

Summary

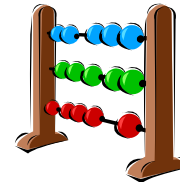
■ There are 3 topics in learning research

- Understanding human brains
- Developing learning machines
- Clarifying essence of learning mathematically



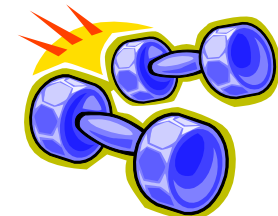
■ There are 3 types of learning.

- Supervised learning
- Unsupervised learning
- Reinforcement learning



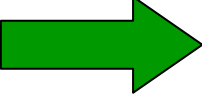
■ Topics of unsupervised learning:

- Dimensionality reduction
- Data clustering
- Blind source separation
- Outlier/novelty detection



Dimensionality Reduction

$$\{\mathbf{x}_i\}_{i=1}^n, \quad \mathbf{x}_i \in \mathbb{R}^d, \quad d \gg 1$$

- High-dimensional data is too complex to analyze:  “Curse of dimensionality”
- We want to reduce the dimensionality of the data while preserving the intrinsic “information” in the data.
- Dimensionality reduction is also called
 - Embedding
 - Data visualization (if the dimension is reduced up to 3)

What Kind of “Information” Do We Want to Preserve?

- There are several criteria: For example, embed the data such that
 - Original data is preserved
 - Local structure is maintained
 - Interesting components are extracted
- Which criterion is the best?
- It depends on the data, purpose...

Methods to be Dealt With

- Original data is preserved:
 - Principal component analysis (PCA)
 - Kernel PCA
- Local structure is maintained:
 - Locality preserving projection
 - Laplacian eigenmap embedding
- Interesting components are extracted:
 - Projection pursuit
 - Non-Gaussian component analysis

Notation

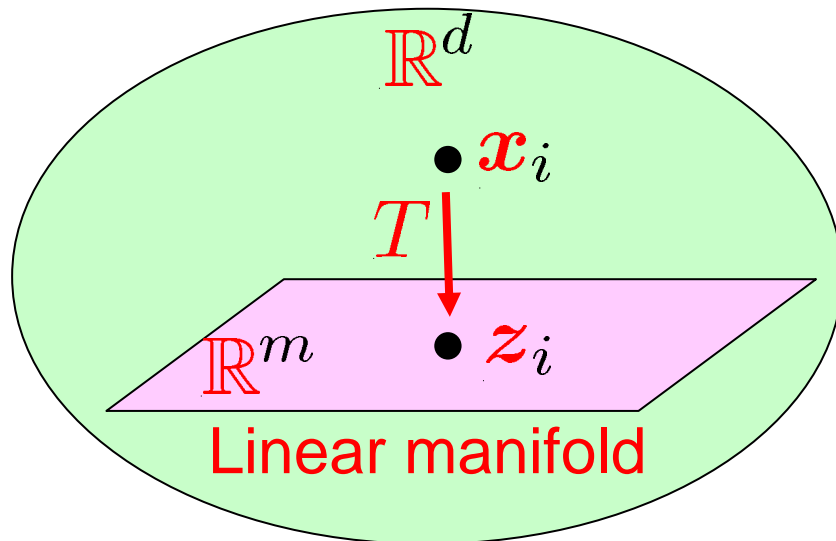
■ **Data:** $\{\mathbf{x}_i\}_{i=1}^n$, $\mathbf{x}_i \in \mathbb{R}^d$, $d \gg 1$

■ **Mapping:** $T : \mathbb{R}^d \rightarrow \mathbb{R}^m$, $1 \leq m \ll d$

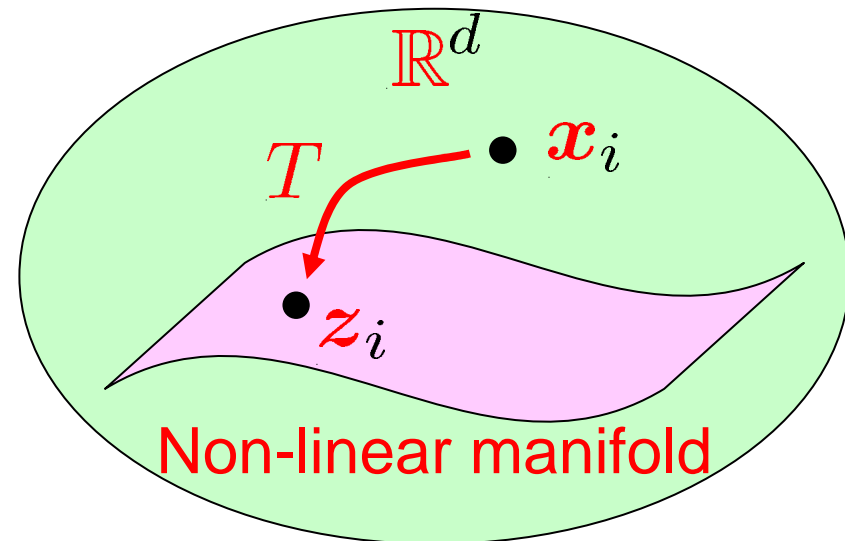
(Note: T is generally non-linear)

■ **Embedded data:** $\{\mathbf{z}_i\}_{i=1}^n$, $\mathbf{z}_i = T\mathbf{x}_i \in \mathbb{R}^m$

Linear embedding



Non-linear embedding



Preserving Original Data

- We want to embed the data so that embedded data is as “close” to the original data as possible.

