# Advanced Data Analysis: Kernel PCA
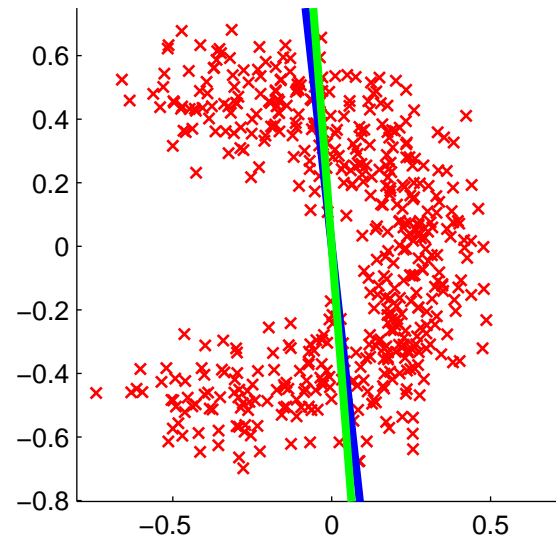
Masashi Sugiyama (Computer Science)

W8E-505, sugi@cs.titech.ac.jp

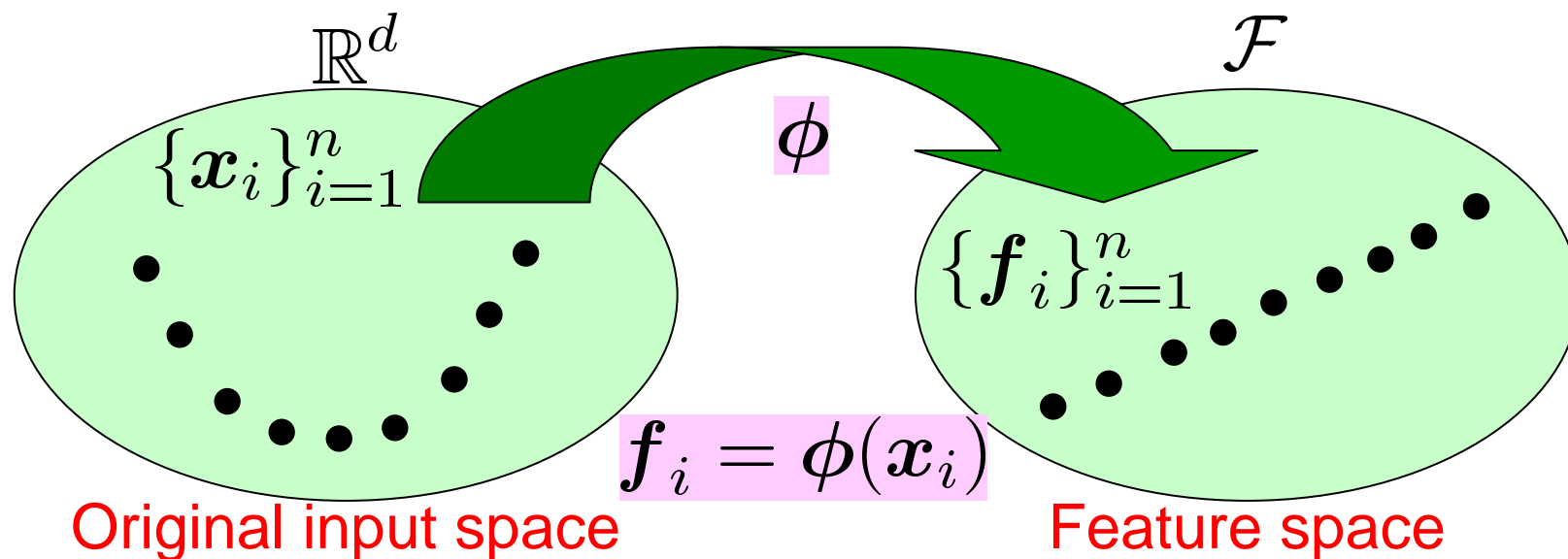http://sugiyama-www.cs.titech.ac.jp/~sugi

# Data with Curved Structures



■ If the data cloud is bent, any linear methods cannot find the curved structure.

