

Advanced Data Analysis: Principal Component Analysis

Masashi Sugiyama (Computer Science)

W8E-406, sugi@cs.titech.ac.jp

<http://sugiyama-www.cs.titech.ac.jp/~sugi>

Curse of Dimensionality

3

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

- If your data samples are high-dimensional, they are often **too complex to directly analyze**.
- Usual geometric intuitions are often only applicable to low-dimensional spaces; such intuitions could be even misleading in high-dimensional spaces.

Curse of Dimensionality (cont.) ⁴

- When the dimensionality increases,
 - Volume of unit hyper-cube V_c is always 1.
 - Volume of inscribed hyper-sphere V_s goes to 0.
- Relative size of hyper-sphere gets small!

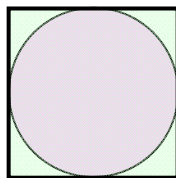
$$\frac{V_s}{V_c} \rightarrow 0$$

$d = 1$



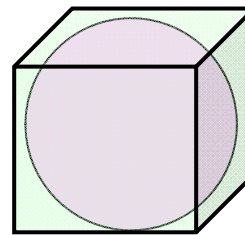
1

$d = 2$



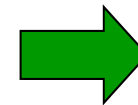
$$\pi(0.5)^2 \approx 0.79$$

$d = 3$



$$4\pi(0.5)^3/3 \approx 0.52$$

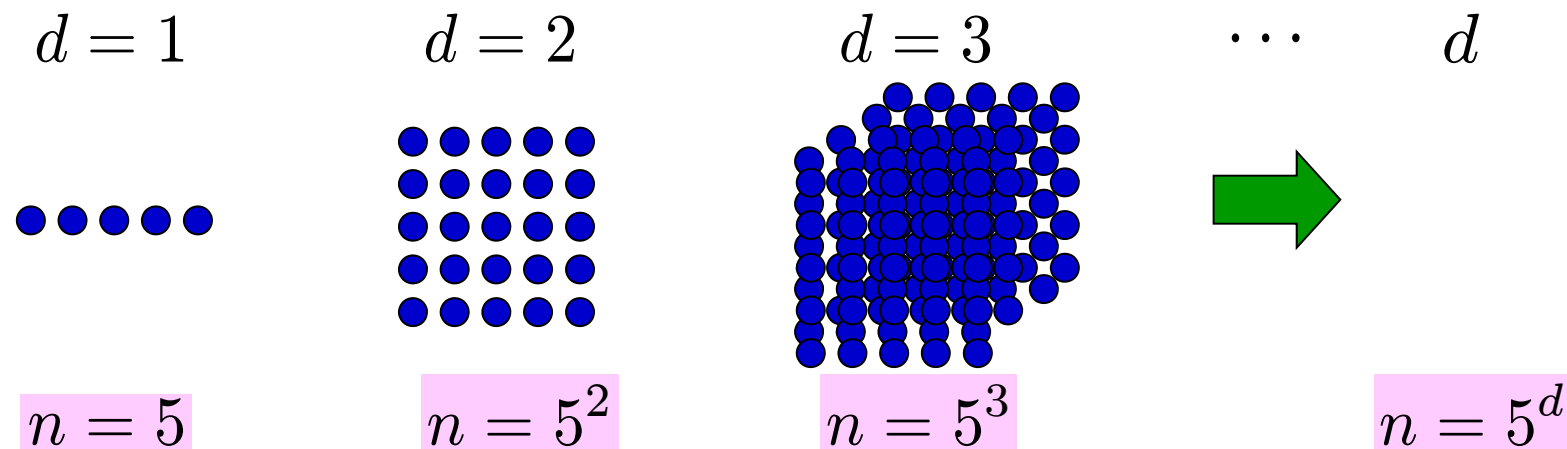
$\dots d = \infty$



0

Curse of Dimensionality (cont.) ⁵

- Grid sampling requires an exponentially large number.



- Unless you have an exponentially large number of samples, your high-dimensional samples are **never dense**.

Dimensionality Reduction

6

- We want to reduce the dimensionality of the data while preserving the intrinsic “**information**” in the data.
- Dimensionality reduction is also called **embedding**; if the dimension is reduced up to 3, it is also called **data visualization**.
- **Basic assumption (or belief)** behind dimensionality reduction: your high-dimensional data is redundant in some sense.

