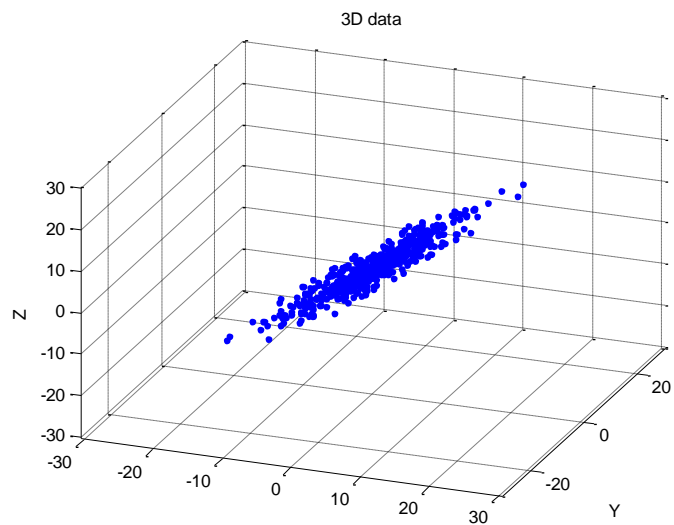
$$\mathbf{u}_1^T \mathbf{u}_1 = 1$$

$$\boxed{\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)} \Rightarrow \mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1 \Rightarrow \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1$$

The variance will be a maximum when we set \mathbf{u}_1 equal to the eigenvector having the largest eigenvalue λ_1 .

- $M > 1$:

– Optimal linear projection is now defined by the M eigenvectors of the data covariance matrix \mathbf{S} corresponding to the M largest eigenvalues.





Minimum-error formulation

- Introduce a complete orthonormal set of D-dimensional basis vector $\{\mathbf{u}_i\}$ where $i=1,\dots,D$ that satisfy $\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$.
- For a data set of observations $\{\mathbf{x}_n\}$, we have $\mathbf{x}_n = \sum_{i=1}^D \alpha_{ni} \mathbf{u}_i \Rightarrow \mathbf{x}_n = \sum_{i=1}^D (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i$
- Goal:** approximate \mathbf{x}_n using a representation involving a restricted number $M < D$ of variables corresponding to a projection onto a lower-dimensional.

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i \xrightarrow{\text{Minimizing it.}} J = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 \xrightarrow{\quad} \begin{aligned} z_{nj} &= \mathbf{x}_n^T \mathbf{u}_j \\ b_j &= \bar{\mathbf{x}}^T \mathbf{u}_j \end{aligned}$$

$$\Rightarrow \mathbf{x}_n - \tilde{\mathbf{x}}_n = \sum_{i=M+1}^D \{(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_i\} \mathbf{u}_i \Rightarrow J = \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i)^2 = \sum_{i=M+1}^D \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i$$

$$\Rightarrow \tilde{J} = \mathbf{u}_2^T \mathbf{S} \mathbf{u}_2 + \lambda_2 (1 - \mathbf{u}_2^T \mathbf{u}_2) \Rightarrow \mathbf{S} \mathbf{u}_i = \lambda_i \mathbf{u}_i \Rightarrow J = \sum_{i=M+1}^D \lambda_i$$



Applications of PCA

- Data compression (offline digits data set):
 - First four eigenvectors and eigenvalues:

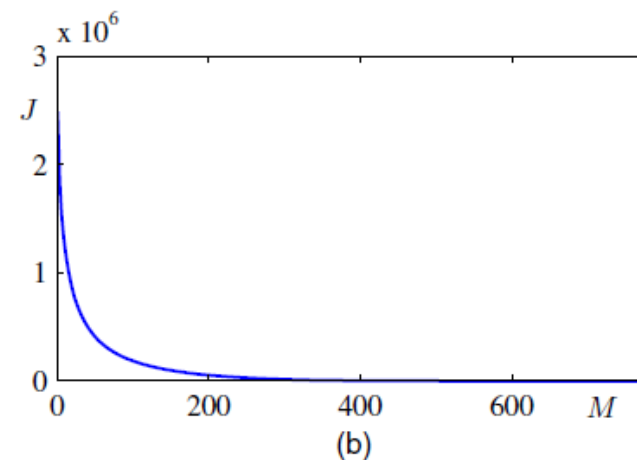
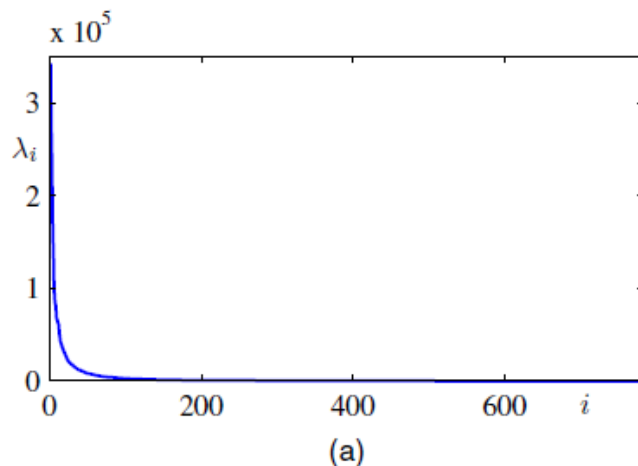
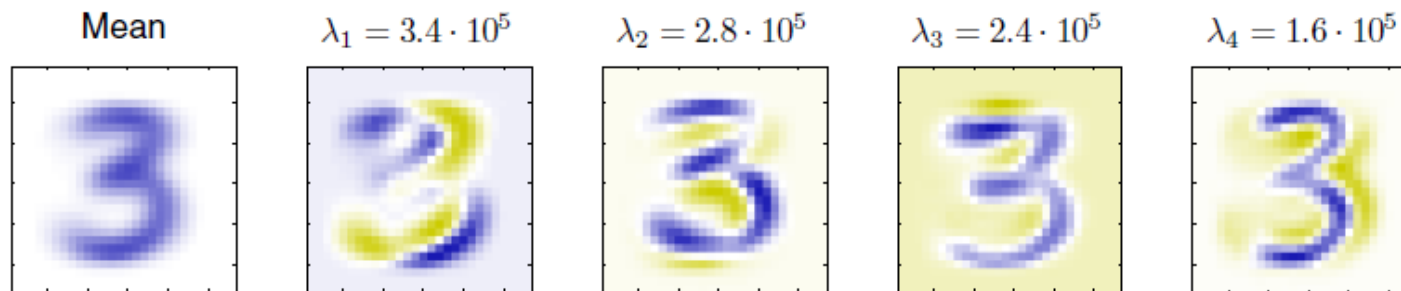


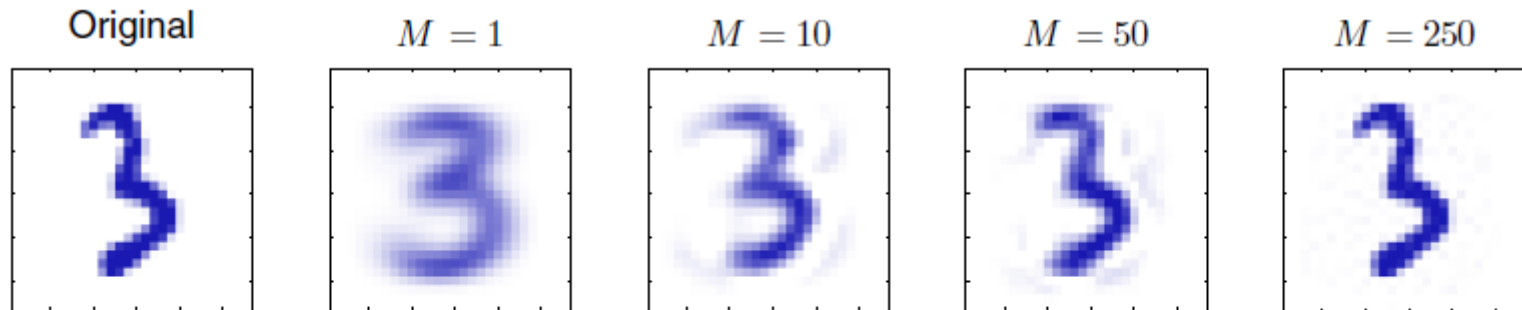
Figure 12.4 (a) Plot of the eigenvalue spectrum for the off-line digits data set. (b) Plot of the sum of the discarded eigenvalues, which represents the sum-of-squares distortion J introduced by projecting the data onto a principal component subspace of dimensionality M .



Applications of PCA

- Data compression (offline digits data set):

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i + \sum_{i=M+1}^D (\bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i = \bar{\mathbf{x}} + \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i$$



An original example from the off-line digits data set together with its PCA reconstructions obtained by retaining M principal components for various values of M . As M increases the reconstruction becomes more accurate and would become perfect when $M = D = 28 \times 28 = 784$.



Applications of PCA

- Data normalization (pre-processing):

