

# Lecture 6

## Kernel Machines – Part 2

(February 13, 2015)

Mu Zhu  
University of Waterloo

# The Kernel “Trick”

- separating hyperplane “merely” a linear classifier ... big deal
- ignited wide recognition of the “kernel trick”
  - statisticians like Grace Wahba ... long before the SVM
  - reproducing kernel Hilbert space,  $\mathcal{H}_K$
- both training and prediction stages depend on  $\mathbf{x}$ ’s only through their inner products, e.g.,  $\mathbf{x}_i^T \mathbf{x}_j$  and  $\mathbf{x}_i^T \mathbf{x}$
- can replace with  $K(\mathbf{x}_i; \mathbf{x}_j)$  and  $K(\mathbf{x}; \mathbf{x}_i)$ , which amounts to using a “different kind” of inner product

# The Kernel “Trick”

## SVM Training

$$\hat{\beta} = \sum_{i=1}^n \hat{\alpha}_i y_i \mathbf{x}_i, \quad \text{where } \hat{\alpha}_i \text{'s obtained from}$$

$$\begin{aligned} \max_{\alpha_i} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \quad [\text{swap } K(\mathbf{x}_i; \mathbf{x}_j)] \\ \text{s.t} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{and} \quad 0 \leq \alpha_i \leq \gamma \end{aligned}$$

## SVM Prediction

$$f(\mathbf{x}) = \hat{\beta}_0 + \hat{\beta}^T \mathbf{x} = \hat{\beta}_0 + \sum_{i=1}^n \hat{\alpha}_i y_i \mathbf{x}_i^T \mathbf{x} \quad [\text{swap } K(\mathbf{x}; \mathbf{x}_i)]$$

# The Kernel “Trick”

- as long as  $K(\cdot; \cdot)$  is **positive semi-definite**,  $\exists \varphi(\cdot)$  such that

$$K(\mathbf{u}; \mathbf{v}) = [\varphi(\mathbf{u})]^T [\varphi(\mathbf{v})],$$

a result due to Mercer dating back to **1909**

- so get decision boundary linear in the space of  $\varphi(\mathbf{x})$ , but nonlinear in original space
- need not know the space of  $\varphi(\mathbf{x})$  explicitly ... sometimes called the **implicit feature space**
- statisticians are used to explicitly adding higher order terms

## Example in $\mathbb{R}^2$

$$\varphi(\mathbf{u}) = \begin{bmatrix} 1 & \sqrt{2}u_1 & \sqrt{2}u_2 & u_1^2 & u_2^2 & \sqrt{2}u_1u_2 \end{bmatrix}$$

$$\begin{aligned} K(\mathbf{u}; \mathbf{v}) &= [\varphi(\mathbf{u})]^\top [\varphi(\mathbf{v})] \\ &= 1 + 2u_1v_1 + 2u_2v_2 + u_1^2v_1^2 + u_2^2v_2^2 + 2(u_1u_2)(v_1v_2) \\ &= (1 + u_1v_1 + u_2v_2)^2 \\ &= (1 + \mathbf{u}^\top \mathbf{v})^2 \end{aligned}$$

# The Radial Basis Kernel in $\mathbb{R}$

$$\phi(u) = \left[ 1 \quad u \quad \frac{u^2}{\sqrt{2!}} \quad \frac{u^3}{\sqrt{3!}} \quad \dots \right], \quad \varphi(u) = \frac{\phi(u)}{\|\phi(u)\|}$$

$$K(u; v) = [\varphi(u)]^T [\varphi(v)]$$

$$\begin{aligned} &= \frac{1 + uv + \frac{u^2 v^2}{2!} + \frac{u^3 v^3}{3!} + \dots}{\sqrt{\left(1 + u^2 + \frac{u^4}{2!} + \frac{u^6}{3!} + \dots\right)} \sqrt{\left(1 + v^2 + \frac{v^4}{2!} + \frac{v^6}{3!} + \dots\right)}} \\ &= \frac{e^{uv}}{e^{u^2/2} e^{v^2/2}} = \exp \left[ -\frac{u^2 - 2uv + v^2}{2} \right] \\ &= \exp \left[ -\frac{(u - v)^2}{2} \right] \end{aligned}$$

