Advanced Data Analysis: More on Kernels

Masashi Sugiyama (Computer Science)

W8E-406, <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi

Kernel Trick

113

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

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

if K is symmetric and positive semi-definite:

$$oldsymbol{K}^ op = oldsymbol{K} \hspace{0.5mm} orall oldsymbol{y}, \hspace{0.5mm} \langle oldsymbol{K}oldsymbol{y}, oldsymbol{y}
angle \geq 0$$

Such K(x, x') is called the reproducing kernel.

Rather than directly specifying $\phi(x)$, we implicitly specify $\phi(x)$ by a reproducing kernel.

Combination of Reproducing Kernels

For any reproducing kernels (RKs) $K^{(1)}(\boldsymbol{x}, \boldsymbol{x}') = \langle \phi^{(1)}(\boldsymbol{x}), \phi^{(1)}(\boldsymbol{x}') \rangle$ $K^{(2)}(\boldsymbol{x}, \boldsymbol{x}') = \langle \phi^{(2)}(\boldsymbol{x}), \phi^{(2)}(\boldsymbol{x}') \rangle$ 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(x, x') = K⁽¹⁾(x, x') + K⁽²⁾(x, x')
 Product of RKs is still RK:

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

Proof

We prove that there exists a feature map $\phi(x)$ such that $\langle \boldsymbol{\phi}(\boldsymbol{x}), \boldsymbol{\phi}(\boldsymbol{x}') \rangle = K(\boldsymbol{x}, \boldsymbol{x}')$. For $\phi(\boldsymbol{x}) = \sqrt{\alpha} \phi^{(1)}(\boldsymbol{x})$ $\langle \boldsymbol{\phi}(\boldsymbol{x}), \boldsymbol{\phi}(\boldsymbol{x}') \rangle = \alpha \langle \boldsymbol{\phi}^{(1)}(\boldsymbol{x}), \boldsymbol{\phi}^{(1)}(\boldsymbol{x}') \rangle = \alpha K^{(1)}(\boldsymbol{x}, \boldsymbol{x}')$ For $\phi(\boldsymbol{x}) = \begin{pmatrix} \phi^{(1)}(\boldsymbol{x}) \\ \phi^{(2)}(\boldsymbol{x}) \end{pmatrix}$ $\langle \boldsymbol{\phi}(\boldsymbol{x}), \overline{\boldsymbol{\phi}(\boldsymbol{x}')}
angle = \langle \boldsymbol{\phi}^{(1)}(\boldsymbol{x}), \boldsymbol{\phi}^{(1)}(\boldsymbol{x}')
angle + \langle \boldsymbol{\phi}^{(2)}(\boldsymbol{x}), \boldsymbol{\phi}^{(2)}(\boldsymbol{x}')
angle$ $= K^{(1)}(\boldsymbol{x}, \boldsymbol{x}') + K^{(2)}(\boldsymbol{x}, \boldsymbol{x}')$ For $[\phi(\boldsymbol{x})]_{i,i} = [\phi^{(1)}(\boldsymbol{x})]_i [\phi^{(2)}(\boldsymbol{x})]_i$, $\langle \boldsymbol{\phi}(\boldsymbol{x}), \boldsymbol{\phi}(\boldsymbol{x}')
angle = \sum [\boldsymbol{\phi}^{(1)}(\boldsymbol{x})]_i [\boldsymbol{\phi}^{(2)}(\boldsymbol{x})]_j [\boldsymbol{\phi}^{(1)}(\boldsymbol{x}')]_i [\boldsymbol{\phi}^{(2)}(\boldsymbol{x}')]_j$ $=\langle \boldsymbol{\phi}^{(1)}(\boldsymbol{x}), \boldsymbol{\phi}^{(1)}(\boldsymbol{x}')
angle \langle \boldsymbol{\phi}^{(2)}(\boldsymbol{x}), \boldsymbol{\phi}^{(2)}(\boldsymbol{x}')
angle$ $= K^{(1)}(\boldsymbol{x}, \boldsymbol{x}') K^{(2)}(\boldsymbol{x}, \boldsymbol{x}')$

Exercise: Playing with Kernel Trick

Norm:

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

Distance:

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

Angle:

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

$$\langle \boldsymbol{f}, \boldsymbol{f}'
angle = \| \boldsymbol{f} \| \| \boldsymbol{f}' \| \cos heta$$



Playing with Kernel Trick (cont.)¹⁷

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

For Gaussian kernels,

•
$$\| \boldsymbol{f} \|^2 = 1$$

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

•
$$\cos \theta = K(\boldsymbol{x}, \boldsymbol{x}')$$



Kernel Trick Revisited ¹¹⁸

 $\langle \boldsymbol{f}, \boldsymbol{f}'
angle = K(\boldsymbol{x}, \boldsymbol{x}')$

- 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 nonlinearlized by the kernel trick:
 - PCA, LPP, FDA, LFDA
 - K-means clustering
 - Perceptron (support vector machine)

Kernel LPP¹¹⁹

Kernel LPP embedding of sample $f \ (= \phi(x))$:

$$\boldsymbol{g} = \boldsymbol{A}^{\top} \boldsymbol{k} \qquad \qquad \boldsymbol{k} = (K(\boldsymbol{x}, \boldsymbol{x}_1), K(\boldsymbol{x}, \boldsymbol{x}_2), \dots, K(\boldsymbol{x}, \boldsymbol{x}_n))^{\top} \\ \boldsymbol{A} = (\boldsymbol{\alpha}_{n-m+1} | \boldsymbol{\alpha}_{n-m+2} | \cdots | \boldsymbol{\alpha}_n)$$