➡ Limitation of linear method!

# Non-Linearizing Linear Methods

- A simple non-linear extension of linear methods while keeping computational advantages of linear methods:
  - Map the original data to a feature space by a non-linear transformation
  - Run linear algorithm in the feature space



$$\mathbb{R}^d \qquad \mathcal{F}$$

$$\{\boldsymbol{x}_i\}_{i=1}^n \qquad \boldsymbol{\phi} \qquad \{\boldsymbol{f}_i\}_{i=1}^n$$

$$\boldsymbol{f}_i = \boldsymbol{\phi}(\boldsymbol{x}_i)$$
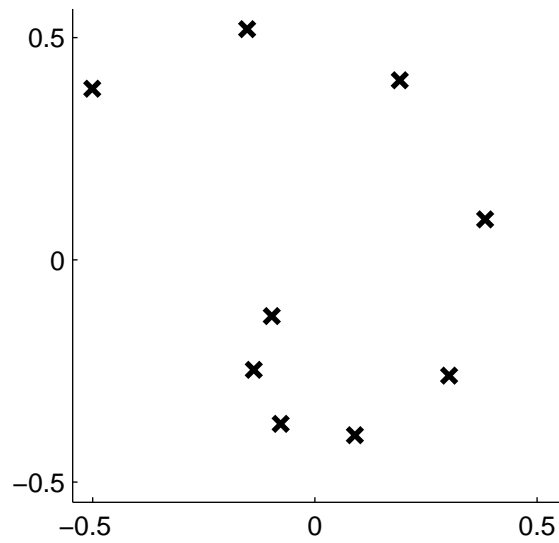
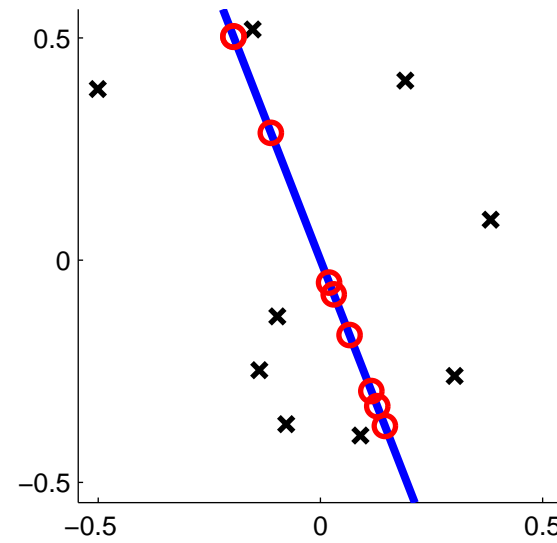Original input space  Feature space

# Example

■ $d = 2$

Centered data
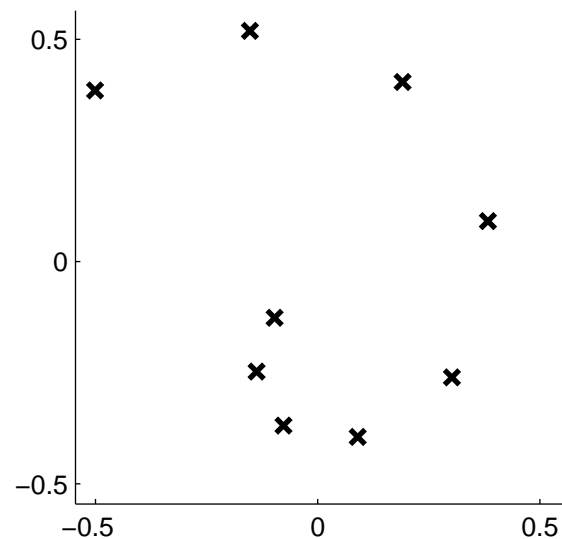in input space

Linear PCA

# Example (cont.)

**Polar coordinate:**

$$\boldsymbol{x} = \begin{pmatrix} a \\ b \end{pmatrix} \longrightarrow \boldsymbol{f} = \begin{pmatrix} r\cos\theta \\ r\sin\theta \end{pmatrix}$$

Centered data
in input space

Centered data
in feature space

# Example (cont.)

■ Run PCA in feature space.

