

# Advanced Data Analysis: Kernel PCA

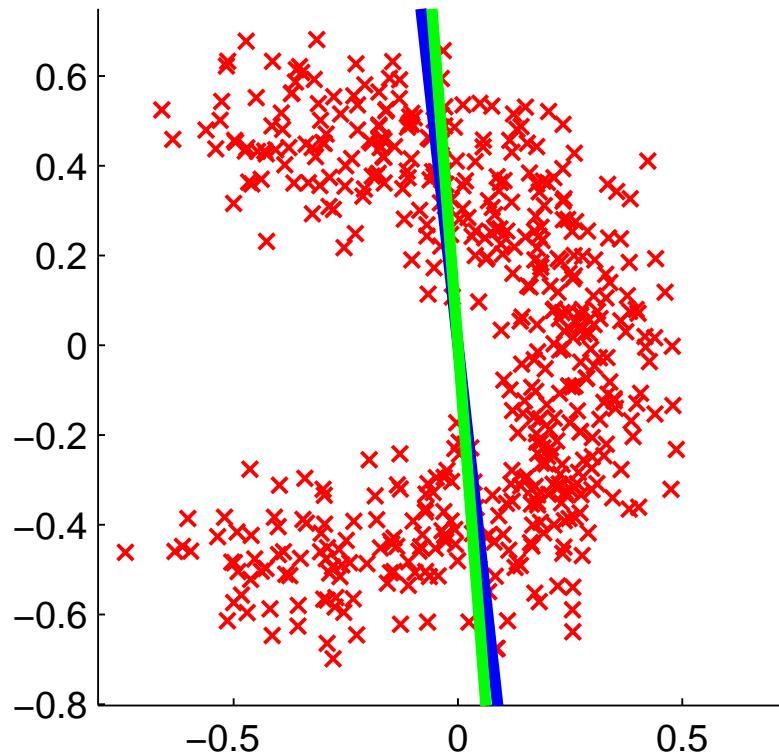
Masashi Sugiyama (Computer Science)

W8E-406, [sugi@cs.titech.ac.jp](mailto:sugi@cs.titech.ac.jp)