# Reproducing Kernel Hilbert Space

- notion of **inner product**  $\sim$  “soul” of Hilbert spaces
- $\mathcal{H}_K \sim$  a “nice” functional (Hilbert) space associated with a kernel  $K(\mathbf{u}; \mathbf{v})$  having the so-called “reproducing property”:

$$\langle f(\mathbf{x}), \underbrace{K(\mathbf{x}; \mathbf{z})}_{\text{think } g(\mathbf{x})} \rangle_{\mathcal{H}_K} = f(\mathbf{z})$$

- quite **abstract** ... a bit “easier” to appreciate if one imagines **eigen-functions** of  $K(\mathbf{u}; \mathbf{v})$  forming a basis of  $\mathcal{H}_K$ , so

$$K(\mathbf{x}; \mathbf{z}) = \sum_{\ell} \lambda_{\ell} \phi_{\ell}(\mathbf{x}) \phi_{\ell}(\mathbf{z}) \quad \text{and} \quad f(\mathbf{x}) = \sum_{\ell} c_{\ell} \phi_{\ell}(\mathbf{x})$$

# Representer Theorem

**Theorem** Consider the regularized functional estimation problem in  $\mathcal{H}_K$ ,

$$\min_{f \in \mathcal{H}_K} \sum_{i=1}^n L[y_i, f(\mathbf{x}_i)] + \lambda \|f(\mathbf{x})\|_{\mathcal{H}_K}^2.$$

The minimizing solution,  $f(\mathbf{x})$ , must have the form,

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}; \mathbf{x}_i).$$



# Inner Products

$$f(\mathbf{x}) = \sum_i a_i K(\mathbf{x}; \mathbf{x}_i), \quad g(\mathbf{x}) = \sum_j b_j K(\mathbf{x}; \mathbf{x}_j)$$

$$\langle f(\mathbf{x}), g(\mathbf{x}) \rangle_{\mathcal{H}_K} = \left\langle \underbrace{\sum_i a_i K(\mathbf{x}; \mathbf{x}_i)}_{f(\mathbf{x})}, g(\mathbf{x}) \right\rangle_{\mathcal{H}_K}$$

$$\begin{aligned} &= \sum_i a_i \underbrace{\langle K(\mathbf{x}; \mathbf{x}_i), g(\mathbf{x}) \rangle_{\mathcal{H}_K}}_{g(\mathbf{x}_i)} = \sum_i a_i \left[ \underbrace{\sum_j b_j K(\mathbf{x}_i, \mathbf{x}_j)}_{g(\mathbf{x}_i)} \right] \\ &= \sum_i \sum_j a_i b_j K(\mathbf{x}_i; \mathbf{x}_j) \end{aligned}$$

# Regularized Functional Estimation

$$\min_{f \in \mathcal{H}_K} \sum_{i=1}^n L[y_i, f(\mathbf{x}_i)] + \lambda \|f(\mathbf{x})\|_{\mathcal{H}_K}^2,$$

$$\min_{\boldsymbol{\alpha}} \sum_{i=1}^n L \left[ y_i, \sum_{j=1}^n \alpha_j K(\mathbf{x}_i; \mathbf{x}_j) \right] + \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i; \mathbf{x}_j)$$

$$\min_{\boldsymbol{\alpha}} L(\mathbf{y}, \mathbf{K}\boldsymbol{\alpha}) + \lambda \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha}$$

SVM does exactly this, with “hinge” loss

# The Gram Matrix

$$\begin{aligned} \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} &\Rightarrow \mathbf{X}\mathbf{X}^T = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_n \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{x}_1^T \mathbf{x}_1 & \mathbf{x}_1^T \mathbf{x}_2 & \dots & \mathbf{x}_1^T \mathbf{x}_n \\ \mathbf{x}_2^T \mathbf{x}_1 & \mathbf{x}_2^T \mathbf{x}_2 & \dots & \mathbf{x}_2^T \mathbf{x}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_n^T \mathbf{x}_1 & \mathbf{x}_n^T \mathbf{x}_2 & \dots & \mathbf{x}_n^T \mathbf{x}_n \end{bmatrix} \end{aligned}$$

# Kernelization

$$\mathbf{X}\mathbf{X}^T = \begin{bmatrix} \mathbf{x}_1^T \mathbf{x}_1 & \mathbf{x}_1^T \mathbf{x}_2 & \dots & \mathbf{x}_1^T \mathbf{x}_n \\ \mathbf{x}_2^T \mathbf{x}_1 & \mathbf{x}_2^T \mathbf{x}_2 & \dots & \mathbf{x}_2^T \mathbf{x}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_n^T \mathbf{x}_1 & \mathbf{x}_n^T \mathbf{x}_2 & \dots & \mathbf{x}_n^T \mathbf{x}_n \end{bmatrix}$$