Centered data
in feature space

PCA projection
in feature space

# Example (cont.)

■ Pull the results back to input space.

Non-linear PCA                    Linear PCA



■ Non-linear PCA describes the original data much better than linear PCA.

# Notation Revisited

- Input samples:
$$\{\boldsymbol{x}_i\}_{i=1}^{n} \qquad \boldsymbol{x}_i \in \mathbb{R}^d$$

- Feature mapping:
$$\boldsymbol{\phi} : \mathbb{R}^d \to \mathcal{F}$$

- Samples in feature space:
$$\boldsymbol{f}_i = \boldsymbol{\phi}(\boldsymbol{x}_i)$$

# Centering in Feature Space

- PCA requires centered samples, thus we need to center samples by

$$\overline{\boldsymbol{f}}_i = \boldsymbol{f}_i - \frac{1}{n} \sum_{j=1}^{n} \boldsymbol{f}_j$$

- In matrix form,

$$\overline{\boldsymbol{F}} = \boldsymbol{F}\boldsymbol{H}$$

$$\boldsymbol{F} = (\boldsymbol{f}_1 | \boldsymbol{f}_2 | \cdots | \boldsymbol{f}_n)$$
$$\overline{\boldsymbol{F}} = (\overline{\boldsymbol{f}}_1 | \overline{\boldsymbol{f}}_2 | \cdots | \overline{\boldsymbol{f}}_n)$$

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

$\boldsymbol{I}_n$: $n$-dimensional identity matrix

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

# PCA in Feature Space (Primal)

$$\overline{C}\psi = \lambda\psi \qquad\qquad \overline{C} = \overline{F}\,\overline{F}^\top$$

■ PCA solution:

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

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

$$\langle \psi_i, \psi_j \rangle = \delta_{i,j} \qquad\qquad \lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_\mu$$

■ PCA embedding of a sample $f$ :

$$\overline{g} = B_{PCA}\left(f - \frac{1}{n}F\mathbf{1}_n\right)$$

$\mu = \dim(\mathcal{F})$ $\qquad \mathbf{1}_n$: $n$-dimensional vector with all ones

# PCA in High-Dimensional Feature Space

$$\mu = \dim(\mathcal{F})$$

- If $\mu$ is high,
  - Description ability of non-linear PCA will increase.
  - However, computational cost increases since the dimension of $\overline{C}$ is $\mu$.

- It would be possible to reduce computational cost since

$$\mathrm{rank}\,(\overline{C}) = \min(\mu, n) \leq \mu$$

$$\overline{C} = \overline{F}\,\overline{F}^{\top} \qquad \overline{F} = (\overline{f}_1 | \overline{f}_2 | \cdots | \overline{f}_n)$$

# Dual Formulation

(A) $\overline{C}\psi = \lambda\psi$ $\qquad$ $\overline{C} = \overline{F}\,\overline{F}^{\top}$

(B) $\overline{K}\alpha = \lambda\alpha$ $\qquad$ $\overline{K} = \overline{F}^{\top}\overline{F}$

■ Solution of (A) can be obtained from (B).

- Proof: If $\alpha$ is a solution of (B), it holds that

$$\overline{C}\,\overline{F}\alpha = \overline{F}\,\overline{F}^{\top}\,\overline{F}\alpha = \overline{F}\,\overline{K}\alpha = \lambda\overline{F}\alpha$$

This implies that $\psi = \overline{F}\alpha$ is a solution of (A).

■ Note: solution of (B) can also be obtained from (A).