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

# Data with Curved Structure

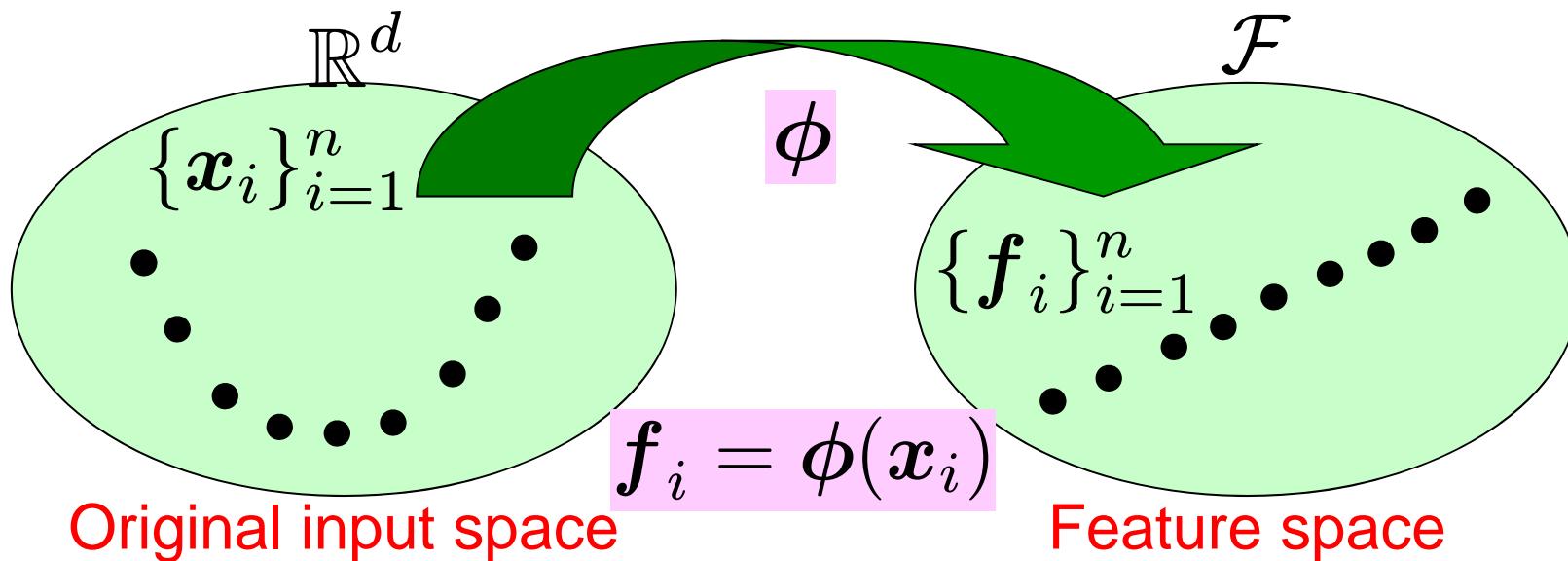
85



- Limitation of linear methods: Any linear methods cannot find curved structure.

# Non-Linearizing Linear Methods<sup>86</sup>

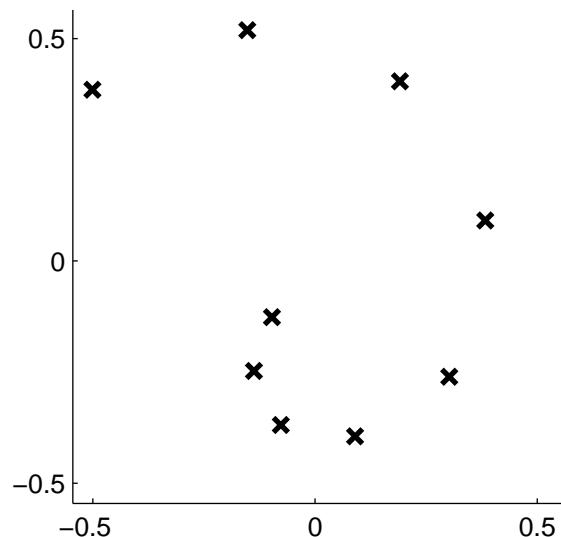
- A simple non-linear extension of linear methods while keeping computational advantages of linear methods:
  - Map original data to a **feature space** by a **non-linear transformation**
  - Run a **linear algorithm** in the feature space