{λ_i, α_i}^m_{i=1}:Sorted generalized eigenvalues and normalized eigenvectors of *KLKα* = λ*KDKα* λ1 ≥ λ2 ≥ ··· ≥ λn (*KDKα_i, α_j*) = δ_{i,j} *K*_{i,j} = *K*(*x*_i, *x_j*) *L* = *D* − *W D* = diag(∑ⁿ_{j=1} *W*_{i,j})
Note: When *KDK* is not full-rank, it may be replaced with *KDK* + ε*I_n*. ε :small positive scalar

120 Kernel LPP Embedding of Given Features Kernel LPP embedding of $\{f_i\}_{i=1}^n$: $\boldsymbol{G} = \boldsymbol{A}^{\top} \boldsymbol{K}$ $\boldsymbol{G} = (\boldsymbol{g}_1 | \boldsymbol{g}_2 | \cdots | \boldsymbol{g}_n)$ **G** can be directly obtained as $oldsymbol{G} = oldsymbol{\Psi}^ op oldsymbol{\Psi} = (oldsymbol{\psi}_{n-m+1} | oldsymbol{\psi}_{n-m+2} | \cdots | oldsymbol{\psi}_n)$ • $\{\gamma_i, \psi_i\}_{i=1}^n$:Sorted eigenvalues and normalized eigenvectors of $L\psi = \gamma D\psi$ $\gamma_1 \geq \gamma_2 \geq \cdots \geq \gamma_n$ $\langle oldsymbol{D}oldsymbol{\psi}_i, oldsymbol{\psi}_i
angle = \delta_{i,j}$ **Note:** When similarity matrix W is sparse, L and D are also sparse. Sparse eigenproblems can be solved efficiently.

Laplacian Eigenmap121 $L\psi = \gamma D\psi$ L = D - W $D = diag(\sum_{j=1}^{n} W_{i,j})$

Definition of L implies L1 = 0.

$$igstarrow \gamma_n=0, \ \ oldsymbol{\psi}_n \propto \mathbf{1}$$

In practice, we remove ${oldsymbol{\psi}}_n$ and use

$$oldsymbol{G} = (oldsymbol{\psi}_{n-m} | oldsymbol{\psi}_{n-m+1} | \cdots | oldsymbol{\psi}_{n-1})^ op$$

This non-linear embedding method is called Laplacian eigenmap.



Original data (3D)



Embedded Data (2D)

122



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

Laplacian eigenmap can successfully unfold the non-linear manifold.

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.

Homework (cont.) ¹²⁴

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

$$\begin{split} \boldsymbol{KL}^{(b)} \boldsymbol{K\alpha} &= \lambda \boldsymbol{KL}^{(w)} \boldsymbol{K\alpha} \\ \boldsymbol{L}^{(b)} &= \boldsymbol{D}^{(b)} - \boldsymbol{W}^{(b)} \\ \boldsymbol{D}^{(b)} &= \operatorname{diag}(\sum_{j=1}^{n} \boldsymbol{W}^{(b)}_{i,j}) \\ \boldsymbol{W}^{(b)}_{i,j} &= \begin{cases} 1/n - 1/n_{\ell} & (y_i = y_j = \ell) \\ 1/n & (y_i \neq y_j) \end{cases} \quad \begin{aligned} \boldsymbol{K\alpha} &= \lambda \boldsymbol{KL}^{(w)} \boldsymbol{K\alpha} \\ \boldsymbol{L}^{(w)} &= \boldsymbol{D}^{(w)} - \boldsymbol{W}^{(w)} \\ \boldsymbol{D}^{(w)} &= \operatorname{diag}(\sum_{j=1}^{n} \boldsymbol{W}^{(w)}_{i,j}) \\ \boldsymbol{W}^{(w)}_{i,j} &= \begin{cases} 1/n_{\ell} & (y_i = y_j = \ell) \\ 0 & (y_i \neq y_j) \end{cases} \end{aligned}$$

Note: When solving the above eigenproblem, we may practically need to regularize it as $KL^{(b)}K\alpha = \lambda(KL^{(w)}K + \epsilon I_n)\alpha$

LFDA can also be kernelized similarly!

Notification of Final Assignment

125

Data Analysis: Apply dimensionality reduction or clustering techniques to your own data set and "mine" something interesting!

Deadline: July 31st (Wed) 17:00 Pring your printed report to W/9E 40

- Bring your printed report to W8E-406.
- E-mail submission is also possible (though not recommended).

Mini-Conference on Data Analysis

126

- On July 16th and 23rd, we have a miniconference on data analysis.
- Some of the students may present their data analysis results.
- Those who give a talk at the conference will have very good grades!

Mini-Conference on Data Analysis

- Application procedure: On June 25th, just say to me "I want to give a talk!".
- Presentation: approx. 10 min (?)
 - Description of your data
 - Methods to be used
 - Outcome
- Slides should be in English.
- Better to speak in English, but Japanese may also be allowed (perhaps your friends will provide simultaneous translation!).