Advanced Data Analysis: Principal Component Analysis

Masashi Sugiyama (Computer Science)

W8E-505, <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi

Curse of Dimensionality

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

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

Length of the diagonal line of the unit hypercube gets large as the dimensionality increases.



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

- 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 highdimensional data is redundant in some sense.

Notation: Linear Embedding⁷

Data samples:

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

Embedding matrix:

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

Embedded data samples:

$$\{oldsymbol{z}_i\}_{i=1}^n, \hspace{0.2cm} oldsymbol{z}_i = oldsymbol{B}oldsymbol{x}_i \in \mathbb{R}^m$$



Principal Component Analysis (PCA)

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

We center the data samples by

$$\overline{\boldsymbol{x}}_i = \boldsymbol{x}_i - rac{1}{n} \sum_{j=1}^n \boldsymbol{x}_j$$

$$\frac{1}{n}\sum_{i=1}^{n}\overline{\boldsymbol{x}}_{i}=0$$

In matrix,

 $\overline{X} = XH$

$$egin{aligned} \overline{oldsymbol{X}} &= (\overline{oldsymbol{x}}_1 | \overline{oldsymbol{x}}_2 | \cdots | \overline{oldsymbol{x}}_n) \ oldsymbol{X} &= (oldsymbol{x}_1 | oldsymbol{x}_2 | \cdots | oldsymbol{x}_n) \ oldsymbol{H} &= oldsymbol{I}_n - rac{1}{n} oldsymbol{1}_{n imes n} \end{aligned}$$

 I_n : *n*-dimensional identity matrix

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

Orthogonal Projection

■ $\{b_i (\in \mathbb{R}^d)\}_{i=1}^m$: Orthonormal basis in *m*-dimensional embedding subspace

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

In matrix,
$$BB^{\top} = I_m$$

 $B = (b_1 | b_2 | \cdots | b_m)^{\top}$

Orthogonal projection of \overline{x}_i is expressed by

$$\sum_{j=1}^m \langle oldsymbol{b}_j, \overline{oldsymbol{x}}_i
angle oldsymbol{b}_j ~~ \left(=oldsymbol{B}^ op oldsymbol{B} \overline{oldsymbol{x}}_i
ight)$$

PCA Criterion

Minimize the sum of squared distances.

$$\sum_{i=1}^{n} \| \boldsymbol{B}^{\top} \boldsymbol{B} \overline{\boldsymbol{x}}_{i} - \overline{\boldsymbol{x}}_{i} \|^{2} \left(= -\operatorname{tr}(\boldsymbol{B} \overline{\boldsymbol{C}} \boldsymbol{B}^{\top}) + \operatorname{tr}(\overline{\boldsymbol{C}}) \right)$$

$$\overline{oldsymbol{C}} = \sum_{i=1}^n \overline{oldsymbol{x}}_i \overline{oldsymbol{x}}_i^{ op} = \overline{oldsymbol{X}} \ \overline{oldsymbol{X}}^{ op}$$

PCA criterion:

 $\boldsymbol{B}_{PCA} = \operatorname*{argmax}_{\boldsymbol{B} \in \mathbb{R}^{m \times d}} \operatorname{tr}(\boldsymbol{B} \overline{\boldsymbol{C}} \boldsymbol{B}^{\top})$ subject to $\boldsymbol{B} \boldsymbol{B}^{\top} = \boldsymbol{I}_{m}$



11

PCA: Summary

A PCA solution:

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

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

 $\lambda_1 \ge \lambda_2 \ge \cdots \ge \lambda_d$

$$\langle oldsymbol{\psi}_i,oldsymbol{\psi}_j
angle=\delta_{i,j}$$

PCA embedding of a sample x:

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

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

Proof13Lagrangian:
$$L(B, \Delta) = \operatorname{tr}(B\overline{C}B^{\top}) - \operatorname{tr}((BB^{\top} - I_m)\Delta)$$

 Δ :Lagrange multipliers (symmetric)Stationary point (necessary condition): $\frac{\partial L}{\partial B} = 2B\overline{C} - 2\Delta B = 0$
 $\frac{\partial L}{\partial \Delta} = BB^{\top} - I_m = 0$ $\frac{\partial L}{\partial \Delta} = BB^{\top} - I_m = 0$ $BB^{\top} = I_m (2)$ Eigendecomposition:
 $\Delta = T\Gamma T^{\top} (3)$ $T^{-1} = T^{\top}$

Proof (cont.)
14
(1) & (3)

$$\overline{C}B^{\top} = B^{\top}T\Gamma T^{\top}(4)$$

 $\overline{C}B^{\top}T = B^{\top}T\Gamma$
 $\overline{C}F = F\Gamma$ (5)
 $F = B^{\top}T$
(5) is an eigensystem
 $\mathcal{R}(F) = \operatorname{span}(\{\psi_{k_i}\}_{i=1}^m)$ (6)
 $\Gamma = \operatorname{diag}(\lambda_{k_1}, \lambda_{k_2}, \dots, \lambda_{k_m})$ (7)
 $k_i \in \{1, 2, \dots, d\}$
 $\mathcal{R}(F) = \mathcal{R}(B^{\top}T) = \mathcal{R}(B^{\top})$ (8)
(6) & (8)
 $\mathcal{R}(B^{\top}) = \operatorname{span}(\{\psi_{k_i}\}_{i=1}^m)$
Solution is expressed as
 $B = (\psi_{k_1} | \psi_{k_2} | \cdots | \psi_{k_m})^{\top}$



Correlation

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

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





4

2

 t_{\circ}

-2

PCA Uncorrelates Data

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

Covariance matrix of the PCAembedded samples is diagonal.

$$\frac{1}{n}\sum_{i=1}^{n} \overline{\boldsymbol{z}}_{i} \overline{\boldsymbol{z}}_{i}^{\top} = \operatorname{diag}\left(\lambda_{1}, \lambda_{2}, \dots, \lambda_{m}\right)$$

(Homework)

17

Each element in \overline{z} is uncorrelated!





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.)



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

Homework

- 1. Implement PCA and reproduce the 2dimensional examples shown in the class.
 - Data sets are available from

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



• You may create similar (or more interesting) data sets by yourself.

Homework (cont.)

2. Let

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

$$\lambda_1 \ge \lambda_2 \ge \dots \ge \lambda_d$$

$$\langle {oldsymbol D} {oldsymbol \psi}_i, {oldsymbol \psi}_j
angle = \delta_{i,j}$$

Prove that a solution of

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

subject to $\boldsymbol{B}\boldsymbol{D}\boldsymbol{B}^{\top}=\boldsymbol{I}_{m}$

is given by

$$\boldsymbol{B}_{min} = (\boldsymbol{\psi}_d | \boldsymbol{\psi}_{d-1} | \cdots | \boldsymbol{\psi}_{d-m+1})^{\mathsf{T}}$$

Homework (cont.)

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

$$\frac{1}{n}\sum_{i=1}^{n} \overline{\boldsymbol{z}}_{i} \overline{\boldsymbol{z}}_{i}^{\top} = \operatorname{diag}\left(\lambda_{1}, \lambda_{2}, \dots, \lambda_{m}\right)$$

$$\overline{\boldsymbol{z}}_i = \boldsymbol{B}_{PCA} \overline{\boldsymbol{x}}_i$$
$$\boldsymbol{B}_{PCA} = (\boldsymbol{\psi}_1 | \boldsymbol{\psi}_2 | \cdots | \boldsymbol{\psi}_m)^\top$$

Suggestion

 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