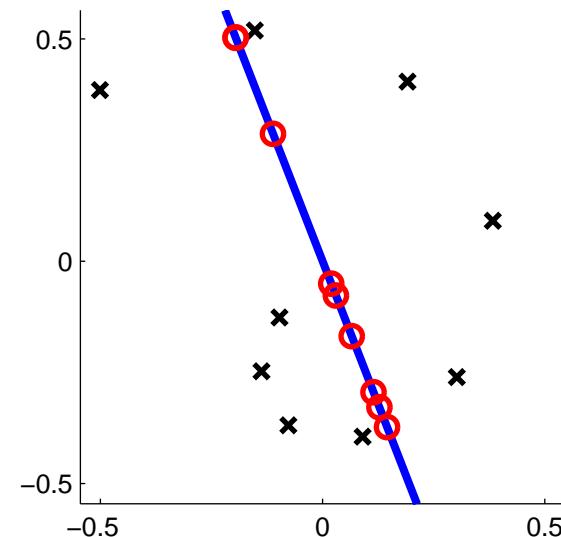
# Example

■  $d = 2$

Data in input space



Linear PCA

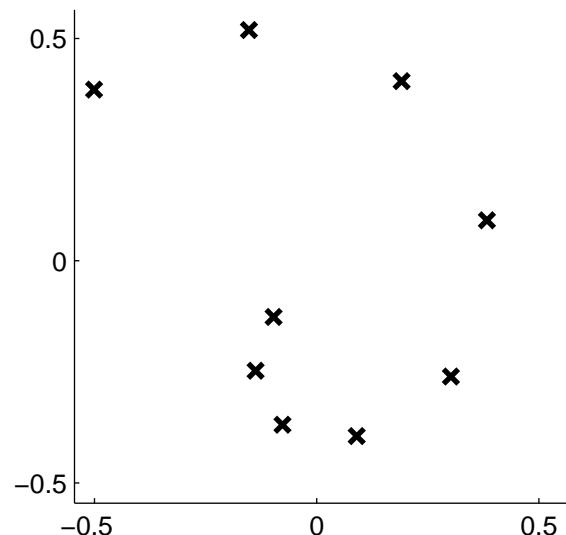


# Example (cont.)

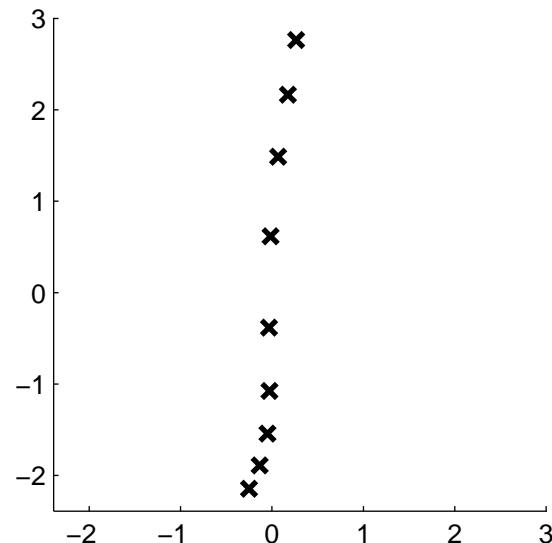
## ■ Polar coordinate:

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

Data in input space



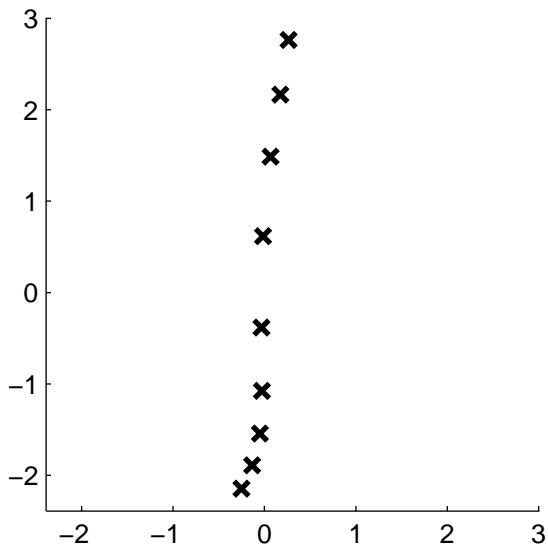
Data in feature space



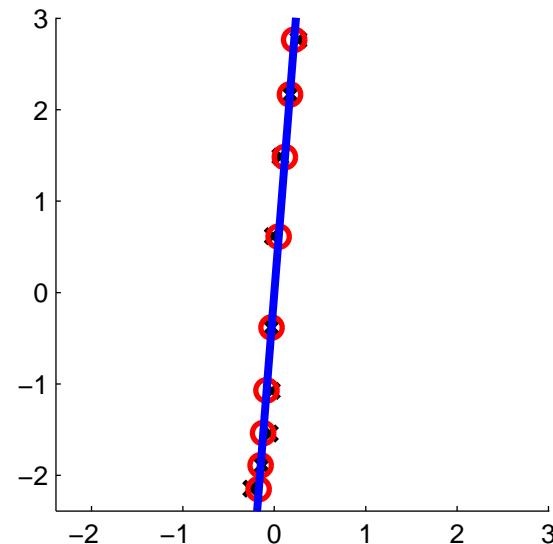
# Example (cont.)

■ Run PCA in feature space.

