# Advanced Data Analysis: More on Kernels

Masashi Sugiyama (Computer Science)

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

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

# Kernel Trick
# with Reproducing Kernel

- 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: $\forall \boldsymbol{y}, \ \langle \boldsymbol{K}\boldsymbol{y}, \boldsymbol{y} \rangle \geq 0$

# Combination of Reproducing Kernels

For any reproducing kernels (RKs)

$$K^{(1)}(\boldsymbol{x}, \boldsymbol{x}'), K^{(2)}(\boldsymbol{x}, \boldsymbol{x}')$$

- Positive scaling of RK is still RK

$$K(\boldsymbol{x}, \boldsymbol{x}') = \alpha K^{(1)}(\boldsymbol{x}, \boldsymbol{x}') \quad \alpha > 0$$

- Sum of RKs is still RK:

$$K(\boldsymbol{x}, \boldsymbol{x}') = K^{(1)}(\boldsymbol{x}, \boldsymbol{x}') + K^{(2)}(\boldsymbol{x}, \boldsymbol{x}')$$

- Product of RKs is still RK:

$$K(\boldsymbol{x}, \boldsymbol{x}') = K^{(1)}(\boldsymbol{x}, \boldsymbol{x}')K^{(2)}(\boldsymbol{x}, \boldsymbol{x}')$$

# Proof

We prove that there exists a feature map $\phi(\boldsymbol{x})$ such that $\langle \phi(\boldsymbol{x}), \phi(\boldsymbol{x}') \rangle = K(\boldsymbol{x}, \boldsymbol{x}')$ .

- For $\phi(\boldsymbol{x}) = \sqrt{\alpha}\phi^{(1)}(\boldsymbol{x})$,

$$\langle \phi(\boldsymbol{x}), \phi(\boldsymbol{x}') \rangle = \alpha \langle \phi^{(1)}(\boldsymbol{x}), \phi^{(1)}(\boldsymbol{x}') \rangle = \alpha K^{(1)}(\boldsymbol{x}, \boldsymbol{x}')$$

- For $\phi(\boldsymbol{x}) = \begin{pmatrix} \phi^{(1)}(\boldsymbol{x}) \\ \phi^{(2)}(\boldsymbol{x}) \end{pmatrix}$, $\qquad K^{(i)}(\boldsymbol{x}, \boldsymbol{x}') = \langle \phi^{(i)}(\boldsymbol{x}), \phi^{(i)}(\boldsymbol{x}') \rangle$

$$\langle \phi(\boldsymbol{x}), \phi(\boldsymbol{x}') \rangle = \langle \phi^{(1)}(\boldsymbol{x}), \phi^{(1)}(\boldsymbol{x}') \rangle + \langle \phi^{(2)}(\boldsymbol{x}), \phi^{(2)}(\boldsymbol{x}') \rangle$$

$$= K^{(1)}(\boldsymbol{x}, \boldsymbol{x}') + K^{(2)}(\boldsymbol{x}, \boldsymbol{x}')$$

- For $[\phi(\boldsymbol{x})]_{i,j} = [\phi^{(1)}(\boldsymbol{x})]_i [\phi^{(2)}(\boldsymbol{x})]_j$,

$$\langle \phi(\boldsymbol{x}), \phi(\boldsymbol{x}') \rangle = \sum_{i,j} [\phi^{(1)}(\boldsymbol{x})]_i [\phi^{(2)}(\boldsymbol{x})]_j [\phi^{(1)}(\boldsymbol{x}')]_i [\phi^{(2)}(\boldsymbol{x}')]_j$$

$$= \langle \phi^{(1)}(\boldsymbol{x}), \phi^{(1)}(\boldsymbol{x}') \rangle \langle \phi^{(2)}(\boldsymbol{x}), \phi^{(2)}(\boldsymbol{x}') \rangle$$

$$= K^{(1)}(\boldsymbol{x}, \boldsymbol{x}') \ K^{(2)}(\boldsymbol{x}, \boldsymbol{x}')$$