Notation: Linear Embedding

7

■ Data samples:

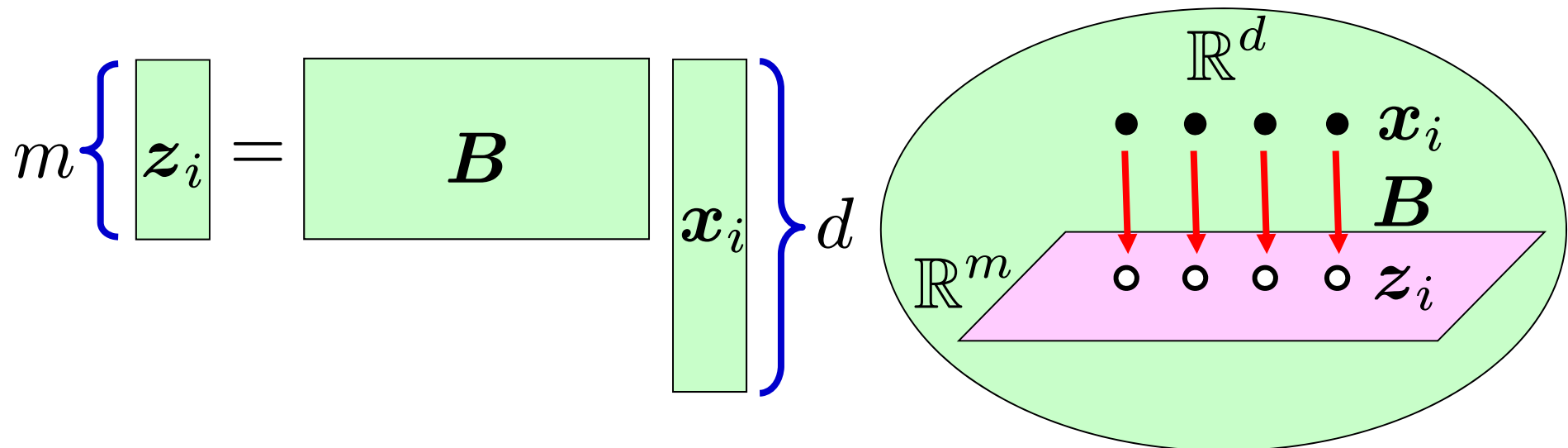
$$\{\mathbf{x}_i\}_{i=1}^n, \quad \mathbf{x}_i \in \mathbb{R}^d, \quad d \gg 1$$

■ Embedding matrix:

$$\mathbf{B} \in \mathbb{R}^{m \times d}, \quad 1 \leq m \ll d$$

■ Embedded data samples:

$$\{\mathbf{z}_i\}_{i=1}^n, \quad \mathbf{z}_i = \mathbf{B}\mathbf{x}_i \in \mathbb{R}^m$$

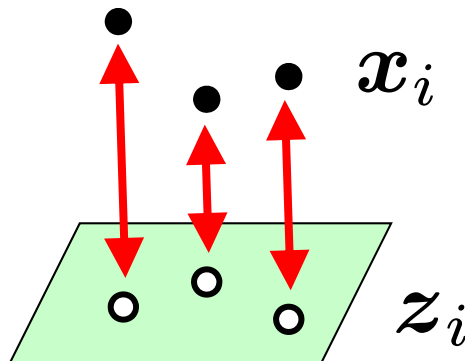


Principal Component Analysis (PCA)⁸

- **Idea**: We want to get rid of a **redundant** dimension of the data samples

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 0.1 \end{pmatrix}, \begin{pmatrix} 3 \\ -0.1 \end{pmatrix}$$

- This could be achieved by **minimizing the distance** between embedded samples and original samples.