Data  
in feature space



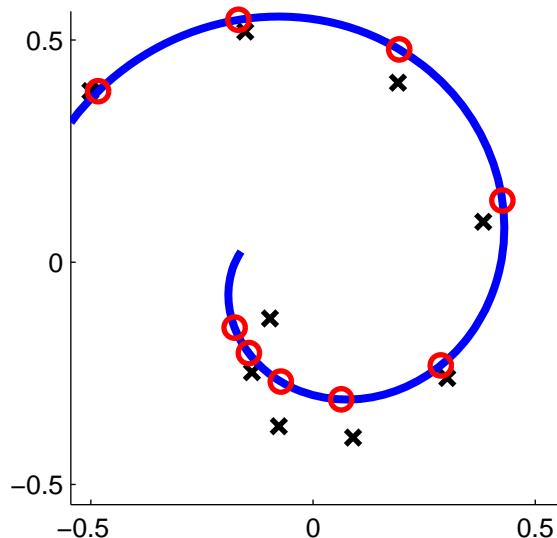
PCA projection  
in feature space



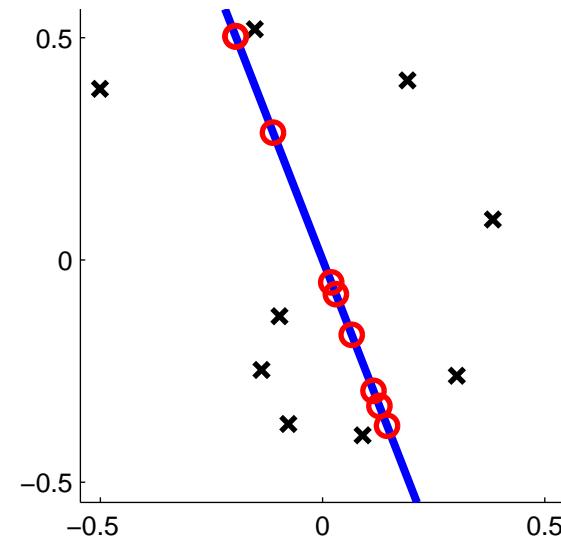
# Example (cont.)

- Pull the projected points 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:

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

■ Feature mapping:

$$\phi : \mathbb{R}^d \rightarrow \mathcal{F}$$

■ Samples in feature space:

$$\mathbf{f}_i = \phi(\mathbf{x}_i)$$

# Centering in Feature Space

- PCA requires centered samples, so we need to explicitly center samples by

$$\bar{f}_i = f_i - \frac{1}{n} \sum_{j=1}^n f_j$$

- In matrix form,

$$\bar{F} = F H$$

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

$$F = (f_1 | f_2 | \cdots | f_n)$$

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

$I_n$ :  $n$ -dimensional identity matrix

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

# PCA in Feature Space (Primal)<sup>93</sup>

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

$$\bar{C} = \bar{F} \bar{F}^\top$$

## ■ PCA solution:

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

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

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

## ■ PCA embedding of sample $f$ :

$$\bar{g} = B_{\text{PCA}}(f - \frac{1}{n}F\mathbf{1}_n)$$

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

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

# PCA in High-Dimensional Feature Space

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

- If  $\mu$  is larger, non-linear PCA will be more expressive.
- However, computational costs increase since the matrix  $\bar{C}$  we eigendecompose is  $\mu \times \mu$ .
- If  $\mu > n$ , it is possible to reduce the computation cost because

$$\text{rank}(\bar{C}) \leq n < \mu$$

$$\bar{C} = \bar{F} \bar{F}^\top$$

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

# Dual Formulation

$$(A) \quad \bar{C}\psi = \lambda\psi$$

$$(B) \quad \bar{K}\alpha = \lambda\alpha$$

$$\bar{C} = \bar{F} \bar{F}^\top$$

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

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

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

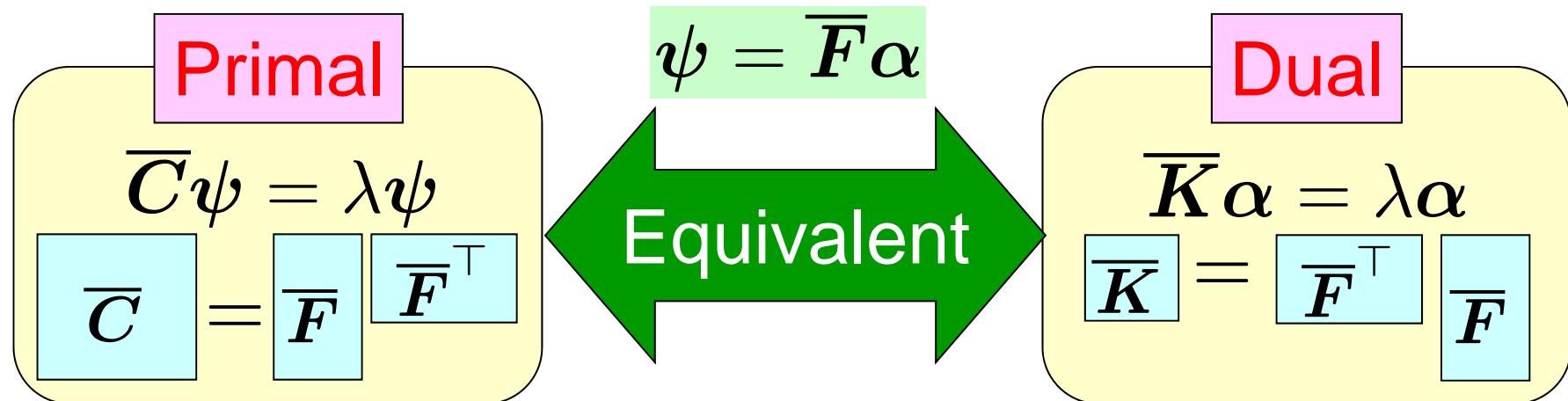
$$\bar{C} \bar{F}\alpha = \bar{F} \bar{F}^\top \bar{F}\alpha = \bar{F} \bar{K}\alpha = \lambda \bar{F}\alpha$$