# Exercise: Playing with Kernel Trick

- **Norm:**

$$\|\boldsymbol{f}_i\| = \sqrt{K(\boldsymbol{x}_i, \boldsymbol{x}_i)}$$

- **Distance:**

$$\|\boldsymbol{f}_i - \boldsymbol{f}_j\|^2 = K(\boldsymbol{x}_i, \boldsymbol{x}_i) - 2K(\boldsymbol{x}_i, \boldsymbol{x}_j) + K(\boldsymbol{x}_j, \boldsymbol{x}_j)$$

- **Angle:**

$$\cos\theta = \frac{K(\boldsymbol{x}_i, \boldsymbol{x}_j)}{\sqrt{K(\boldsymbol{x}_i, \boldsymbol{x}_i)K(\boldsymbol{x}_j, \boldsymbol{x}_j)}}$$

$$\langle \boldsymbol{f}_i, \boldsymbol{f}_j \rangle = \|\boldsymbol{f}_i\|\|\boldsymbol{f}_j\| \cos\theta$$

# Playing with Kernel Trick (cont.)

- In particular, for Gaussian kernels,

  - $\|\boldsymbol{f}_i\|^2 = 1$

  - $\|\boldsymbol{f}_i - \boldsymbol{f}_j\|^2 = 2 - 2K(\boldsymbol{x}_i, \boldsymbol{x}_j)$

  - $\cos\theta = K(\boldsymbol{x}_i, \boldsymbol{x}_j)$

$$K(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\|\boldsymbol{x} - \boldsymbol{x}'\|^2 / c^2\right)$$
$$c > 0$$

# Kernel Trick Revisited

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

■ An inner product in the feature space can be efficiently computed by the kernel function.

■ If a linear algorithm is expressed only in terms of the inner product, it can be non-linearlized by the kernel trick:

- PCA, LPP, FDA, LFDA
- K-means clustering
- Perceptron (support vector machine)

# Kernel LPP

■ Kernel LPP embedding of a sample $\boldsymbol{f}$ :

$$g = A^\top k$$

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

$$A = (\boldsymbol{\alpha}_{n-m+1} | \boldsymbol{\alpha}_{n-m+2} | \cdots | \boldsymbol{\alpha}_n)$$

● $\{\lambda_i, \boldsymbol{\alpha}_i\}_{i=1}^m$ :Sorted generalized eigenvalues and normalized eigenvectors of $\boldsymbol{KLK\alpha} = \lambda \boldsymbol{KDK\alpha}$

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

$$\boldsymbol{K} = \boldsymbol{F}^\top \boldsymbol{F} \qquad \boldsymbol{L} = \boldsymbol{D} - \boldsymbol{W}$$

$$\boldsymbol{F} = (\boldsymbol{f}_1 | \boldsymbol{f}_2 | \cdots | \boldsymbol{f}_n) \qquad \boldsymbol{D} = \mathrm{diag}(\textstyle\sum_{j=1}^n \boldsymbol{W}_{i,j})$$

■ Note: When $\boldsymbol{KDK}$ is not full-rank, it should be replaced by $\boldsymbol{KDK} + \varepsilon \boldsymbol{I}_n$ . $\varepsilon$ :small positive scalar

# Kernel LPP Embedding of Given Features

■ Kernel LPP embedding of $\{\boldsymbol{f}_i\}_{i=1}^n$ :

$$\boldsymbol{G} = \boldsymbol{A}^\top \boldsymbol{K} \qquad \boldsymbol{G} = (\boldsymbol{g}_1 | \boldsymbol{g}_2 | \cdots | \boldsymbol{g}_n)$$

■ $\boldsymbol{G}$ can be directly obtained as

$$\boldsymbol{G} = \boldsymbol{\Psi}^\top \qquad \boldsymbol{\Psi} = (\boldsymbol{\psi}_{n-m+1} | \boldsymbol{\psi}_{n-m+2} | \cdots | \boldsymbol{\psi}_n)$$

●   $\{\gamma_i, \boldsymbol{\psi}_i\}_{i=1}^n$   :Sorted eigenvalues and normalized eigenvectors of $\boldsymbol{L}\boldsymbol{\psi} = \gamma \boldsymbol{D}\boldsymbol{\psi}$

$$\gamma_1 \geq \gamma_2 \geq \cdots \geq \gamma_n \qquad \langle \boldsymbol{D}\boldsymbol{\psi}_i, \boldsymbol{\psi}_j \rangle = \delta_{i,j}$$

■ Note: When similarity matrix $\boldsymbol{W}$ is sparse, $\boldsymbol{L}$ and $\boldsymbol{D}$ are also sparse!