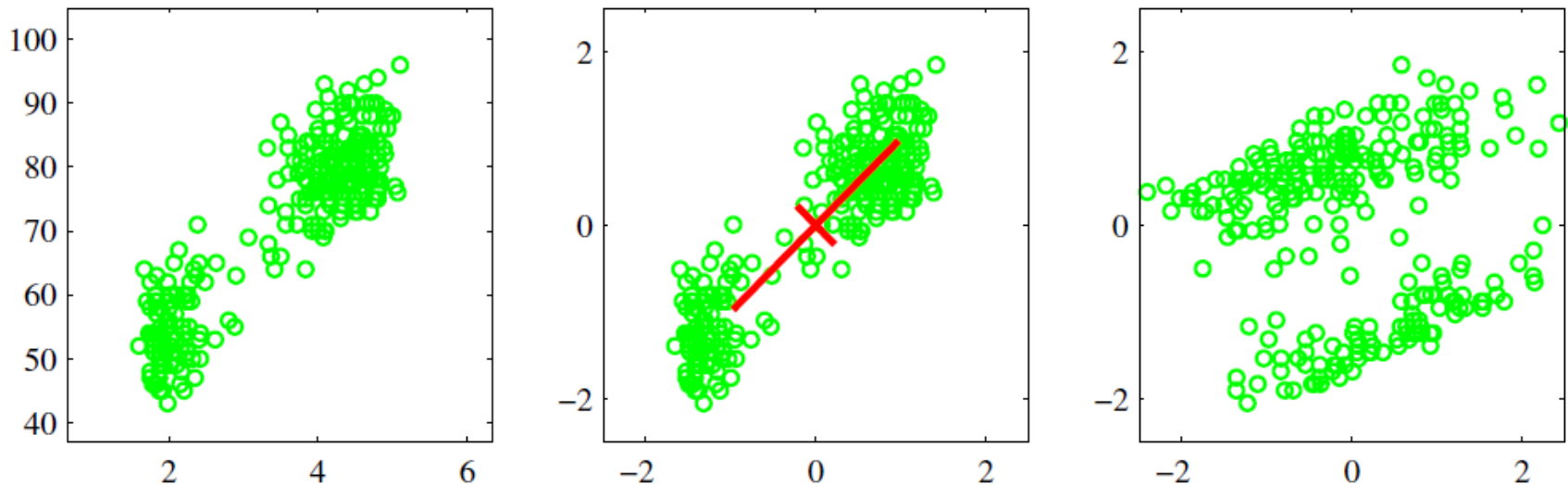
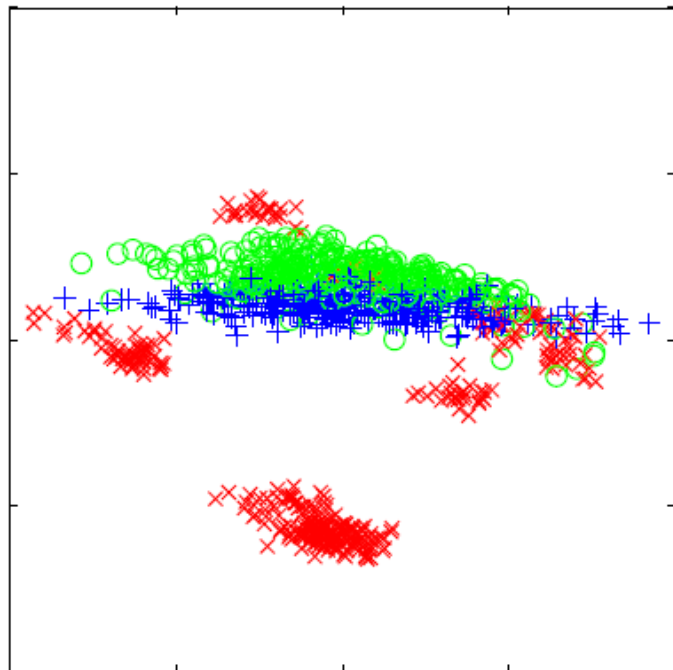


Figure 12.6 Illustration of the effects of linear pre-processing applied to the Old Faithful data set. The plot on the left shows the original data. The centre plot shows the result of standardizing the individual variables to zero mean and unit variance. Also shown are the principal axes of this normalized data set, plotted over the range $\pm\lambda_i^{1/2}$. The plot on the right shows the result of whitening of the data to give it zero mean and unit covariance.

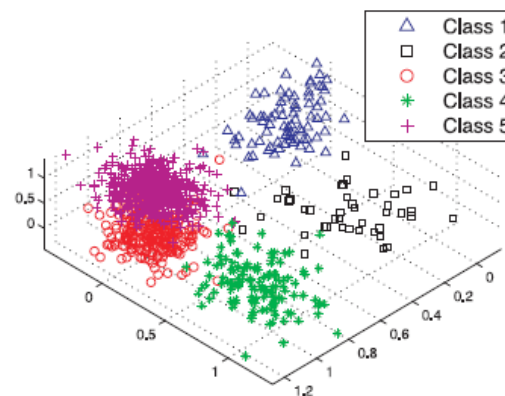


Applications of PCA

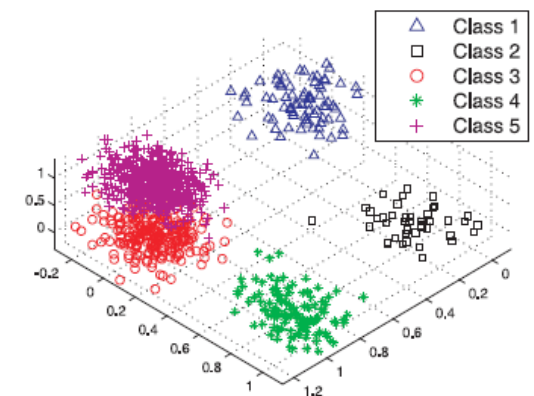
- Data visualization:



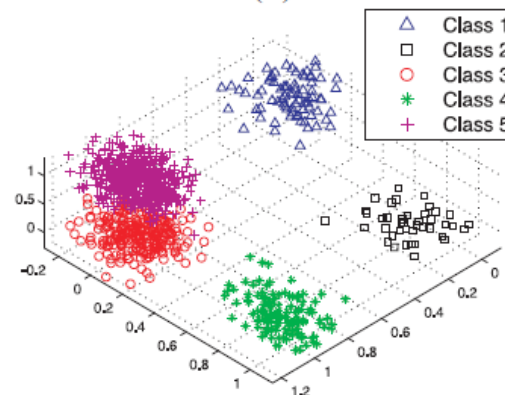
Visualization of the oil flow data set obtained by projecting the data onto the first two principal components. The red, blue, and green points correspond to the 'laminar', 'homogeneous', and 'annular' flow configurations respectively.



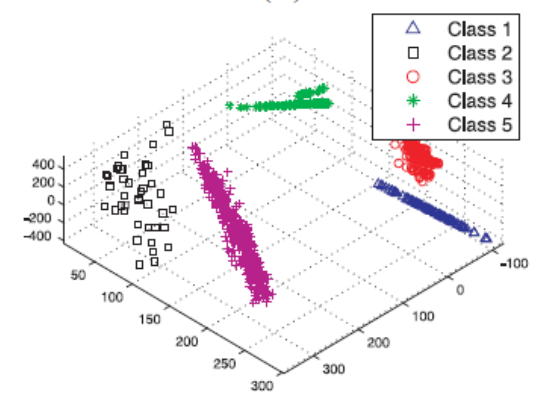
(a)



(b)



(c)

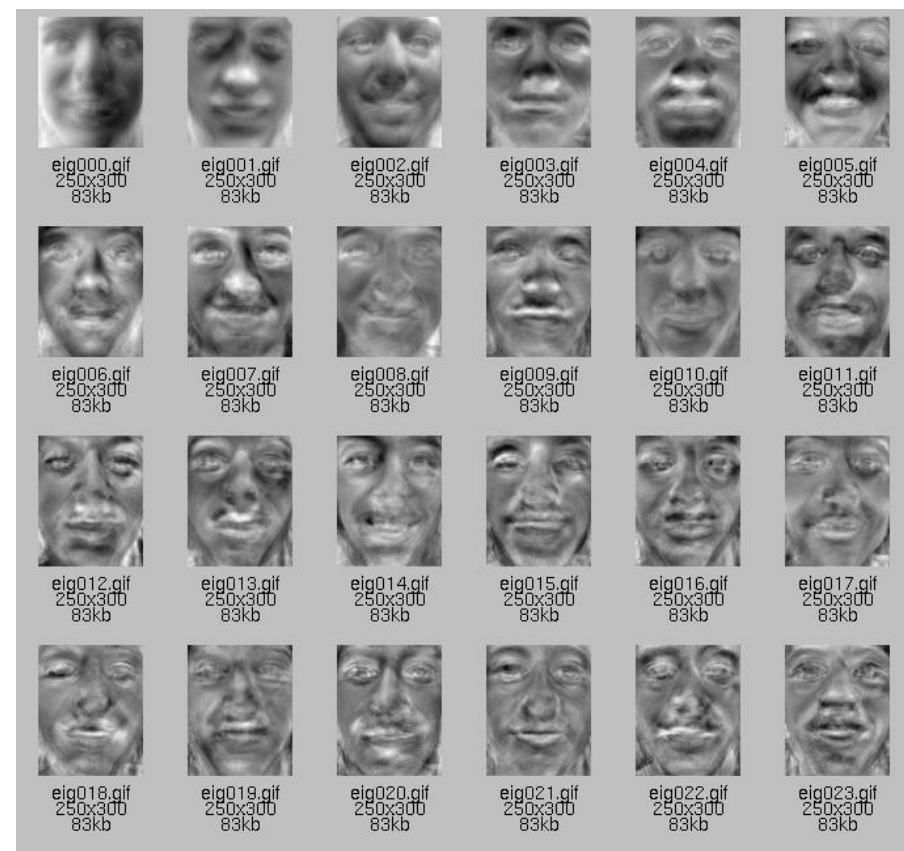
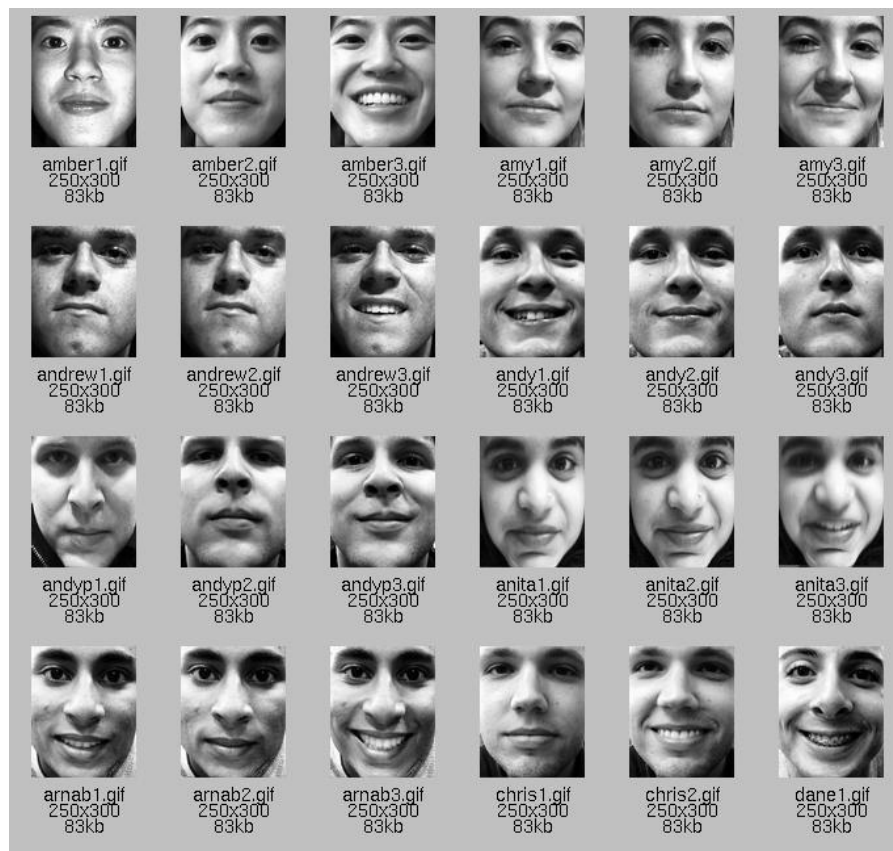