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

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

■ Given  $\bar{K}$ , solving (B) is faster than (A) when  $\mu > n$  since  $\text{rank}(\bar{C}) \leq n < \mu$

# Primal and Dual Formulations 96



$$\bar{C} = \bar{F} \bar{F}^\top$$

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

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

# Renormalization of Eigenvectors<sup>97</sup>

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

- Standard eigensolvers output 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  $\langle \psi_i, \psi_j \rangle = \langle \bar{K}\alpha_i, \alpha_j \rangle = \lambda_i \delta_{i,j}$ , we need to explicitly renormalize  $\{\psi_i\}_{i=1}^m$  as

$$\psi_i \leftarrow \frac{\psi_i}{\|\psi_i\|} = \frac{1}{\sqrt{\lambda_i}} \bar{F} \alpha_i$$

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

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

# PCA in Feature Space (Dual) 98

■ PCA embedding of sample  $f$  :

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

- $\{\lambda_i, \alpha_i\}_{i=1}^n$ : Sorted eigenvalues and normalized eigenvectors of  $H K H \alpha = \lambda \alpha$

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n \quad \langle \alpha_i, \alpha_j \rangle = \delta_{i,j}$$

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$$

$$A = (\alpha_1 | \alpha_2 | \cdots | \alpha_m)$$

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

$$K = F^\top F$$

$$k = F^\top f$$

$I_n$ :  $n$ -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) 99

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

- In the dual formulation, the computation cost depends not on  $\mu$  but only on  $n$ , if  $K$  and  $k$  are given.
- However, the computation costs of  $K$  and  $k$  still depend on  $\mu$ .

$$K = F^\top F$$

$$k = f^\top F$$

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

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

$$k_i = \langle f, f_i \rangle$$

# Kernel Trick

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

$$K_{i,j} = \langle \mathbf{f}_i, \mathbf{f}_j \rangle = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = K(\mathbf{x}_i, \mathbf{x}_j)$$

if  $K$  is symmetric and positive semi-definite:

$$\mathbf{K}^\top = \mathbf{K} \quad \forall \mathbf{y}, \quad \langle \mathbf{K}\mathbf{y}, \mathbf{y} \rangle \geq 0$$

- Such  $K(\mathbf{x}, \mathbf{x}')$  is called the **reproducing kernel**.
- Rather than directly specifying  $\phi(\mathbf{x})$  , we implicitly specify  $\phi(\mathbf{x})$  by a reproducing kernel.

# Examples of Kernels

101

## ■ Polynomial kernel:

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

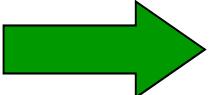
$$K(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle^c$$

$$c \in \mathbb{N}$$

- When  $d = 2$  and  $c = 2$ ,

$$\begin{aligned}\langle \mathbf{x}, \mathbf{x}' \rangle^2 &= (ss' + tt')^2 \\ &= sss's' + 2ss'tt' + ttt't'\end{aligned}$$

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


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

$$\mu = 3$$

- In general,  $\mu = \binom{c+d-1}{c}$

# Examples of Kernels (cont.) <sup>102</sup>

■ Gaussian kernel:

$$K(x, x') = \exp(-\|x - x'\|^2/c^2)$$

$$c > 0$$

Note:  $\mu = \infty$  !