Data Centering

9

- We **center** the data samples by

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

$$\frac{1}{n} \sum_{i=1}^n \bar{\mathbf{x}}_i = 0$$

- In matrix,

$$\bar{\mathbf{X}} = \mathbf{X} \mathbf{H}$$

$$\bar{\mathbf{X}} = (\bar{\mathbf{x}}_1 | \bar{\mathbf{x}}_2 | \cdots | \bar{\mathbf{x}}_n)$$

$$\mathbf{X} = (\mathbf{x}_1 | \mathbf{x}_2 | \cdots | \mathbf{x}_n)$$

$$\mathbf{H} = \mathbf{I}_n - \frac{1}{n} \mathbf{1}_{n \times n}$$

\mathbf{I}_n : n -dimensional identity matrix

$\mathbf{1}_{n \times n}$: $n \times n$ matrix with all ones

Orthogonal Projection

10

- $\{\mathbf{b}_i (\in \mathbb{R}^d)\}_{i=1}^m$: Orthonormal basis in m -dimensional embedding subspace

$$\langle \mathbf{b}_i, \mathbf{b}_j \rangle = \delta_{i,j} = \begin{cases} 1 & (i = j) \\ 0 & (i \neq j) \end{cases}$$

- In matrix, $\mathbf{B}\mathbf{B}^\top = \mathbf{I}_m$

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

- Orthogonal projection of $\bar{\mathbf{x}}_i$ is expressed by

$$\sum_{j=1}^m \langle \mathbf{b}_j, \bar{\mathbf{x}}_i \rangle \mathbf{b}_j \quad \left(= \mathbf{B}^\top \mathbf{B} \bar{\mathbf{x}}_i \right)$$

PCA Criterion

11

- Minimize the sum of squared distances.

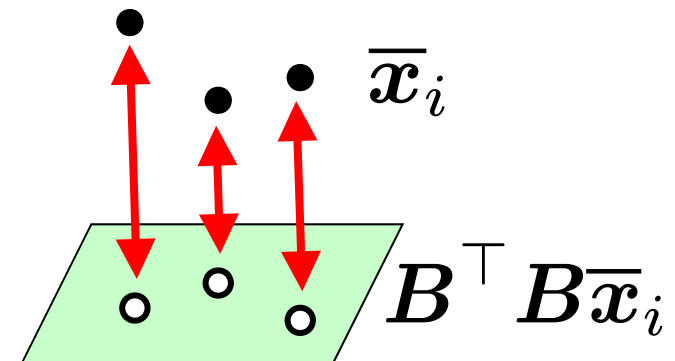
$$\sum_{i=1}^n \|B^\top B \bar{x}_i - \bar{x}_i\|^2 \quad \left(= -\text{tr}(B \bar{C} B^\top) + \text{tr}(\bar{C}) \right)$$

$$\bar{C} = \sum_{i=1}^n \bar{x}_i \bar{x}_i^\top = \bar{X} \bar{X}^\top$$

- PCA criterion:

$$B_{PCA} = \underset{B \in \mathbb{R}^{m \times d}}{\text{argmax}} \text{tr}(B \bar{C} B^\top)$$

$$\text{subject to } B B^\top = I_m$$



PCA: Summary

12

■ A PCA solution:

$$\mathbf{B}_{PCA} = (\boldsymbol{\psi}_1 | \boldsymbol{\psi}_2 | \cdots | \boldsymbol{\psi}_m)^\top$$

$\{\lambda_i, \boldsymbol{\psi}_i\}_{i=1}^m$: Sorted eigenvalues and normalized eigenvectors of $\overline{\mathbf{C}}\boldsymbol{\psi} = \lambda\boldsymbol{\psi}$

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$$

$$\langle \boldsymbol{\psi}_i, \boldsymbol{\psi}_j \rangle = \delta_{i,j}$$

■ PCA embedding of a sample \mathbf{x} :

$$\overline{\mathbf{z}} = \mathbf{B}_{PCA} \left(\mathbf{x} - \frac{1}{n} \mathbf{X} \mathbf{1}_n \right)$$

$\mathbf{1}_n$: n -dimensional vector with all ones

■ Lagrangian:

$$L(\mathbf{B}, \mathbf{\Delta}) = \text{tr}(\mathbf{B}\overline{\mathbf{C}}\mathbf{B}^\top) - \text{tr}((\mathbf{B}\mathbf{B}^\top - \mathbf{I}_m)\mathbf{\Delta})$$

$\mathbf{\Delta}$: Lagrange multipliers (symmetric)

■ Stationary point (necessary condition):

$$\bullet \frac{\partial L}{\partial \mathbf{B}} = 2\mathbf{B}\overline{\mathbf{C}} - 2\mathbf{\Delta}\mathbf{B} = 0$$

$$\longrightarrow \overline{\mathbf{C}}\mathbf{B}^\top = \mathbf{B}^\top\mathbf{\Delta} \quad (1)$$

$$\bullet \frac{\partial L}{\partial \mathbf{\Delta}} = \mathbf{B}\mathbf{B}^\top - \mathbf{I}_m = 0$$

$$\longrightarrow \mathbf{B}\mathbf{B}^\top = \mathbf{I}_m \quad (2)$$

■ Eigendecomposition:

$$\mathbf{\Delta} = \mathbf{T}\mathbf{\Gamma}\mathbf{T}^\top \quad (3)$$

\mathbf{T} : orthogonal matrix
 $\mathbf{\Gamma}$: diagonal matrix

$$\mathbf{T}^{-1} = \mathbf{T}^\top$$

Proof (cont.)