■ Given $\overline{K}$, solving (B) is faster than (A) when $\mu > n$ since

$$\mathrm{rank}\,(\overline{C}) = n < \mu$$

# Primal and Dual Formulations

$$\psi = \overline{F}\alpha$$

**Primal**

$$\overline{C}\psi = \lambda\psi$$

$$\boxed{\overline{C}} = \boxed{\overline{F}}\;\boxed{\overline{F}^{\top}}$$

**Equivalent**

**Dual**

$$\overline{K}\alpha = \lambda\alpha$$

$$\boxed{\overline{K}} = \boxed{\overline{F}^{\top}}\;\boxed{\overline{F}}$$

$$\overline{C} = \overline{F}\,\overline{F}^{\top}$$

$$\overline{K} = \overline{F}^{\top}\overline{F}$$

$$\overline{F} = (\overline{f}_1|\overline{f}_2|\cdots|\overline{f}_n)$$

# Renormalization of Eigenvectors

$$\overline{K}\alpha = \lambda\alpha$$

■ Standard eigensolvers output an orthonormal eigenvectors.

$$\langle \alpha_i, \alpha_j \rangle = \delta_{i,j}$$

■ However, PCA requires the primal eigenvectors $\{\psi_i\}_{i=1}^m$ to be orthonormal.

■ Since $\langle \psi_i, \psi_j \rangle = \langle \overline{K}\alpha_i, \alpha_j \rangle = \lambda_i \delta_{i,j}$ , we need to renormalize $\{\psi_i\}_{i=1}^m$ by

$$\psi_i \longleftarrow \frac{\psi_i}{\|\psi_i\|} = \frac{1}{\sqrt{\lambda_i}}\overline{F}\alpha_i$$

$$\psi_i = \overline{F}\alpha_i$$

$$\overline{K}\alpha_i = \lambda_i \alpha_i$$

# PCA in Feature Space (Dual)

■ PCA embedding of a sample $f$ :

$$\overline{g} = \Lambda^{-\frac{1}{2}} A^\top H(k - \frac{1}{n}K\mathbf{1}_n)$$  (Homework)

- $\{\lambda_i, \boldsymbol{\alpha}_i\}_{i=1}^m$ :Sorted eigenvalues and normalized eigenvectors of $\overline{K}\boldsymbol{\alpha} = \lambda\boldsymbol{\alpha}$

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n \qquad \langle \boldsymbol{\alpha}_i, \boldsymbol{\alpha}_j \rangle = \delta_{i,j}$$

$\Lambda = \operatorname{diag}(\lambda_1, \lambda_2, \ldots, \lambda_m)$

$A = (\boldsymbol{\alpha}_1 | \boldsymbol{\alpha}_2 | \cdots | \boldsymbol{\alpha}_m)$

$\overline{K} = HKH \quad K = F^\top F$

$H = I_n - \frac{1}{n}\mathbf{1}_{n\times n} \quad k = F^\top f$

$I_n$: $n$-dimensional identity matrix

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

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

# PCA in Feature Space (Dual)

$$\mu = \dim(\mathcal{F})$$

- In the dual formulation, the computational complexity depends not on $\mu$ but only on $n$, if $K$ and $k$ are given.

- However, the computation of $K$ and $k$ still depends on $\mu$.

$$K = F^\top F \qquad k = f^\top F$$

- Note: $K$ and $k$ depend on $\mu$ only through the inner product between samples.

$$K_{i,j} = \langle f_i, f_j \rangle \qquad k_i = \langle f, f_i \rangle$$

# Kernel Trick

- For some transformation $\phi(\boldsymbol{x})\ (=\boldsymbol{f})$, there exists a bivariate function $K(\boldsymbol{x}, \boldsymbol{x}')$ such that

$$\boldsymbol{K}_{i,j} = \langle \boldsymbol{f}_i, \boldsymbol{f}_j \rangle = K(\boldsymbol{x}_i, \boldsymbol{x}_j)$$

- Such implicit mapping $\phi(\boldsymbol{x})$ exists if
  - $\boldsymbol{K}$ is symmetric: $\boldsymbol{K}^\top = \boldsymbol{K}$
  - $\boldsymbol{K}$ is positive semi-definite: $\forall \boldsymbol{y}, \ \langle \boldsymbol{K}\boldsymbol{y}, \boldsymbol{y} \rangle \geq 0$

- Such $K(\boldsymbol{x}, \boldsymbol{x}')$ is called the <span style="color:red">reproducing kernel</span>.