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

# Kernel PCA: Summary

■ Kernel PCA embedding of sample  $f$  is

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

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

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

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

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$$

$$A = (\alpha_1 | \alpha_2 | \cdots | \alpha_m)$$

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

$$k = (K(x, x_1), K(x, x_2), \dots, K(x, x_n))^\top$$

$I_n$ :  $n$ -dimensional identity matrix

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

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

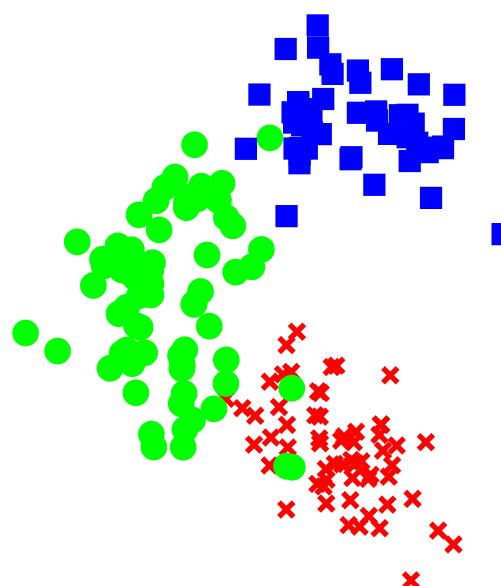
$$K_{i,j} = K(x_i, x_j)$$

# Examples

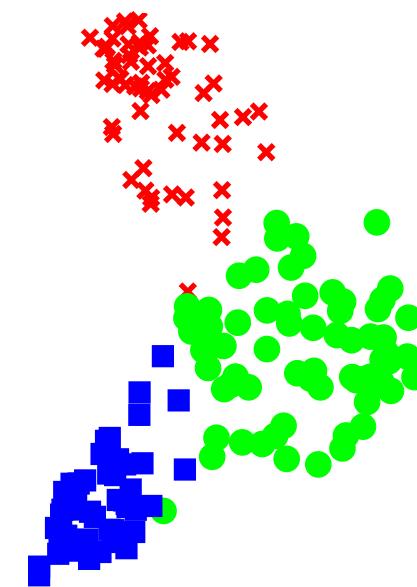
104

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

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



Linear PCA

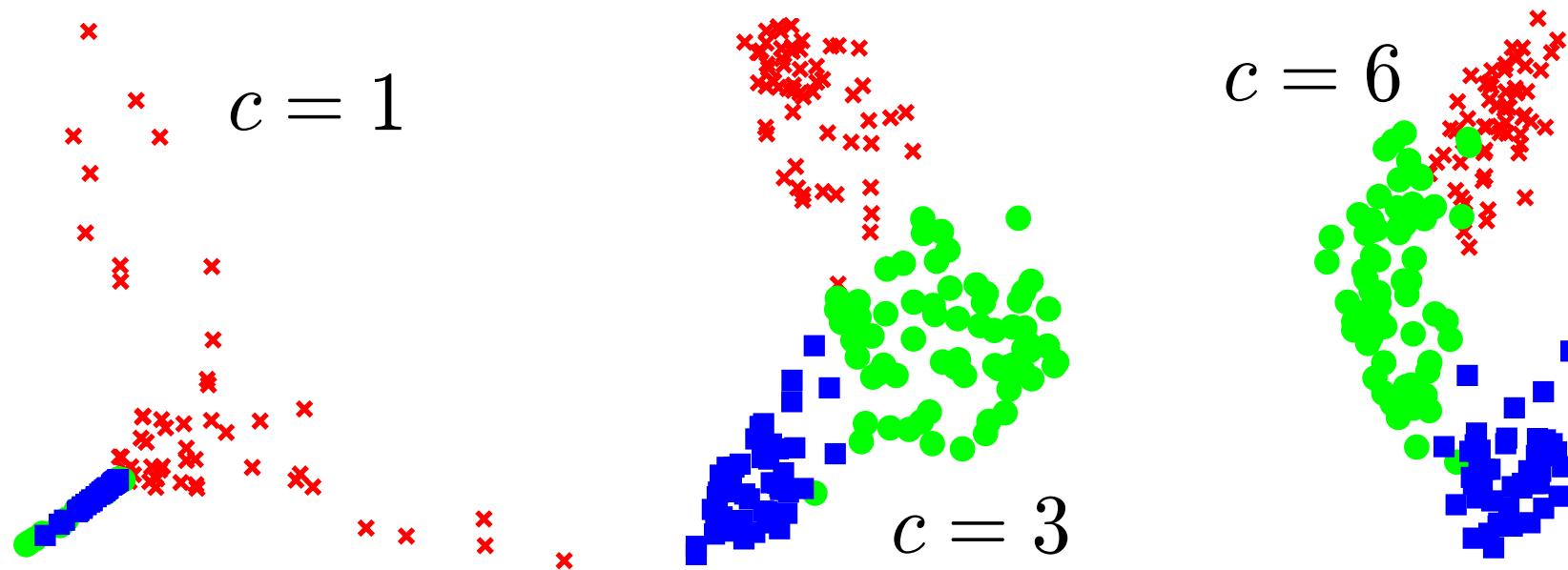


Gaussian KPCA

$$c = 3$$

# Examples (cont.)

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



- 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

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

# 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>