14

■ (1) & (3) $\longrightarrow \overline{C}B^\top = B^\top T\Gamma T^\top$ (4)

$\longrightarrow \overline{C}B^\top T = B^\top T\Gamma$

$\longrightarrow \overline{C}F = F\Gamma$ (5) $F = B^\top T$

■ (5) is an eigensystem

$\longrightarrow \mathcal{R}(F) = \text{span}(\{\psi_{k_i}\}_{i=1}^m)$ (6)

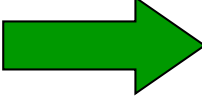
$\Gamma = \text{diag}(\lambda_{k_1}, \lambda_{k_2}, \dots, \lambda_{k_m})$ (7)

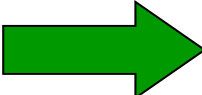
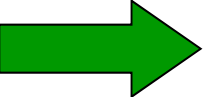
$k_i \in \{1, 2, \dots, d\}$

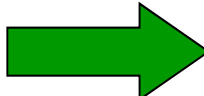
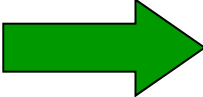
■ $\mathcal{R}(F) = \mathcal{R}(B^\top T) = \mathcal{R}(B^\top)$ (8)

■ (6) & (8) $\longrightarrow \mathcal{R}(B^\top) = \text{span}(\{\psi_{k_i}\}_{i=1}^m)$ (9)

Proof (cont.)

- (2)  $\text{rank}(\mathbf{B}) = m$

 all $\{k_i\}_{i=1}^m$ are **distinct**
- We should choose the best $\{k_i\}_{i=1}^m$ that maximizes $\text{tr}(\mathbf{B}\overline{\mathbf{C}}\mathbf{B}^\top)$.
- (4) & (7) 
$$\begin{aligned}\text{tr}(\mathbf{B}\overline{\mathbf{C}}\mathbf{B}^\top) &= \text{tr}(\mathbf{B}\mathbf{B}^\top \mathbf{T}\mathbf{T}^\top) \\ &= \text{tr}(\mathbf{T}\mathbf{T}^\top) \\ &= \text{tr}(\mathbf{T}^\top \mathbf{T}) \\ &= \sum_{i=1}^m \lambda_{k_i}\end{aligned}$$
- $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$

 $k_i = i$ gives a solution.
- (9)  $\mathbf{B} = (\psi_1 | \psi_2 | \dots | \psi_m)^\top$ (Q.E.D.)

Correlation

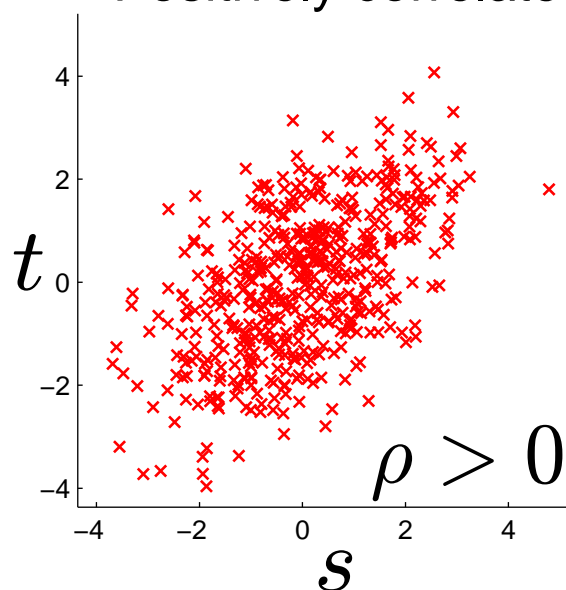
■ Correlation coefficient for $\{s_i, t_i\}_{i=1}^n$:

$$\rho = \frac{\sum_{i=1}^n (s_i - \bar{s})(t_i - \bar{t})}{\sqrt{(\sum_{i=1}^n (s_i - \bar{s})^2) (\sum_{i=1}^n (t_i - \bar{t})^2)}}$$