(d)



Applications of PCA

- Face recognition (eigenface):

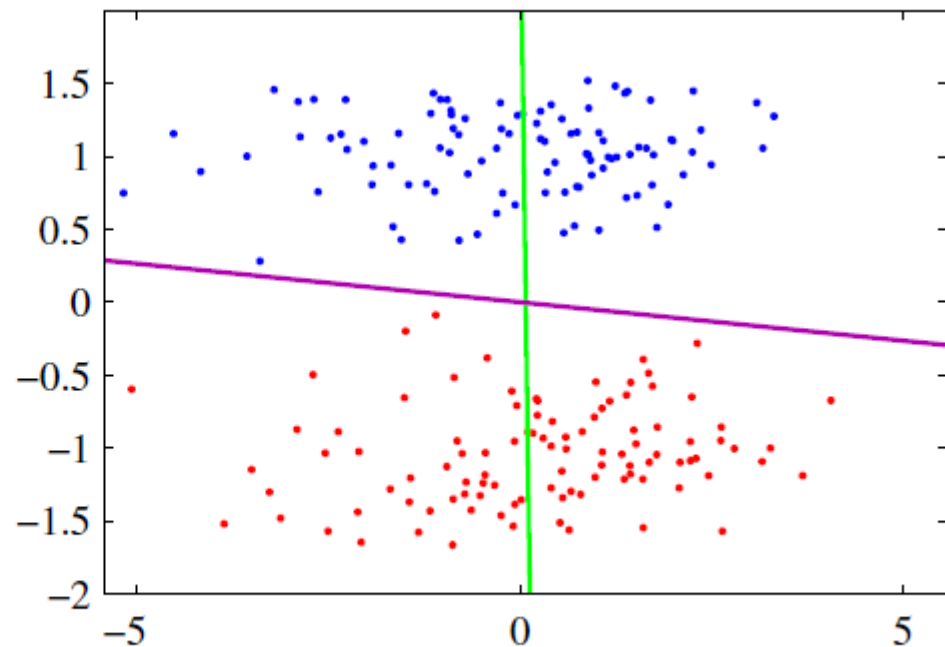




Applications of PCA

- PCA vs. Fisher's LDA:

A comparison of principal component analysis with Fisher's linear discriminant for linear dimensionality reduction. Here the data in two dimensions, belonging to two classes shown in red and blue, is to be projected onto a single dimension. PCA chooses the direction of maximum variance, shown by the magenta curve, which leads to strong class overlap, whereas the Fisher linear discriminant takes account of the class labels and leads to a projection onto the green curve giving much better class separation.



<http://www.face-rec.org/algorithms/PCA/jcn.pdf>

<http://www.cs.jhu.edu/~hager/Public/teaching/cs461/pami97-eigenfaces.pdf>



浙江大学

ZheJiang University



人工智能研究所

Institute of Artificial Intelligence

Probabilistic PCA



Why Probabilistic PCA?

- Probabilistic PCA vs. conventional PCA:
 - We can derive an EM algorithm for PCA that is computationally efficient in situations where only a few leading eigenvectors are required and that avoids having to evaluate the data covariance matrix as an intermediate step.
 - The combination of a probabilistic model and EM allows us to deal with missing values in the data set.
 - Probabilistic PCA forms the basis for a Bayesian treatment of PCA in which the dimensionality of the principal subspace can be found automatically from the data.
 - The existence of a likelihood function allows direct comparison with other probabilistic density models.
 - Probabilistic PCA can be used to model class-conditional densities and hence be applied to classification problems.
 - The probabilistic PCA model can be run generatively to provide samples from the distribution.



Probabilistic PCA

- PCA can also be expressed as the maximum likelihood solution of a probabilistic latent variable model.

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}) \quad p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I}) \quad \mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}$$

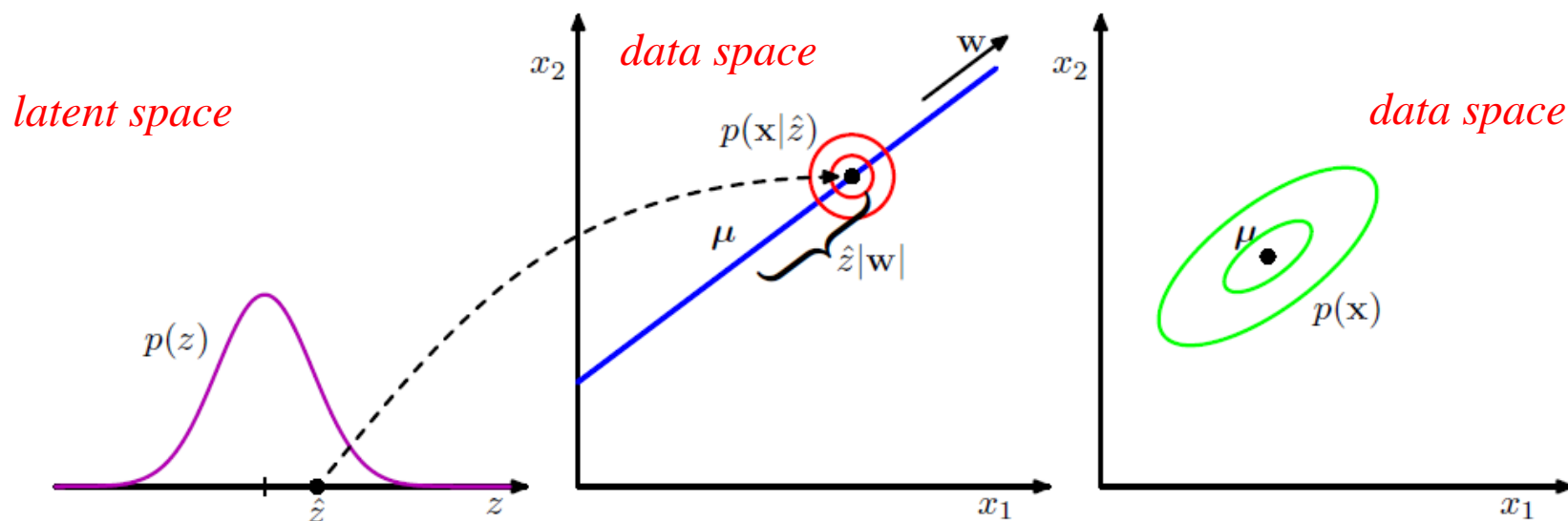


Figure 12.9 An illustration of the generative view of the probabilistic PCA model for a two-dimensional data space and a one-dimensional latent space. An observed data point \mathbf{x} is generated by first drawing a value $\hat{\mathbf{z}}$ for the latent variable from its prior distribution $p(\mathbf{z})$ and then drawing a value for \mathbf{x} from an isotropic Gaussian distribution (illustrated by the red circles) having mean $\mathbf{W}\hat{\mathbf{z}} + \boldsymbol{\mu}$ and covariance $\sigma^2\mathbf{I}$. The green ellipses show the density contours for the marginal distribution $p(\mathbf{x})$.



Probabilistic PCA

- PCA can also be expressed as the maximum likelihood solution of a probabilistic latent variable model.

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}) \quad p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I}) \quad \mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}$$

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) d\mathbf{z} \quad \Rightarrow \quad p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C}) \quad \mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$$

$$\begin{aligned} \Rightarrow \quad \mathbb{E}[\mathbf{x}] &= \mathbb{E}[\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}] = \boldsymbol{\mu} \\ \text{cov}[\mathbf{x}] &= \mathbb{E}[(\mathbf{W}\mathbf{z} + \boldsymbol{\epsilon})(\mathbf{W}\mathbf{z} + \boldsymbol{\epsilon})^T] \\ &= \mathbb{E}[\mathbf{W}\mathbf{z}\mathbf{z}^T\mathbf{W}^T] + \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T] = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I} \end{aligned} \quad \begin{aligned} \mathbf{C}^{-1} &= \sigma^{-1}\mathbf{I} - \sigma^{-2}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^T \\ \mathbf{M} &= \mathbf{W}^T\mathbf{W} + \sigma^2\mathbf{I} \end{aligned}$$

$$\Rightarrow \quad p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mathbf{M}^{-1}\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}), \sigma^{-2}\mathbf{M})$$



Maximum likelihood PCA

- We have: $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C})$ $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$

- Then,

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{n=1}^N \ln p(\mathbf{x}_n|\mathbf{W}, \boldsymbol{\mu}, \sigma^2)$$

$$= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$

$$\xrightarrow{\boldsymbol{\mu} = \bar{\mathbf{x}}} \ln p(\mathbf{X}|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = -\frac{N}{2} \{ D \ln(2\pi) + \ln |\mathbf{C}| + \text{Tr}(\mathbf{C}^{-1}\mathbf{S}) \}$$

- Solving \mathbf{W}_{ML} : $\mathbf{W}_{\text{ML}} = \mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}$

\mathbf{R} is an arbitrary $M \times M$ orthogonal matrix

- Solving σ_{ML} :
$$\sigma_{\text{ML}}^2 = \frac{1}{D-M} \sum_{i=M+1}^D \lambda_i$$

- For $M = D$:
$$\mathbf{C} = \mathbf{U}(\mathbf{L} - \sigma^2 \mathbf{I})^{1/2} \mathbf{R} \mathbf{R}^T (\mathbf{L} - \sigma^2 \mathbf{I})^{1/2} \mathbf{U}^T + \sigma^2 \mathbf{I} = \mathbf{U} \mathbf{L} \mathbf{U}^T = \mathbf{S}$$



EM algorithm for PCA

- We have:
$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{n=1}^N \{ \ln p(\mathbf{x}_n | \mathbf{z}_n) + \ln p(\mathbf{z}_n) \}$$

➡
$$\mathbb{E}[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2)] = - \sum_{n=1}^N \left\{ \frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2} \text{Tr}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T]) \right. \\ \left. + \frac{1}{2\sigma^2} \|\mathbf{x}_n - \boldsymbol{\mu}\|^2 - \frac{1}{\sigma^2} \mathbb{E}[\mathbf{z}_n]^T \mathbf{W}^T (\mathbf{x}_n - \boldsymbol{\mu}) + \frac{1}{2\sigma^2} \text{Tr}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}^T \mathbf{W}) \right\}$$

- E Step:
$$\begin{aligned} \mathbb{E}[\mathbf{z}_n] &= \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x}_n - \bar{\mathbf{x}}) \\ \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] &= \sigma^2 \mathbf{M}^{-1} + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^T \end{aligned}$$

- M Step:

$$\mathbf{W}_{\text{new}} = \left[\sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) \mathbb{E}[\mathbf{z}_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \right]^{-1}$$
$$\sigma_{\text{new}}^2 = \frac{1}{ND} \sum_{n=1}^N \left\{ \|\mathbf{x}_n - \bar{\mathbf{x}}\|^2 - 2 \mathbb{E}[\mathbf{z}_n]^T \mathbf{W}_{\text{new}}^T (\mathbf{x}_n - \bar{\mathbf{x}}) + \text{Tr}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}_{\text{new}}^T \mathbf{W}_{\text{new}}) \right\}$$



浙江大学

ZheJiang University



人工智能研究所

Institute of Artificial Intelligence

Kernel PCA



Kernel PCA

- Consider a data set $\{\mathbf{x}_n\}$ of observations, where $n = 1, \dots, N$, in a space of dimensionality D .

$$\sum_n \mathbf{x}_n = \mathbf{0} \quad \longrightarrow \quad \mathbf{S} \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad \mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \quad \mathbf{u}_i^T \mathbf{u}_i = 1$$

- Consider a nonlinear transformation $\phi(\mathbf{x})$ into an M -dimensional feature space, so that each data point \mathbf{x}_n is thereby projected onto a point $\phi(\mathbf{x}_n)$.