# Laplacian Eigenmap Embedding

$$L\psi = \gamma D\psi$$

$$L = D - W$$
$$D = \mathrm{diag}(\sum_{j=1}^{n} W_{i,j})$$

■ Definition of $L$ implies $L\mathbf{1} = \mathbf{0}$

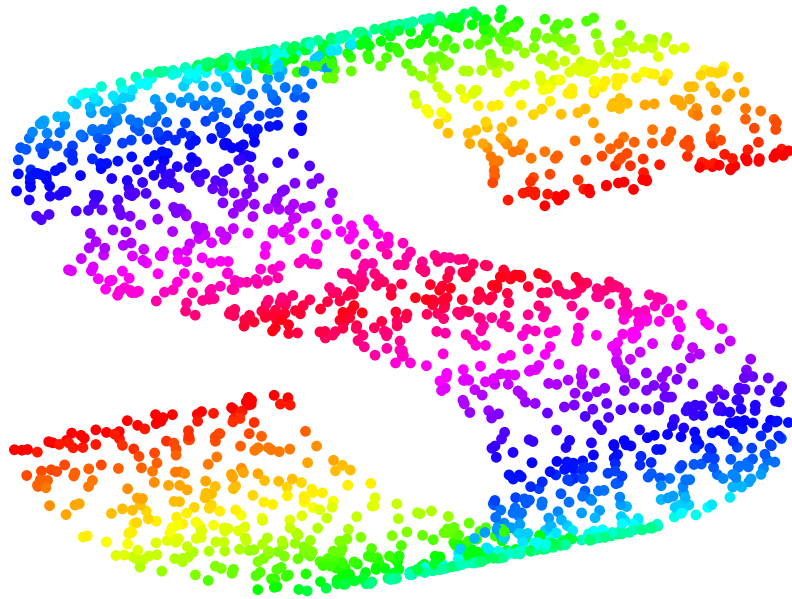$$\psi_n \propto \mathbf{1}$$

■ In practice, we remove $\psi_n$ and use

$$G = (\psi_{n-m}|\psi_{n-m+1}|\cdots|\psi_{n-1})^\top$$

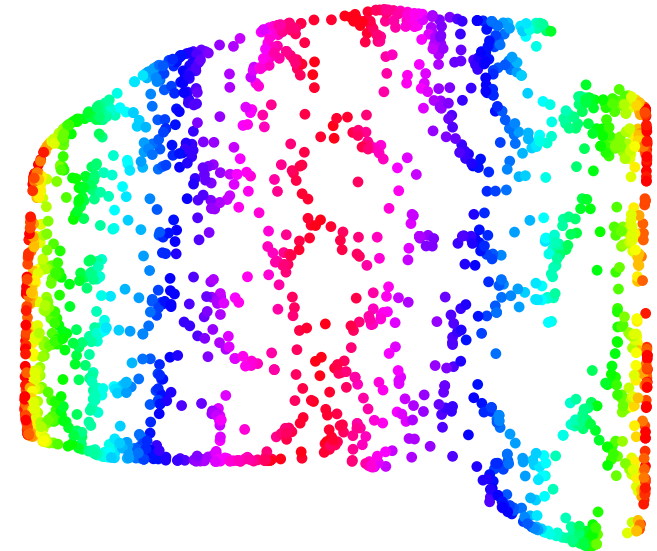■ This non-linear embedding method is called Laplacian eigenmap embedding.
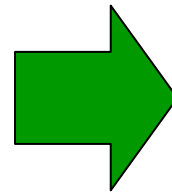
# Example

## Original data (3D)

## Embedded Data (2D)



Note: Similarity matrix is defined by the nearest-neighbor-based method with 10 nearest neighbors.