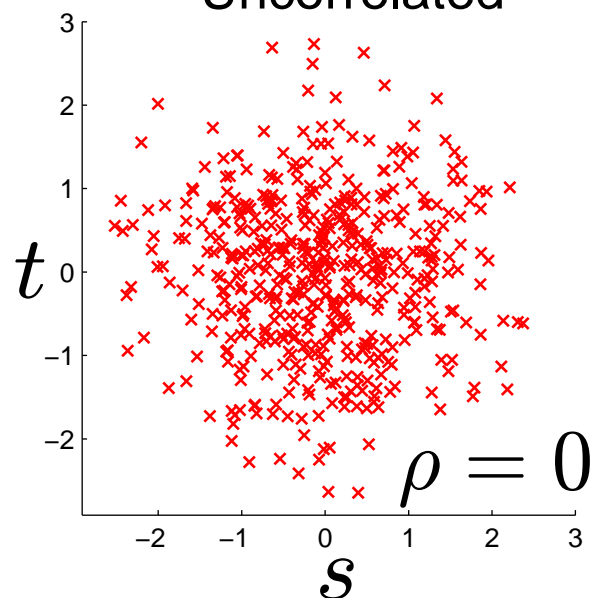
$$\bar{s} = \sum_{i=1}^n s_i$$

$$\bar{t} = \sum_{i=1}^n t_i$$

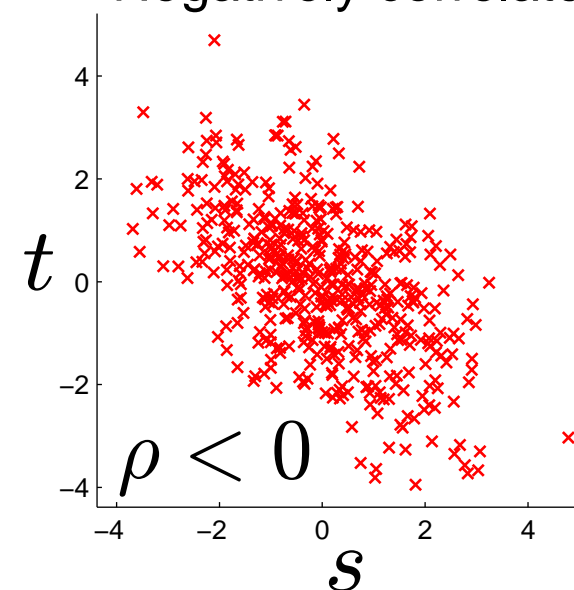
Positively correlated



Uncorrelated



Negatively correlated



PCA Uncorrelates Data

17

$$\mathbf{B}_{PCA} = (\boldsymbol{\psi}_1 | \boldsymbol{\psi}_2 | \cdots | \boldsymbol{\psi}_m)^\top$$

- Covariance matrix of the PCA-embedded samples is diagonal.

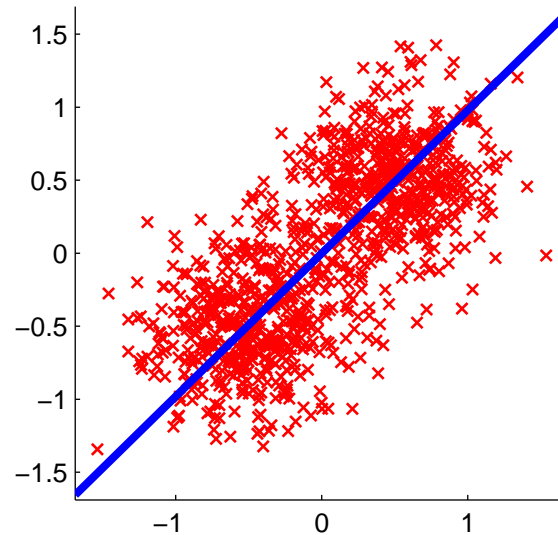
$$\frac{1}{n} \sum_{i=1}^n \bar{\mathbf{z}}_i \bar{\mathbf{z}}_i^\top = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$$

(Homework)

 Each element in $\bar{\mathbf{z}}$ is uncorrelated!