- $d = 2, m = 1$

$$\{\mathbf{x}_i\}_{i=1}^n = \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0.1 \\ 2 \end{pmatrix}, \begin{pmatrix} -0.1 \\ 3 \end{pmatrix} \right\}$$

- For better description of data, throwing the first element away would be good.

- How about

$$\{\mathbf{x}_i\}_{i=1}^n = \left\{ \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \begin{pmatrix} 2 \\ 4.1 \end{pmatrix}, \begin{pmatrix} 3 \\ 5.9 \end{pmatrix} \right\}$$

Principal Component Analysis

9

- Principal component analysis (PCA) tries to automatically find the best description of data.

- Assume

- Mapping $T : \mathbf{x}_i \longrightarrow \mathbf{z}_i$ is linear.
- Data is centered:

$$\mathbf{x}_i \longleftarrow \mathbf{x}_i - \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$$

- Linear mapping is essentially a projection onto a subspace.
- Find the **subspace** which best describes data!

PCA Criterion

- Projection: $\tilde{x}_i = Px_i$
- “Closeness” : Squared Euclidean distance

$$\|\tilde{x}_i - x_i\|^2$$

- PCA criterion:

$$\min_{P \in \mathcal{P}_m} \left[\sum_{i=1}^n \|\tilde{x}_i - x_i\|^2 \right]$$

- \mathcal{P}_m : Set of all orthogonal projection matrices with rank m

How to Obtain Solution

- $\{\mathbf{b}_i\}_{i=1}^m$: Orthonormal basis in a subspace

$$\mathbf{P} = \sum_{i=1}^m \mathbf{b}_i \mathbf{b}_i^\top = \mathbf{B}^\top \mathbf{B}$$

$$\mathbf{B} = (\mathbf{b}_1 | \mathbf{b}_2 | \cdots | \mathbf{b}_m)^\top$$

- PCA criterion is equivalently expressed as

$$\mathbf{B}_{PCA} = \underset{\mathbf{B} \in \mathbb{R}^{m \times d}}{\operatorname{argmin}} [J_{PCA}(\mathbf{B})]$$

$$\text{subject to } \mathbf{B}\mathbf{B}^\top = \mathbf{I}_m$$

$$J_{PCA}(\mathbf{B}) = \sum_{i=1}^n \|\mathbf{B}^\top \mathbf{B} \mathbf{x}_i - \mathbf{x}_i\|^2$$

How to Obtain Solution (cont.)¹²

■ $J_{PCA}(\mathbf{B}) = -\text{tr}(\mathbf{B}^\top \mathbf{B} \mathbf{C}) + \text{tr}(\mathbf{C})$

$$\mathbf{C} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$$

■ PCA criterion is equivalently expressed as

$$\mathbf{B}_{PCA} = \underset{\mathbf{B} \in \mathbb{R}^{m \times d}}{\text{argmax}} \left[\text{tr}(\mathbf{B}^\top \mathbf{B} \mathbf{C}) \right]$$

subject to $\mathbf{B} \mathbf{B}^\top = \mathbf{I}_m$

Lemma

■ $B : m \times d, (1 \leq m \leq d)$

■ $C : d \times d$, positive, symmetric

■ **Problem:**

$$B_{max} = \operatorname{argmax}_{B \in \mathbb{R}^{m \times d}} \left[\operatorname{tr}(BCB^\top) \right]$$

subject to $BB^\top = I_m$

■ **Solution:**

$$B_{max} = (\psi_1 | \psi_2 | \cdots | \psi_m)^\top$$

■ $\{\lambda_i, \psi_i\}_{i=1}^m$: Eigenvalues and eigenvectors of $C\psi = \lambda\psi$

$$(\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d) \quad (\langle \psi_i, \psi_j \rangle = \delta_{i,j})$$

Proof

■ Lagrangian:

$$L(B, \Lambda) = \text{tr}(BCB^\top) - \text{tr}((BB^\top - I)\Lambda)$$

(Note: Λ is symmetric)

■ Stationary point:

$$\frac{\partial L}{\partial B} = 2BC - 2\Lambda B = 0$$

$$\Rightarrow CB^\top = B^\top \Lambda$$

$$\Rightarrow \mathcal{R}(CB^\top) = \mathcal{R}(B^\top \Lambda) \subset \mathcal{R}(B^\top)$$

Proof (cont.)