■ Laplacian eigenmap can successfully unfold the non-linear manifold.

# Kernel Tricks for Measuring Independence

- $x, y$: one-dimensional random variables.

- For a Gaussian RKHS $\mathcal{H}$, $x, y$ are independent if and only if $\rho = 0$.

$$\rho = \max_{f,g \in \mathcal{H}, \|f\|=\|g\|=1} \text{covariance}(f(x), g(y))$$

$$= \max_{f,g \in \mathcal{H}, \|f\|=\|g\|=1} \mathbb{E}[\langle f, \overline{\phi}(x) \rangle \langle g, \overline{\phi}(y) \rangle]$$

$$\overline{\phi}(x) = \phi(x) - \mathbb{E}[\phi(x)] \quad \overline{\phi}(y) = \phi(y) - \mathbb{E}[\phi(y)]$$

- Note: $\overline{\phi}(\cdot)$ also induces a reproducing kernel

# Kernel Tricks for Measuring Independence (cont.)

■ If we have samples $\{x_i\}_{i=1}^n, \{y_i\}_{i=1}^n$ ,

$$\rho \approx \max_{f,g \in \mathcal{H}, \|f\|=\|g\|=1} \left[ \frac{1}{n} \sum_{i=1}^n \langle f, \overline{\phi}(x_i) \rangle \langle g, \overline{\phi}(y_i) \rangle \right] \equiv \widehat{\rho}$$

■ Let

$$f = \sum_{i=1}^n \alpha_i \overline{\phi}(x_i) + f^\perp \qquad g = \sum_{i=1}^n \beta_i \overline{\phi}(y_i) + g^\perp$$

Then

$$\widehat{\rho} = \max_{\{\alpha_i\}_{i=1}^n, \{\beta_i\}_{i=1}^n} \left[ \frac{1}{n} \sum_{i,j,k=1}^n \alpha_i \beta_j \overline{K}(x_i, x_k) \overline{K}(y_j, y_k) \right]$$

$$\text{subject to } \sum_{i=1}^n \alpha_i^2 = \sum_{i=1}^n \beta_i^2 = 1 \qquad \overline{K}(x, x') = \langle \overline{\phi}(x), \overline{\phi}(x') \rangle$$

# Homework

1. Implement Laplacian eigenmap and unfold the 3-dimensional S-curve data.

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

   Test Laplacian eigenmap with your own (artificial or real) data and analyze its characteristics.

2. Prove that the dual eigenvalue problem of (local) Fisher discriminant analysis is given by

$$\boldsymbol{KL}^{(b)}\boldsymbol{K\alpha} = \lambda\boldsymbol{KL}^{(w)}\boldsymbol{K\alpha}$$

$$\boldsymbol{L}^{(b)} = \boldsymbol{D}^{(b)} - \boldsymbol{W}^{(b)} \qquad \boldsymbol{L}^{(w)} = \boldsymbol{D}^{(w)} - \boldsymbol{W}^{(w)}$$

$$\boldsymbol{D}^{(b)} = \operatorname{diag}(\textstyle\sum_{j=1}^{n} \boldsymbol{W}_{i,j}^{(b)}) \qquad \boldsymbol{D}^{(w)} = \operatorname{diag}(\textstyle\sum_{j=1}^{n} \boldsymbol{W}_{i,j}^{(w)})$$

$$\boldsymbol{W}_{i,j}^{(b)} = \begin{cases} 1/n - 1/n_\ell & (y_i = y_j = \ell) \\ 1/n & (y_i \neq y_j) \end{cases} \qquad \boldsymbol{W}_{i,j}^{(w)} = \begin{cases} 1/n_\ell & (y_i = y_j = \ell) \\ 0 & (y_i \neq y_j) \end{cases}$$

Note that when solving the above eigenproblem, we may need to regularize it as

$$\boldsymbol{KL}^{(b)}\boldsymbol{K\alpha} = \lambda(\boldsymbol{KL}^{(w)}\boldsymbol{K} + \epsilon\boldsymbol{I}_n)\boldsymbol{\alpha}$$

■ LFDA can also be kernelized similarly!