Examples

18

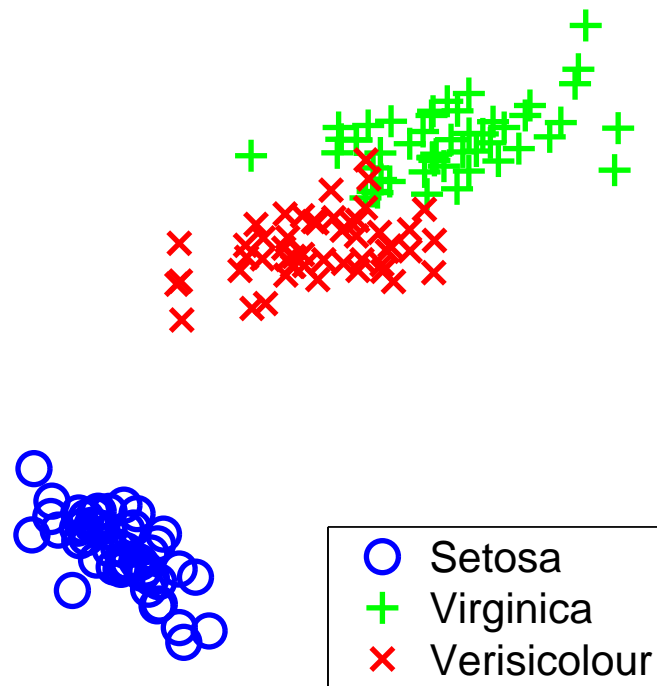


- Data is well described.
- PCA is intuitive, easy to implement, analytic solution available, and fast.

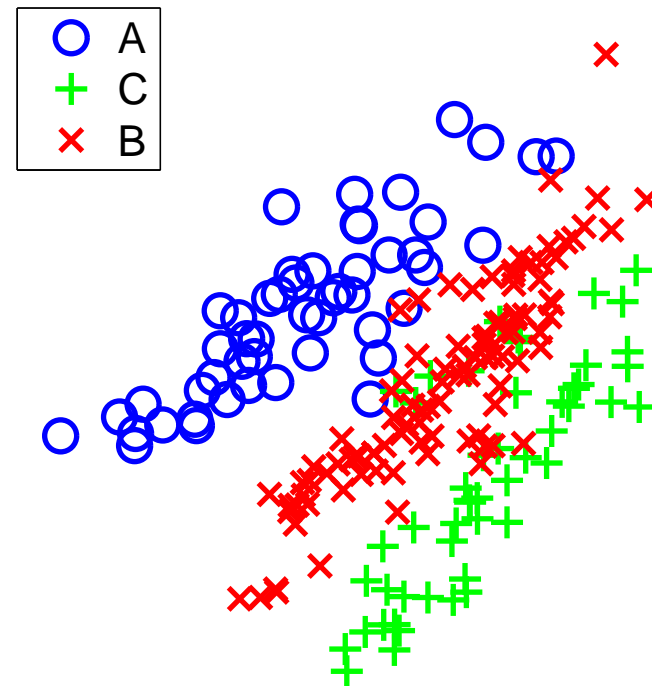
Examples (cont.)

19

Iris data (4d->2d)



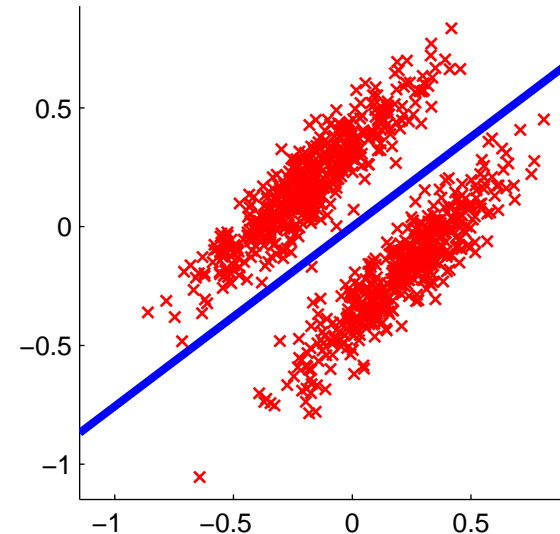
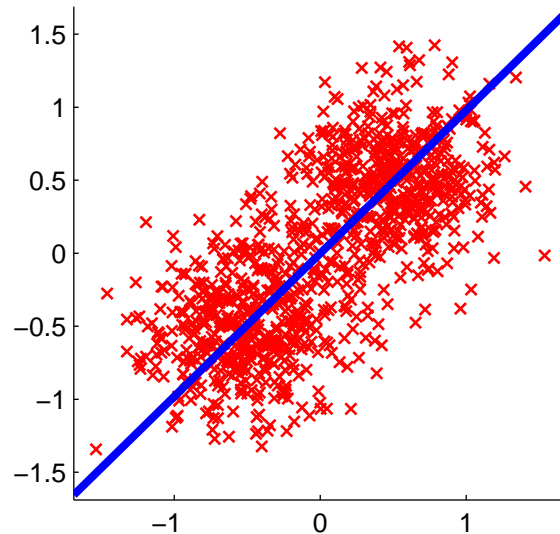
Letter data (16d->2d)