■ $\mathcal{R}(CB^\top) \subset \mathcal{R}(B^\top)$

➡ $\mathcal{R}(B^\top)$ is an invariant subspace of C

➡ $\mathcal{R}(B^\top) = \text{span}(\{\psi_{k_j}\}_{j=1}^m)$

■ $BB^\top = I_m$

➡ $\text{rank}(B^\top) = m$

➡ All $\{k_j\}_{j=1}^m$ are distinct

Proof (cont.)

- $B^\top B$ is the orthogonal projection onto $\text{span}(\{\psi_{k_j}\}_{j=1}^m)$ because

$$(B^\top B)^2 = B^\top B B^\top B = B^\top B$$

$$(B^\top B)^\top = B^\top B$$

- Since $\{\psi_{k_j}\}_{j=1}^m$ are orthonormal,

$$B^\top B = \sum_{j=1}^m \psi_{k_j} \psi_{k_j}^\top$$

Proof (cont.)

■ Eigendecomposition:

$$\mathbf{C} = \sum_{i=1}^d \lambda_i \boldsymbol{\psi}_i \boldsymbol{\psi}_i^\top$$

$$\text{tr}(\mathbf{B}\mathbf{C}\mathbf{B}^\top) = \text{tr}(\mathbf{C}\mathbf{B}^\top\mathbf{B}) = \sum_{j=1}^m \lambda_{k_j}$$

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$$

 $k_j = j$ gives a solution

PCA Embedding

$$B_{PCA} = \underset{B \in \mathbb{R}^{m \times d}}{\operatorname{argmax}} \left[\operatorname{tr}(B^\top BC) \right]$$

subject to $BB^\top = I_m$

- $\{\lambda_i, \psi_i\}_{i=1}^m$: Eigenvalues and eigenvectors of $C\psi = \lambda\psi$
 $(\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n) \quad (\langle \psi_i, \psi_j \rangle = \delta_{i,j})$

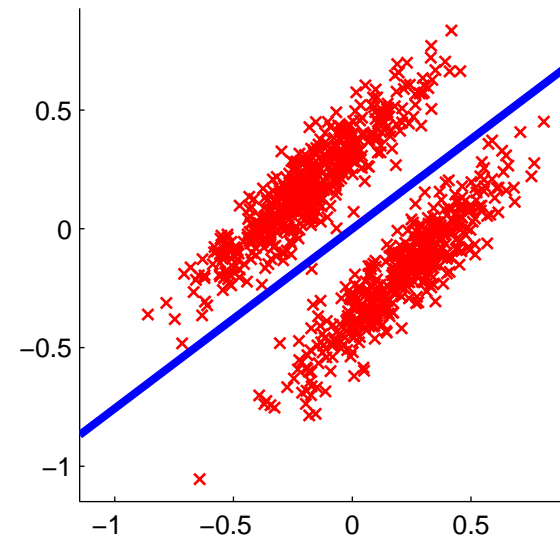
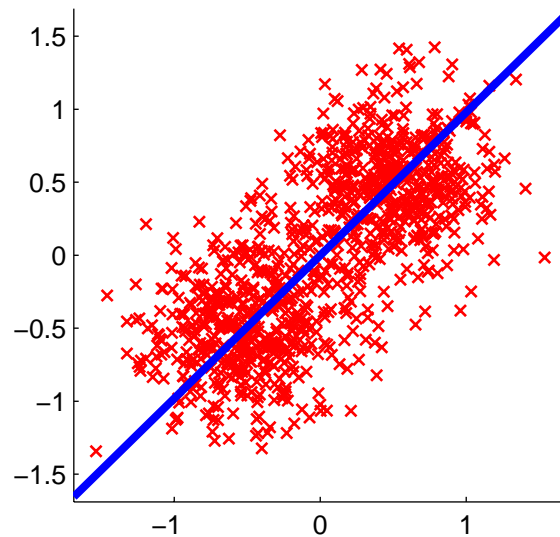
■ PCA solution:

$$B_{PCA} = (\psi_1 | \psi_2 | \dots | \psi_m)^\top$$

- PCA Embedding of a sample x' :

$$z' = B_{PCA}x'$$

Examples



- Data is well described
- PCA is intuitive, easy to implement, analytic solution available, and fast.
- However, PCA does not necessarily preserve interesting information such as clusters.

Homework

- $B : m \times d, (1 \leq m \leq d)$
- $C, D : d \times d$, positive, symmetric
- $B_{min} = \operatorname{argmin}_{B \in \mathbb{R}^{m \times d}} \left[\operatorname{tr}(BCB^\top) \right]$
subject to $BDB^\top = I_m$

■ **Prove:**

$$B_{min} = (\psi_{d-m+1} | \psi_{d-m+2} | \cdots | \psi_d)^\top$$

- $\{\lambda_i, \psi_i\}_{i=1}^m$: **Generalized** eigenvalues and eigenvectors of $C\psi = \lambda D\psi$

$$(\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d) \quad (\langle D\psi_i, \psi_j \rangle = \delta_{i,j})$$

Homework (cont.)

- Read the following article for upcoming classes:
- X. He & P. Niyogi: Locality preserving projections, In *Advances in Neural Information Processing Systems 16*, MIT Press, Cambridge, MA, 2004.