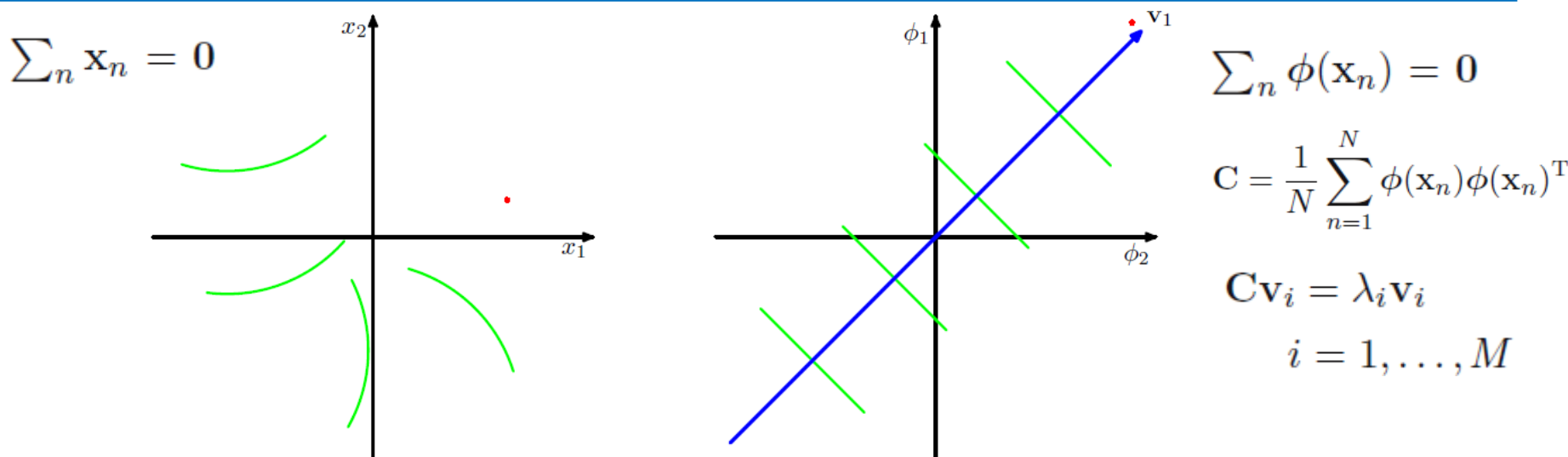


Figure 12.16 Schematic illustration of kernel PCA. A data set in the original data space (left-hand plot) is projected by a nonlinear transformation $\phi(\mathbf{x})$ into a feature space (right-hand plot). By performing PCA in the feature space, we obtain the principal components, of which the first is shown in blue and is denoted by the

Kernel PCA

- In the M-dimensional feature space:

$$\boxed{\sum_n \phi(\mathbf{x}_n) = \mathbf{0}} \quad \mathbf{C} = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \quad \mathbf{C} \mathbf{v}_i = \lambda_i \mathbf{v}_i \quad i = 1, \dots, M$$

$$\Rightarrow \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \{ \phi(\mathbf{x}_n)^T \mathbf{v}_i \} = \lambda_i \mathbf{v}_i \quad \Rightarrow \mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n)$$

- Substituting \mathbf{v}_i back into \mathbf{C} :

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \sum_{m=1}^N a_{im} \phi(\mathbf{x}_m) = \lambda_i \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n)$$

- Multiplying both sides by $\phi(\mathbf{x}_l)^T$:

$$\xrightarrow{k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)} \frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_l, \mathbf{x}_n) \sum_{m=1}^N a_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N a_{in} k(\mathbf{x}_l, \mathbf{x}_n)$$

$$\Rightarrow \mathbf{K}^2 \mathbf{a}_i = \lambda_i N \mathbf{K} \mathbf{a}_i \quad \Rightarrow \mathbf{K} \mathbf{a}_i = \lambda_i N \mathbf{a}_i$$

$$\Rightarrow 1 = \mathbf{v}_i^T \mathbf{v}_i = \sum_{n=1}^N \sum_{m=1}^N a_{in} a_{im} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = \mathbf{a}_i^T \mathbf{K} \mathbf{a}_i = \lambda_i N \mathbf{a}_i^T \mathbf{a}_i$$

$$\Rightarrow y_i(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x})^T \phi(\mathbf{x}_n) = \sum_{n=1}^N a_{in} k(\mathbf{x}, \mathbf{x}_n)$$

Kernel PCA

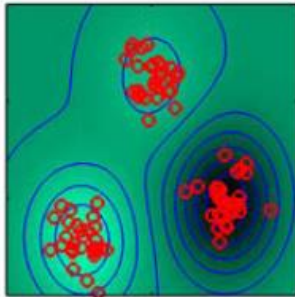
- When has not zero mean in the feature space:

$$\tilde{\phi}(\mathbf{x}_n) = \phi(\mathbf{x}_n) - \frac{1}{N} \sum_{l=1}^N \phi(\mathbf{x}_l)$$

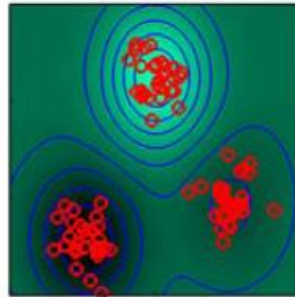
$$\begin{aligned} \tilde{K}_{nm} &= \tilde{\phi}(\mathbf{x}_n)^T \tilde{\phi}(\mathbf{x}_m) \\ &= \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) - \frac{1}{N} \sum_{l=1}^N \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_l) - \frac{1}{N} \sum_{l=1}^N \phi(\mathbf{x}_l)^T \phi(\mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{l=1}^N \phi(\mathbf{x}_j)^T \phi(\mathbf{x}_l) \\ &= k(\mathbf{x}_n, \mathbf{x}_m) - \frac{1}{N} \sum_{l=1}^N k(\mathbf{x}_l, \mathbf{x}_m) - \frac{1}{N} \sum_{l=1}^N k(\mathbf{x}_n, \mathbf{x}_l) + \frac{1}{N^2} \sum_{j=1}^N \sum_{l=1}^N k(\mathbf{x}_j, \mathbf{x}_l). \end{aligned}$$

➡ $\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{1}_N \mathbf{K} - \mathbf{K} \mathbf{1}_N + \mathbf{1}_N \mathbf{K} \mathbf{1}_N$ $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 0.1)$

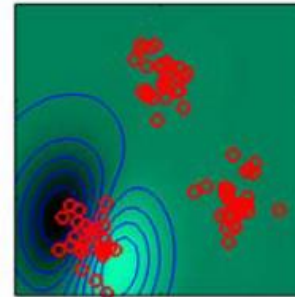
Eigenvalue=21.72



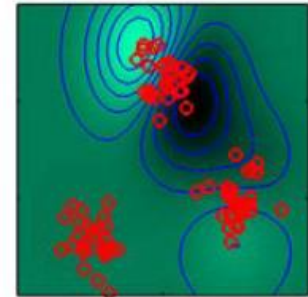
Eigenvalue=21.65



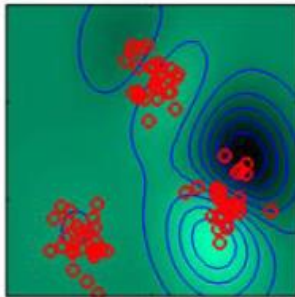
Eigenvalue=4.11



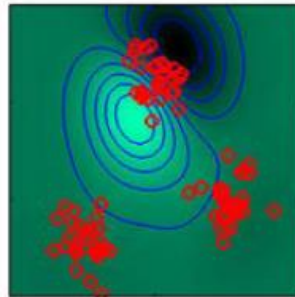
Eigenvalue=3.93



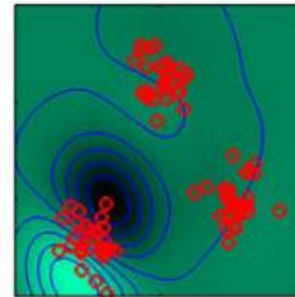
Eigenvalue=3.66



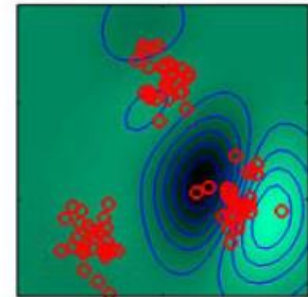
Eigenvalue=3.09



Eigenvalue=2.60



Eigenvalue=2.53





浙江大学

ZheJiang University



人工智能研究所

Institute of Artificial Intelligence

Nonlinear Latent Variable Models

(*Modeling nonlinear manifolds / manifold learning*)



Modeling nonlinear manifolds

- Two nonprobabilistic methods for dimensionality reduction and data visualization:
 - *Isometric feature mapping (ISOMAP): global method*
 - project the data to a lower-dimensional space using MDS, but where the dissimilarities are defined in terms of the *geodesic distances* measured along the manifold.
 - *Locally linear embedding (LLE): local method*
 - Map the high-dimensional data points down to a lower dimensional space while preserving coefficients.