■ Embedded samples seem informative.

Examples (cont.)

20



- However, PCA does not necessarily preserve interesting information such as clusters.

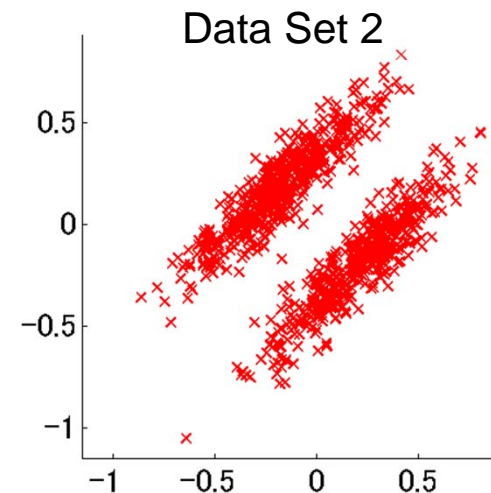
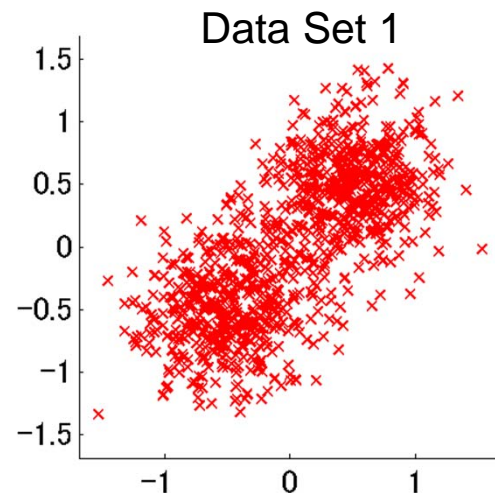
Homework

21

1. Implement PCA and reproduce the 2-dimensional examples shown in the class.

- Data sets 1 and 2 are available from

<http://sugiyama-www.cs.titech.ac.jp/~sugi/data/DataAnalysis>



- Test PCA on your own (artificial or real) data and analyze the characteristics of PCA.

Homework (cont.)

22

2. Let

- $B : m \times d, (1 \leq m \leq d)$
- $C, D : d \times d$, positive definite, symmetric
- $\{\lambda_i, \psi_i\}_{i=1}^m$: Sorted **generalized** eigenvalues and normalized eigenvectors of $C\psi = \lambda D\psi$

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$$

$$\langle D\psi_i, \psi_j \rangle = \delta_{i,j}$$

Prove that a solution of

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

$$\text{subject to } BDB^\top = I_m$$

is given by

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

Homework (cont.)

23

3. Prove that PCA uncorrelates the samples; more specifically, prove that the covariance matrix of the PCA-embedded samples is the following diagonal matrix:

$$\sum_{i=1}^n \bar{\mathbf{z}}_i \bar{\mathbf{z}}_i^\top = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$$

$$\bar{\mathbf{z}}_i = \mathbf{B}_{PCA} \bar{\mathbf{x}}_i$$

$$\mathbf{B}_{PCA} = (\boldsymbol{\psi}_1 | \boldsymbol{\psi}_2 | \dots | \boldsymbol{\psi}_m)^\top$$

Suggestion

24

- 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.

http://books.nips.cc/papers/files/nips16/NIPS2003_AA20.pdf