↓ replace  $\mathbf{x}_i^T \mathbf{x}_j$  with  $K(\mathbf{x}_i; \mathbf{x}_j)$  ↓

$$\mathbf{K} = \begin{bmatrix} K(\mathbf{x}_1; \mathbf{x}_1) & K(\mathbf{x}_1; \mathbf{x}_2) & \dots & K(\mathbf{x}_1; \mathbf{x}_n) \\ K(\mathbf{x}_2; \mathbf{x}_1) & K(\mathbf{x}_2; \mathbf{x}_2) & \dots & K(\mathbf{x}_2; \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ K(\mathbf{x}_n; \mathbf{x}_1) & K(\mathbf{x}_n; \mathbf{x}_2) & \dots & K(\mathbf{x}_n; \mathbf{x}_n) \end{bmatrix}$$

# Ex I: Ridge Regression

## Training

$$\min_{\beta \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \underbrace{\|\beta\|^2}_{\beta^T \beta}$$

↓ reparameterize  $\beta$  as  $\mathbf{X}^T \alpha$  ↓

$$\min_{\alpha \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{X} \underbrace{\mathbf{X}^T \alpha}_{\beta}\|^2 + \lambda \underbrace{\alpha^T \mathbf{X}}_{\beta^T} \underbrace{\mathbf{X}^T \alpha}_{\beta}$$

↓ replace Gram matrix  $\mathbf{X}\mathbf{X}^T$  with  $\mathbf{K}$  ↓

$$\min_{\alpha \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{K}\alpha\|^2 + \lambda \alpha^T \mathbf{K} \alpha$$

# Ex I: Ridge Regression

## Prediction

$$\begin{aligned}\hat{y}_{new} &= \mathbf{x}_{new}^T \hat{\boldsymbol{\beta}} = \mathbf{x}_{new}^T \mathbf{X}^T \hat{\boldsymbol{\alpha}} \\ &= \mathbf{x}_{new}^T \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_n \end{bmatrix} \hat{\boldsymbol{\alpha}} = \sum_{i=1}^n (\mathbf{x}_{new}^T \mathbf{x}_i) \hat{\alpha}_i \\ &= \sum_{i=1}^n \hat{\alpha}_i K(\mathbf{x}_{new}; \mathbf{x}_i)\end{aligned}$$

**Remark** Though it appears to be optional, the ridge penalty term  $\boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha}$  is necessary because, once kernelized, the regression will invariably overfit without it.

## Ex II: Canonical Correlation Analysis

- for random variables  $x, y \in \mathbb{R}$ , their correlation coefficient is

$$\rho = \text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x)\text{Var}(y)}}$$

- for  $\mathbf{x} \in \mathbb{R}^p$  and  $\mathbf{y} \in \mathbb{R}^q$ ,

$$\rho(\mathbf{x}, \mathbf{y}) \equiv \max_{\mathbf{u} \in \mathbb{R}^p, \mathbf{v} \in \mathbb{R}^q} \text{Corr}(\mathbf{u}^T \mathbf{x}, \mathbf{v}^T \mathbf{y})$$

is called the **canonical correlation coefficient** (Hotelling, 1936; *Biometrika*)

- classic ... many (almost all) techniques in a standard textbook on “multivariate analysis” can be reduced/converted to **canonical correlation analysis** (CCA)

## Ex II: Canonical Correlation Analysis

- since

$$\begin{aligned}\text{Corr}(\mathbf{u}^T \mathbf{x}, \mathbf{v}^T \mathbf{y}) &= \frac{\text{Cov}(\mathbf{u}^T \mathbf{x}, \mathbf{v}^T \mathbf{y})}{\sqrt{\text{Var}(\mathbf{u}^T \mathbf{x}) \text{Var}(\mathbf{v}^T \mathbf{y})}} \\ &= \frac{\mathbf{u}^T \text{Cov}(\mathbf{x}, \mathbf{y}) \mathbf{v}}{\sqrt{[\mathbf{u}^T \text{Var}(\mathbf{x}) \mathbf{u}] [\mathbf{v}^T \text{Var}(\mathbf{y}) \mathbf{v}]}}\end{aligned}$$