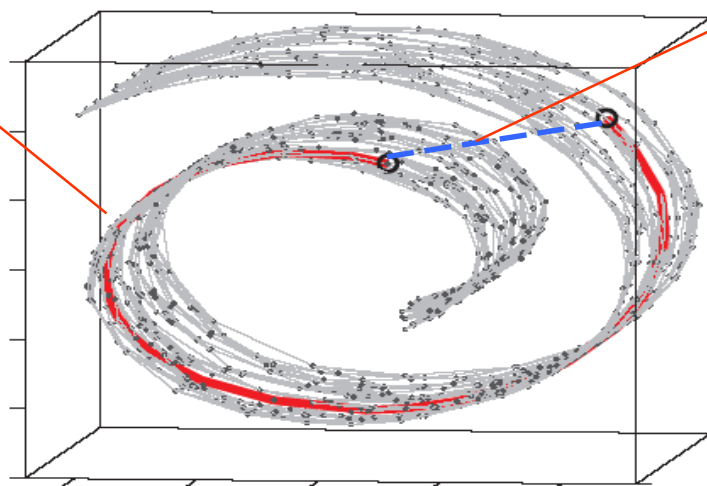
-
1. Tenenbaum J.B., Silva V. De , Langford J. C., A global geometric framework for nonlinear dimensionality reduction, *Science*, 2000, 290 (5500): 2219-2323
 2. Sam Roweis, Lawrence Saul, Nonlinear dimensionality reduction by locally linear embedding, *Science*, 2000, 290(5500):2323-2326



Problem description

- Given a set $\mathbf{x}_1, \dots, \mathbf{x}_k$ of k points in \mathbf{R}^n , find a set of points $\mathbf{y}_1, \dots, \mathbf{y}_k$ in \mathbf{R}^m ($m \ll n$) such that \mathbf{y}_i “represents” \mathbf{x}_i as accurately as possible.
- If the data \mathbf{x}_i is placed in a *super plane* in high dimensional space, the traditional algorithms, such as PCA, work well.
- However, when the data \mathbf{x}_i is placed in a *nonlinear manifold* in high dimensional space, then the linear models can not work any more.
 - A nonlinear manifold can be roughly understood as a distorted super plane, which may be twisted, folded, or curved.

True distance

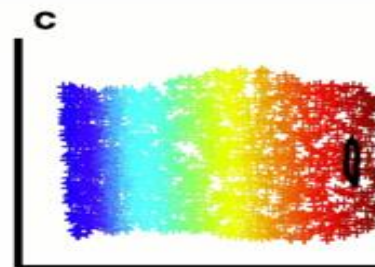
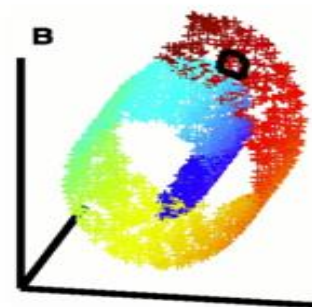
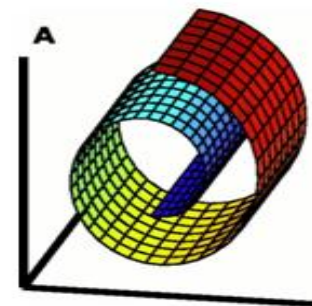
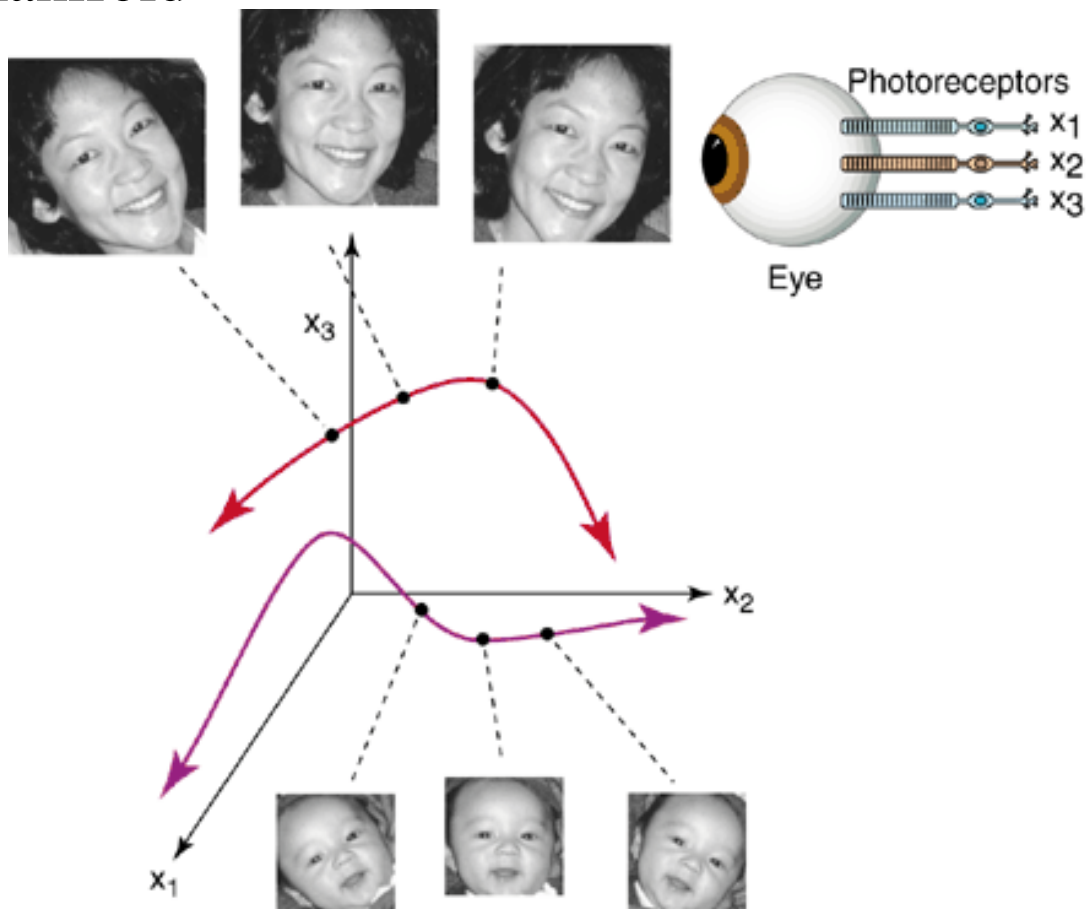


Euclidean distance



Manifold examples

- Embed data in a higher dimensional space to a lower dimensional manifold





Isometric feature mapping (ISOMAP)

- Idea:
 - Estimate the geodesic distance between faraway points.
 - For neighboring points Euclidean distance is a good approximation to the geodesic distance.
 - For farway points estimate the distance by a series of short hops between neighboring points.
 - Find shortest paths in a graph with edges connecting neighboring data points



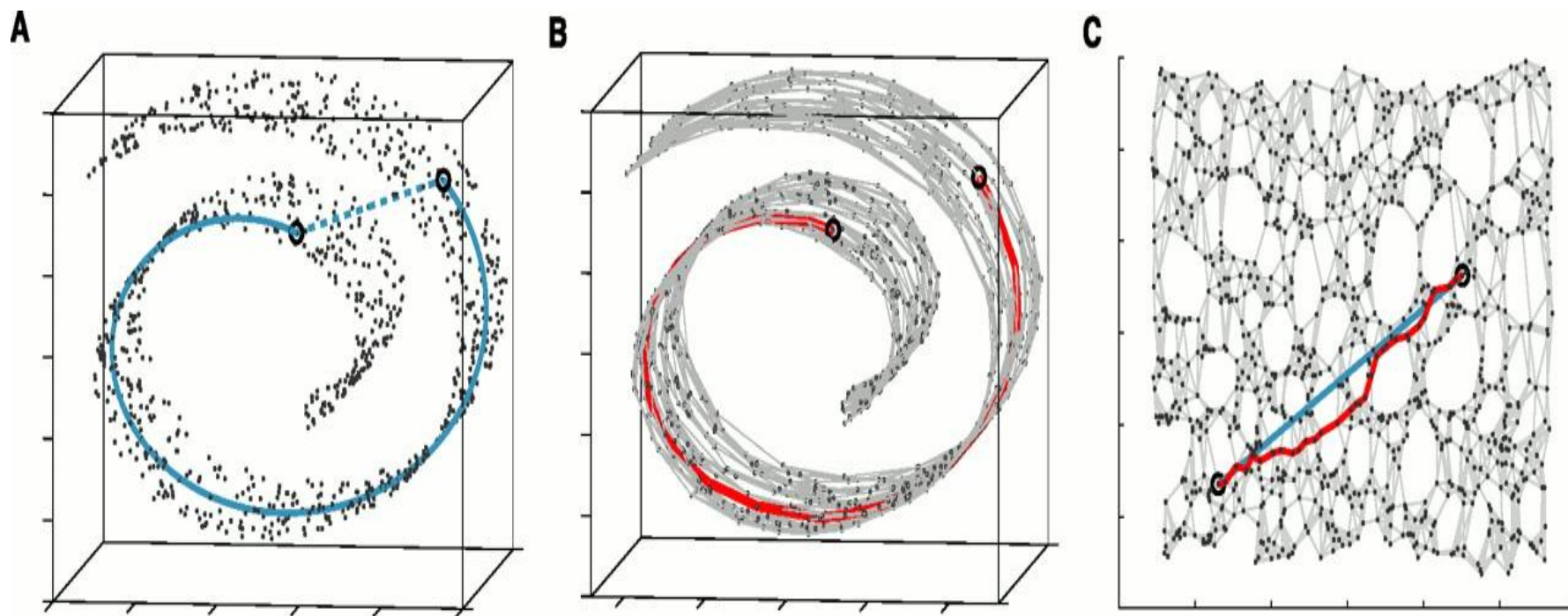
Isometric feature mapping (ISOMAP)

- Algorithm:
 - ① Determine the neighbors.
 - All points in a fixed radius.
 - K nearest neighbors
 - ② Construct a neighborhood graph.
 - Each point is connected to the other if it is a K nearest neighbor.
 - Edge Length equals the Euclidean distance
 - ③ Compute the shortest paths between two nodes
 - Floyd's Algorithm
 - Djkastra's AAlgorithm
 - ④ Construct a lower dimensional embedding.
 - Classical MDS



Iso**metric** feature **map**ping (ISOMAP)