- Rather than directly defining $\phi(\boldsymbol{x})$, we implicitly specify $\phi(\boldsymbol{x})$ by a reproducing kernel.

# Examples of Kernels
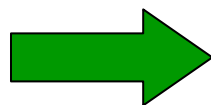
■ Polynomial kernel:

$$\mu = \dim(\mathcal{F})$$

$$K(\boldsymbol{x}, \boldsymbol{x}') = \langle \boldsymbol{x}, \boldsymbol{x}' \rangle^c \qquad c \in \mathbb{N}$$

● When $d = 2$ and $c = 2$,

$$\langle \boldsymbol{x}, \boldsymbol{x}' \rangle^2 = (ss' + tt')^2$$
$$= sss's' + 2ss'tt' + ttt't'$$

$$\boldsymbol{x} = \begin{pmatrix} s \\ t \end{pmatrix}$$

$$\Longrightarrow \quad \boldsymbol{f} = \boldsymbol{\phi}(\boldsymbol{x}) = \begin{pmatrix} s^2 \\ \sqrt{2}st \\ t^2 \end{pmatrix}$$

$$\mu = 3$$

● In general,

$$\mu = {}_{c+d-1}C_c$$

■ Gaussian kernel:

$$K(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\|\boldsymbol{x} - \boldsymbol{x}'\|^2 / c^2\right)$$

$$c > 0$$

Note: $\mu = \infty$ !

$$\mu = \dim(\mathcal{F})$$

# Kernel PCA: Summary

■ Kernel PCA embedding of a sample $f$ is

$$\overline{g} = \Lambda^{-\frac{1}{2}} A^\top H (k - \frac{1}{n} K \mathbf{1}_n)$$

- $\{\lambda_i, \boldsymbol{\alpha}_i\}_{i=1}^m$ :Sorted eigenvalues and normalized eigenvectors of $HKH\boldsymbol{\alpha} = \lambda\boldsymbol{\alpha}$

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n \qquad \langle \boldsymbol{\alpha}_i, \boldsymbol{\alpha}_j \rangle = \delta_{i,j}$$

$\Lambda = \mathrm{diag}(\lambda_1, \lambda_2, \ldots, \lambda_m)$

$A = (\boldsymbol{\alpha}_1 | \boldsymbol{\alpha}_2 | \cdots | \boldsymbol{\alpha}_m)$

$I_n$: $n$-dimensional identity matrix

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

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

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

$k = (K(\boldsymbol{x}, \boldsymbol{x}_1), K(\boldsymbol{x}, \boldsymbol{x}_2), \ldots, K(\boldsymbol{x}, \boldsymbol{x}_n))^\top$

$K_{i,j} = K(\boldsymbol{x}_i, \boldsymbol{x}_j)$

# Examples

■ Wine data (UCI): 13-dim, 178 samples

$$K(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\|\boldsymbol{x} - \boldsymbol{x}'\|^2 / c^2\right)$$



$c = 3$

Linear PCA      Gaussian KPCA

# Examples (cont.)

$$K(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\|\boldsymbol{x} - \boldsymbol{x}'\|^2 / c^2\right)$$

$c = 1$

$c = 6$

$c = 3$

- Choice of kernels (type and parameter) depends on the result.
- Appropriately choosing kernels is not straightforward in practice.

# Homework

1. Implement kernel PCA with Gaussian kernels and reproduce the embedding result of the Wine data set.

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

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

2. Prove that kernel PCA embedding of a sample $f$ is given by

$$\overline{g} = \Lambda^{-\frac{1}{2}} A^\top H \left( k - \frac{1}{n} K \mathbf{1}_n \right)$$

# Suggestion

■ Read the following article for the next class:

- M. Belkin & P. Niyogi: Laplacian eigenmaps for dimensionality reduction and data representation, Neural Computation, 15(6), 1373-1396, 2003.

  http://neco.mitpress.org/cgi/reprint/15/6/1373.pdf