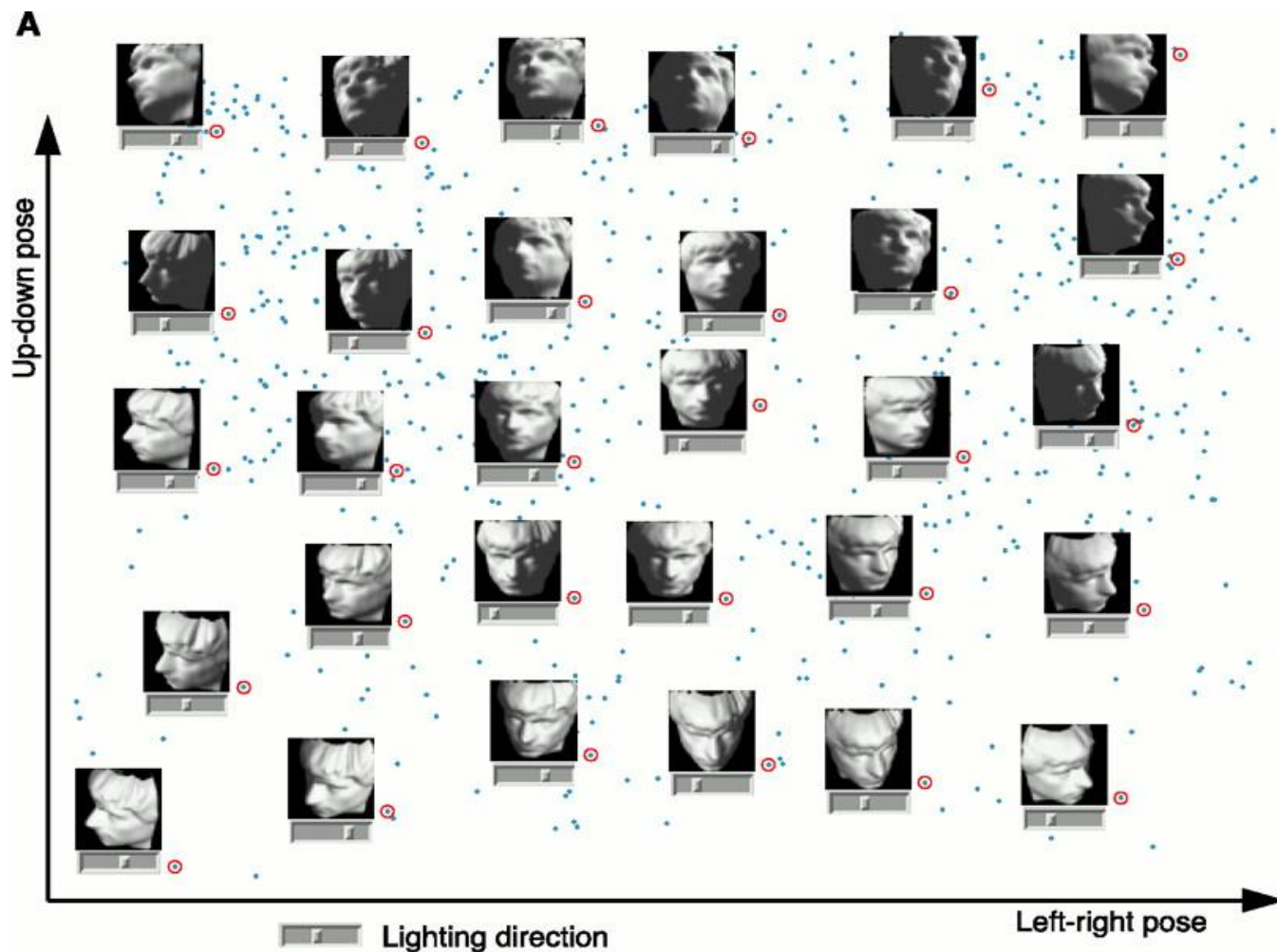
- Algorithm:





Iso**metric** feature **map**ping (ISOMAP)

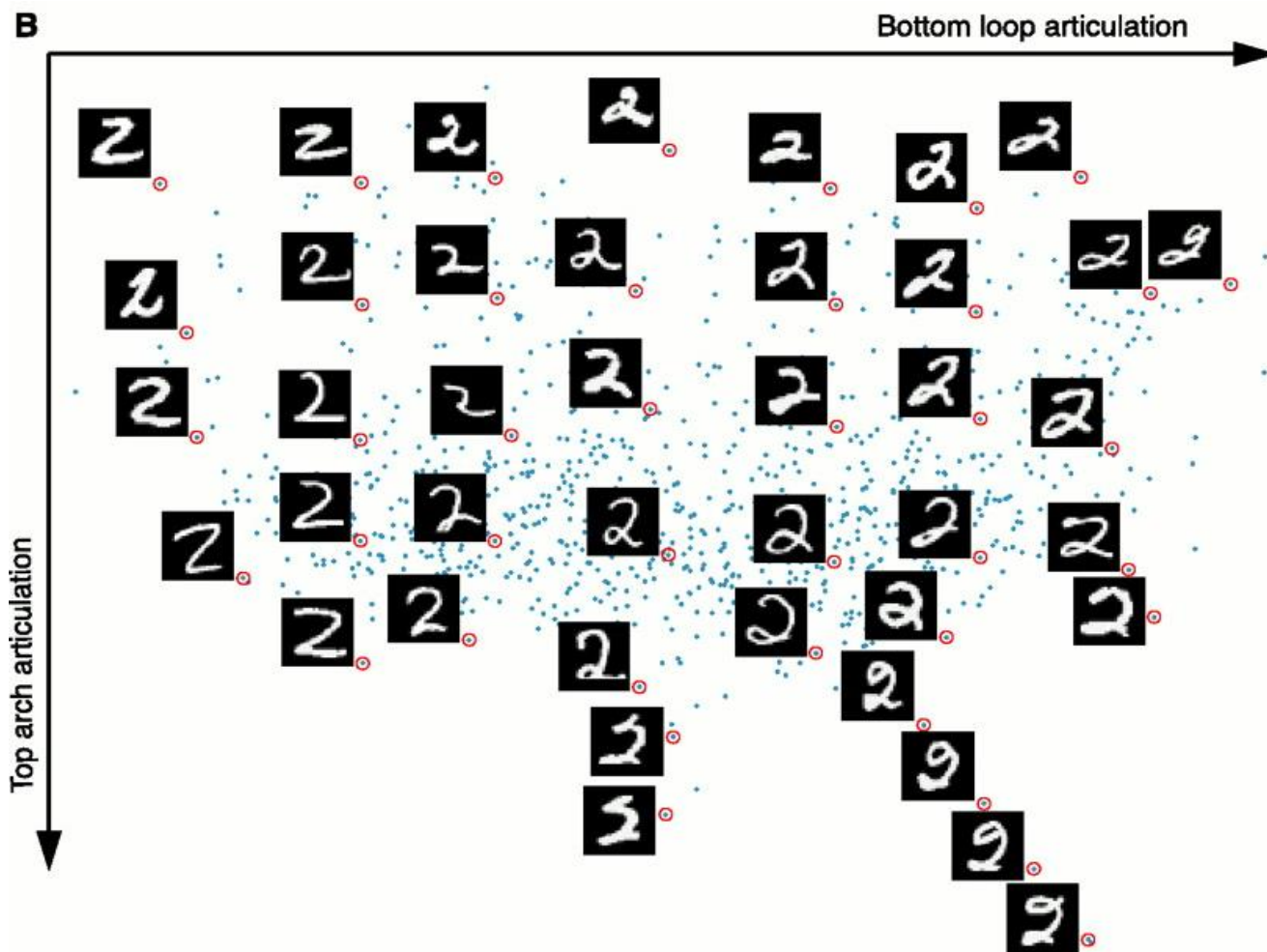
- Results:





Iso_{metric} feature _{map}ping (ISOMAP)

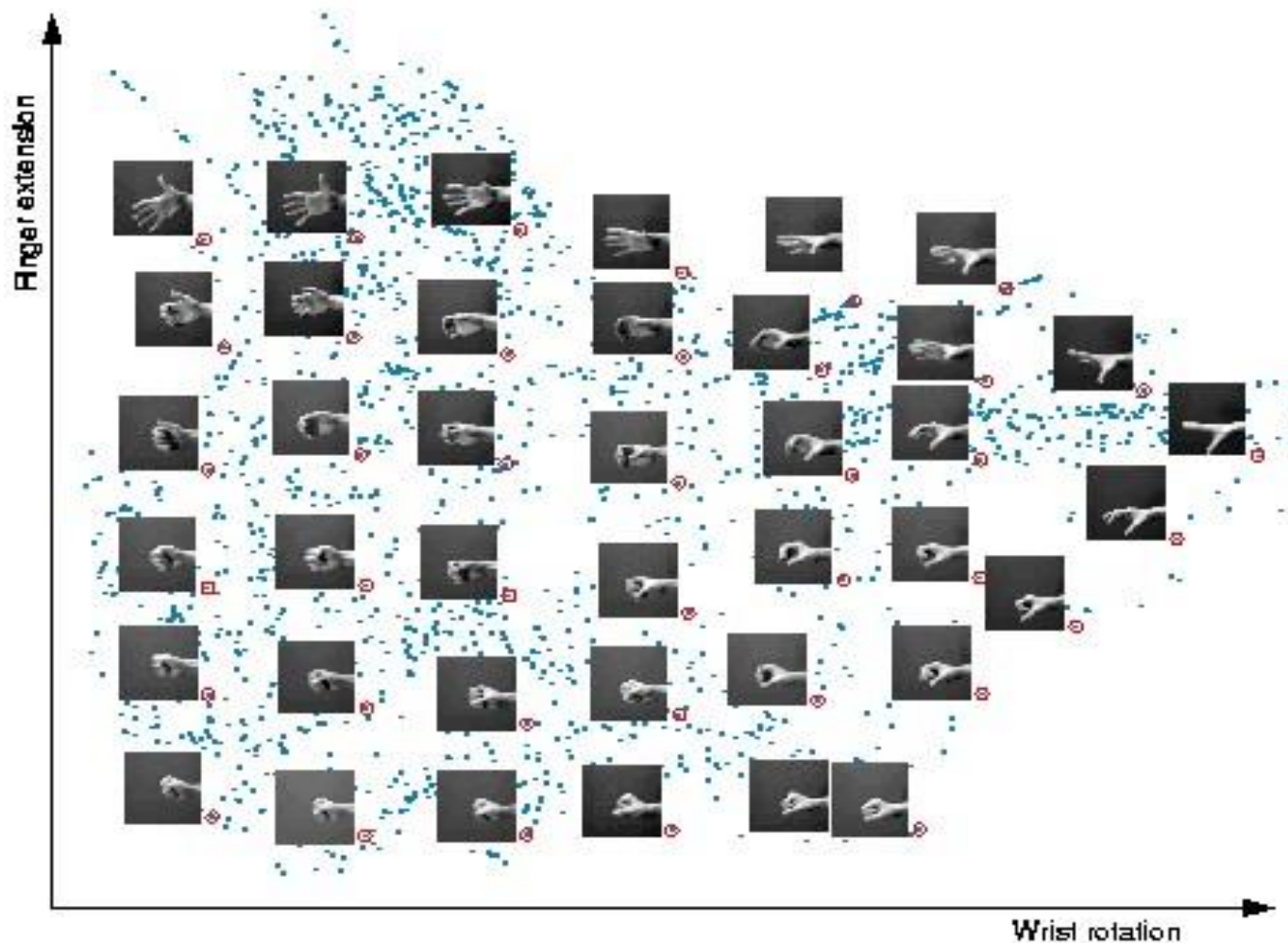
- Results:





Iso**metric** feature **map**ping (ISOMAP)

- Results:



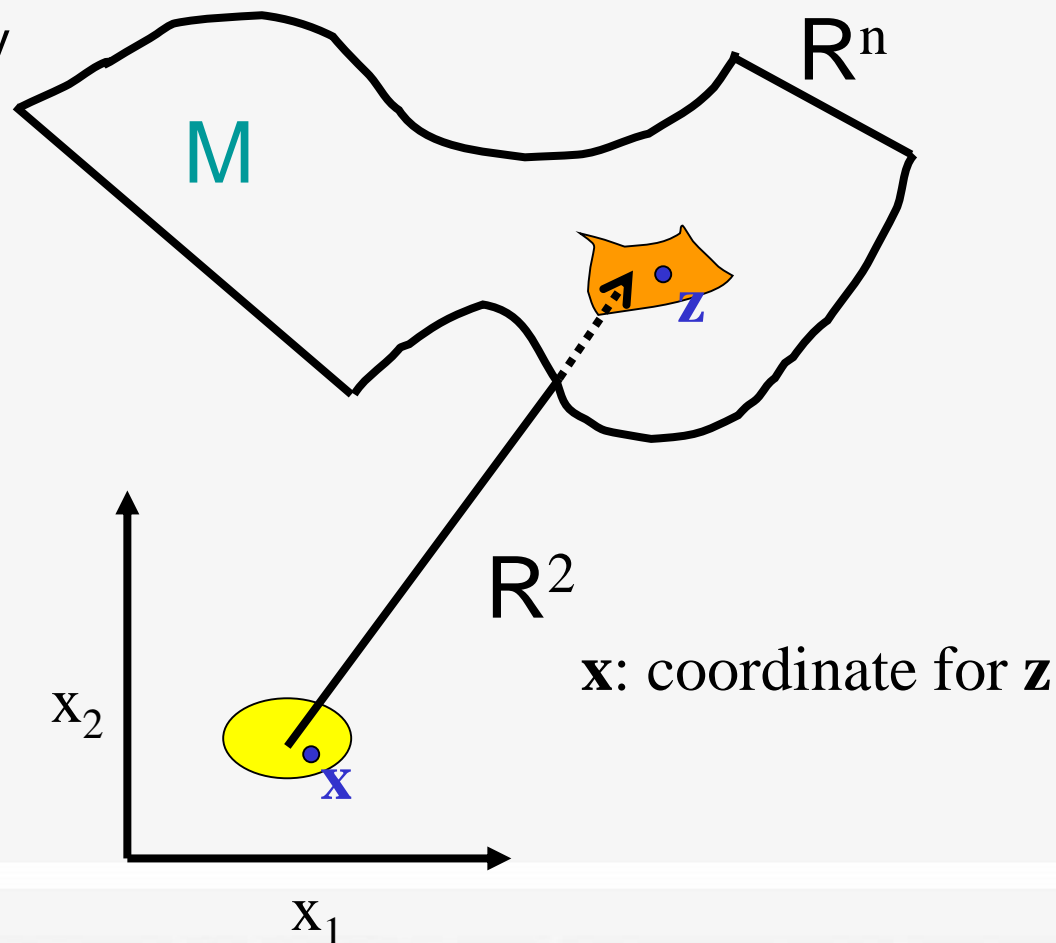


Locally Linear Embedding (LLE)

- Idea:
 - Fit Locally , Think Globally

Locally it is a linear patch

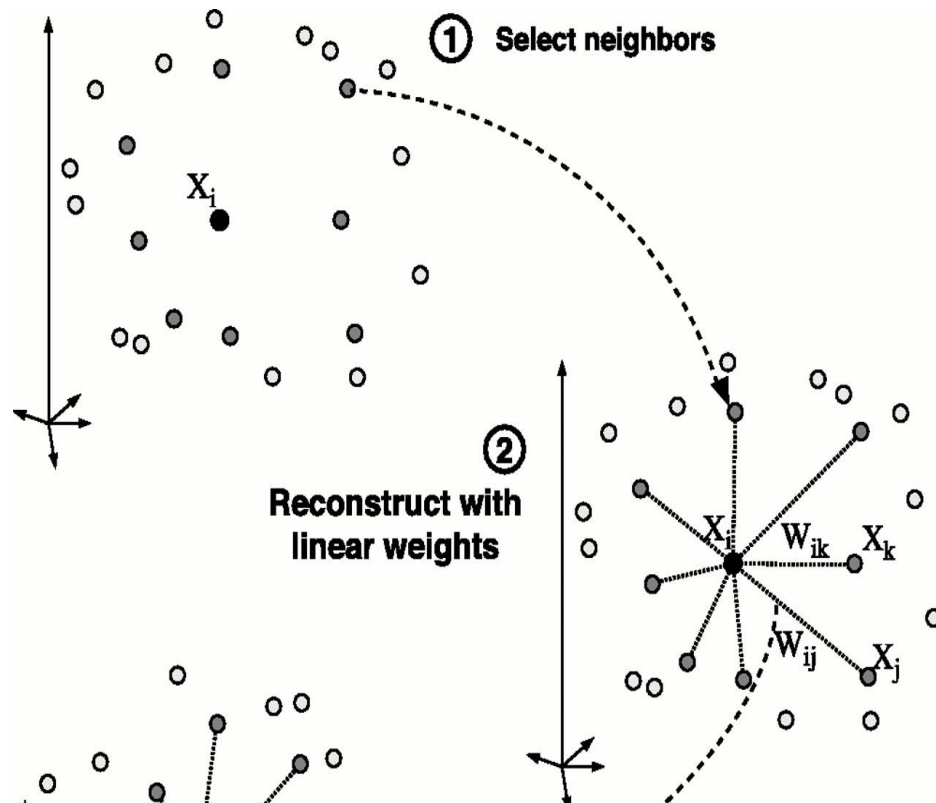
Key: how to combine all local patches together?





Locally Linear Embedding (LLE)

- Algorithm:
 - Fit Locally



We expect each data point and its neighbors to lie on or close to a locally linear patch of the manifold.

Each point can be written as a linear combination of its neighbors.
The weights chosen to minimize the reconstruction Error.

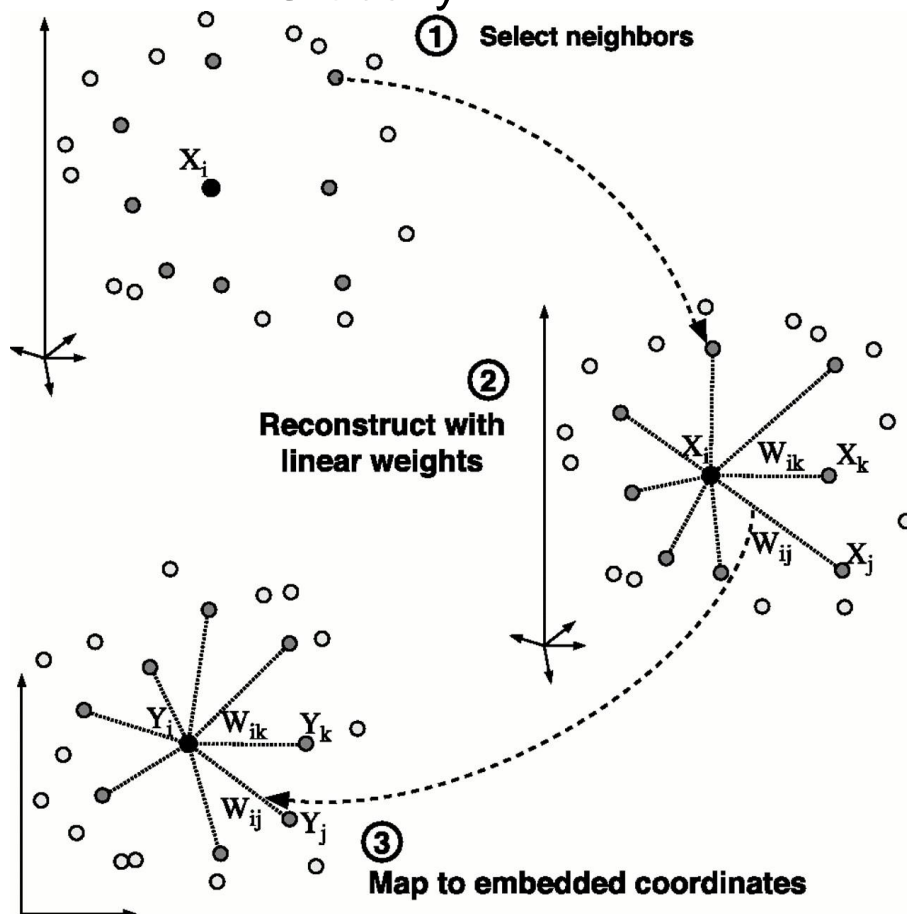
$$\min_W \| X_i - \sum_{j=1}^K W_{ij} X_j \|^2 \quad (1)$$



Locally Linear Embedding (LLE)

- Algorithm:

- Think Globally



$$\min_W \| X_i - \sum_{j=1}^K W_{ij} X_j \|^2 \quad (1)$$



Low-dimensional embedding:

$$Y_{d \times N} = [Y_1 | Y_2 | \dots | Y_N]$$

Use the same weights from the original space:

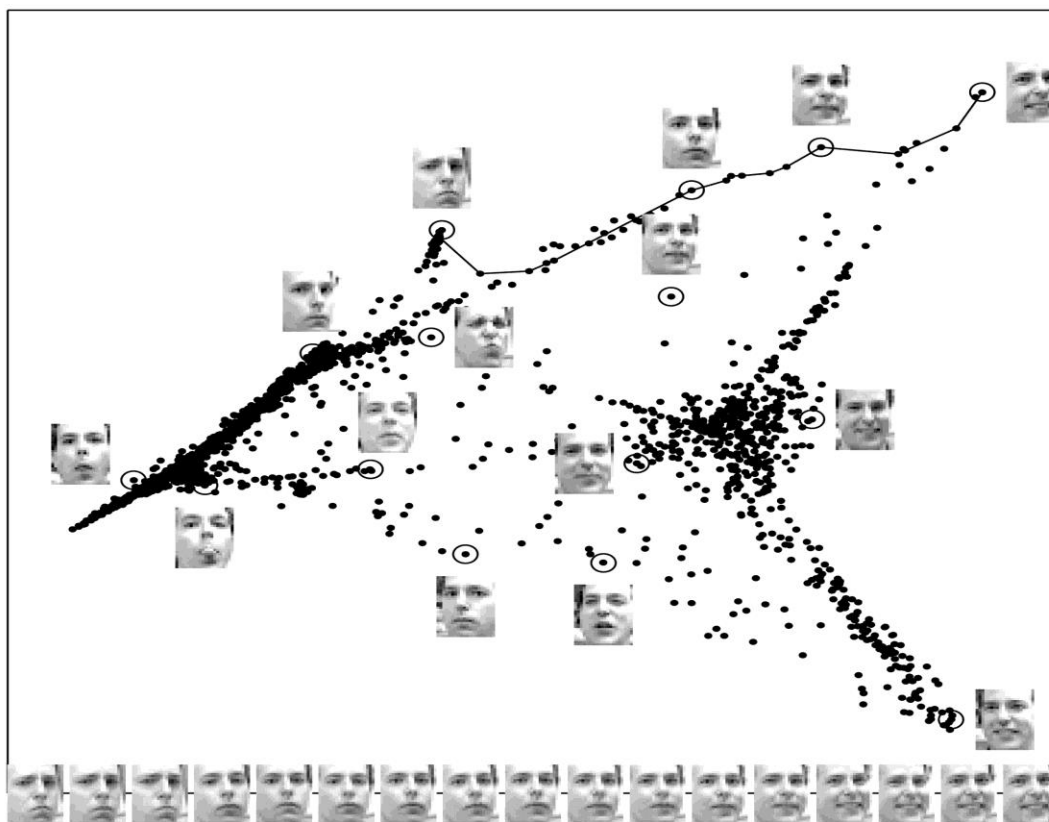
$$\min_Y \sum_{i=1}^N \| Y_i - Y W_i \|^2$$

***Y** is given by the eigenvectors of the lowest d non-zero eigenvalues of the matrix $\mathbf{M} = (\mathbf{I} - \mathbf{W})^T (\mathbf{I} - \mathbf{W})$*



Locally Linear Embedding (LLE)

- **Results:** Images of faces mapped into the embedding space described by the first two coordinates of LLE. Representative faces are shown next to circled points. The bottom images correspond to points along the top-right path (linked by solid line) illustrating one particular mode of variability in pose and expression.





Locally Linear Embedding (LLE)

- Limitations:
 - require dense data points on the manifold for good estimation
 - A good neighborhood seems essential to their success
 - How to choose k ?
 - Too few neighbors
 - Result in rank deficient tangent space and lead to over-fitting
 - Too many neighbors
 - Tangent space will not match local geometry well



Locally Linear Embedding (LLE)

- Limitations:

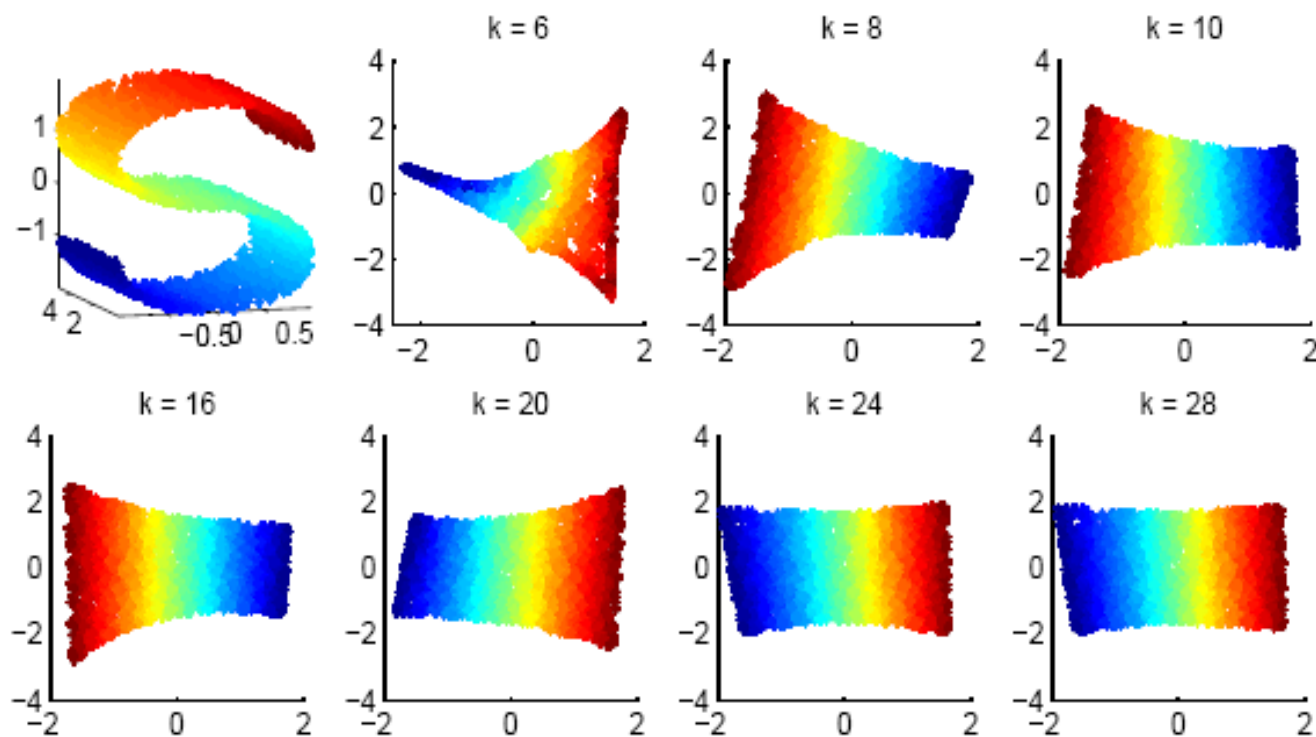


FIG. 5. S-curve (top left) and computed 2D coordinates by LLE with various neighborhood size k .

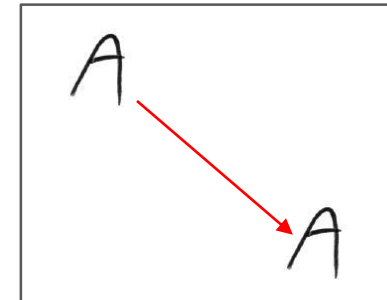
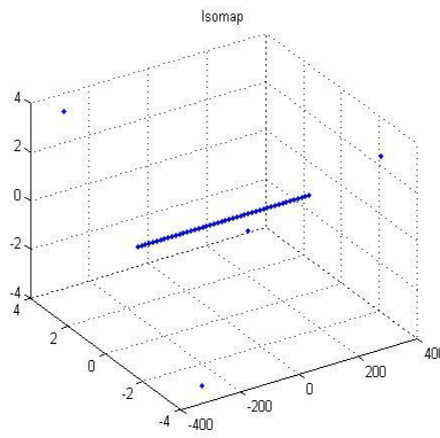
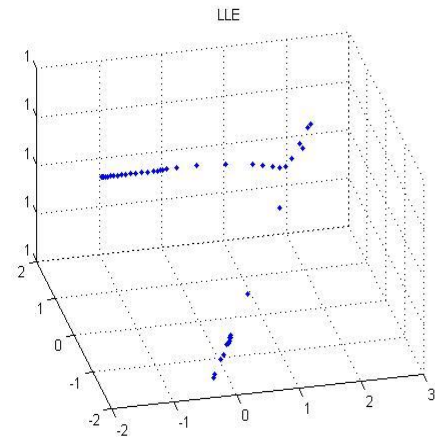
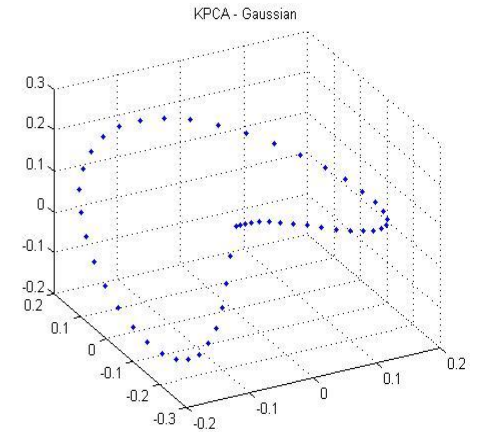
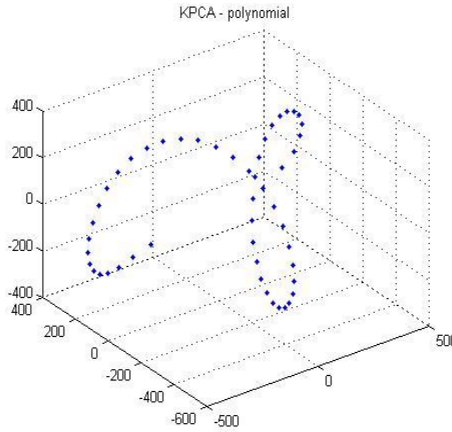
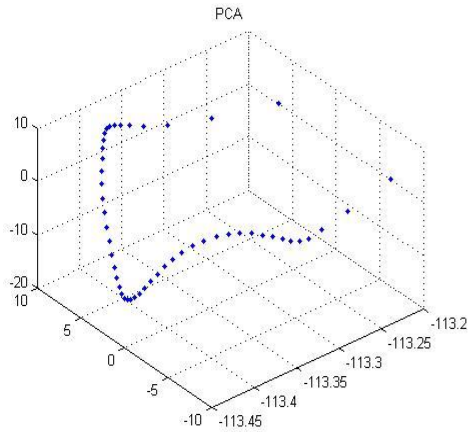


ISOMAP vs. LLE

ISOMAP	LLE
Do MDS on the geodesic distance matrix.	Model local neighborhoods as linear patches and then embed in a lower dimensional manifold.
Global approach	Local approach
Dynamic programming approaches	Computationally efficient..sparse matrices
Convergence limited by the manifold curvature and number of points.	Good representational capacity



Experimental Results





浙江大学

ZheJiang University



人工智能研究所

Institute of Artificial Intelligence

Next: Combining models