- maximization problem for computing  $\rho(\mathbf{x}, \mathbf{y})$  equivalent to

$$\begin{aligned}& \max_{\mathbf{u} \in \mathbb{R}^p, \mathbf{v} \in \mathbb{R}^q} \quad \mathbf{u}^T [\text{Cov}(\mathbf{x}, \mathbf{y})] \mathbf{v}, \\ & \text{subject to} \quad \mathbf{u}^T [\text{Var}(\mathbf{x})] \mathbf{u} = 1 \quad \text{and} \\ & \quad \quad \quad \mathbf{v}^T [\text{Var}(\mathbf{y})] \mathbf{v} = 1\end{aligned}$$



## Ex II: Canonical Correlation Analysis

If both data matrices

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}_{n \times p} \quad \text{and} \quad \mathbf{Y} = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_n^T \end{bmatrix}_{n \times q}$$

are column centered, then, the “usual” empirical estimates are

$$\widehat{\text{Cov}}(\mathbf{x}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{y}_i^T = \frac{1}{n} \mathbf{X}^T \mathbf{Y},$$

$$\widehat{\text{Var}}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{n} \mathbf{X}^T \mathbf{X}, \quad \widehat{\text{Var}}(\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^T = \frac{1}{n} \mathbf{Y}^T \mathbf{Y}.$$

## Ex II: Canonical Correlation Analysis

- thus, empirical version of the problem is

$$\hat{\rho}(x, y) = \max_{\substack{\mathbf{u}^T \mathbf{X}^T \mathbf{X} \mathbf{u} = 1 \\ \mathbf{v}^T \mathbf{Y}^T \mathbf{Y} \mathbf{v} = 1}} \mathbf{u}^T \mathbf{X}^T \mathbf{Y} \mathbf{v},$$

a generalized **eigenvalue** problem

- compare with **principal component analysis** (PCA)

$$\underbrace{\max_{\|\mathbf{u}\|=1} \text{Var}(\mathbf{u}^T \mathbf{x})}_{\text{population version}} \quad \text{or} \quad \underbrace{\max_{\mathbf{u}^T \mathbf{u} = 1} \mathbf{u}^T \mathbf{X}^T \mathbf{X} \mathbf{u}}_{\text{sample version}},$$

a “regular” **eigenvalue** problem

## Ex II: Canonical Correlation Analysis

- eigenvalue problems:

$$\max_{\|\mathbf{a}\|=1} \mathbf{a}^T \mathbf{S} \mathbf{a} \quad \text{for symmetric } \mathbf{S}$$

$$\max_{\mathbf{a}^T \mathbf{M} \mathbf{a} = 1} \mathbf{a}^T \mathbf{S} \mathbf{a} \quad \text{for symmetric } \mathbf{S}, \mathbf{M}$$

$$\max_{\substack{\mathbf{a}^T \mathbf{M} \mathbf{a} = 1 \\ \mathbf{b}^T \mathbf{N} \mathbf{b} = 1}} \mathbf{a}^T \mathbf{S} \mathbf{b} \quad \text{for symmetric } \mathbf{M}, \mathbf{N}$$

- e.g., PCA (top), Fisher discriminants (middle); CCA (bottom)

## Ex II: Canonical Correlation Analysis

$$\hat{\rho}(x, y) = \max_{\substack{\mathbf{u}^\top (\mathbf{X}^\top \mathbf{X}) \mathbf{u} = 1 \\ \mathbf{v}^\top (\mathbf{Y}^\top \mathbf{Y}) \mathbf{v} = 1}} \mathbf{u}^\top \mathbf{X}^\top \mathbf{Y} \mathbf{v}$$

↓ reparameterize  $\mathbf{u} = \mathbf{X}^\top \boldsymbol{\alpha}$ ,  $\mathbf{v} = \mathbf{Y}^\top \boldsymbol{\theta}$  [ $\mathbf{u} \in \mathbb{R}^p$ ,  $\mathbf{v} \in \mathbb{R}^q$ ,  $\boldsymbol{\alpha}, \boldsymbol{\theta} \in \mathbb{R}^n$ ] ↓

$$\hat{\rho}(x, y) = \max_{\substack{\boldsymbol{\alpha}^\top \mathbf{X} \mathbf{X}^\top \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha} = 1 \\ \boldsymbol{\theta}^\top \mathbf{Y} \mathbf{Y}^\top \mathbf{Y} \mathbf{Y}^\top \boldsymbol{\theta} = 1}} \boldsymbol{\alpha}^\top \mathbf{X} \mathbf{X}^\top \mathbf{Y} \mathbf{Y}^\top \boldsymbol{\theta}$$

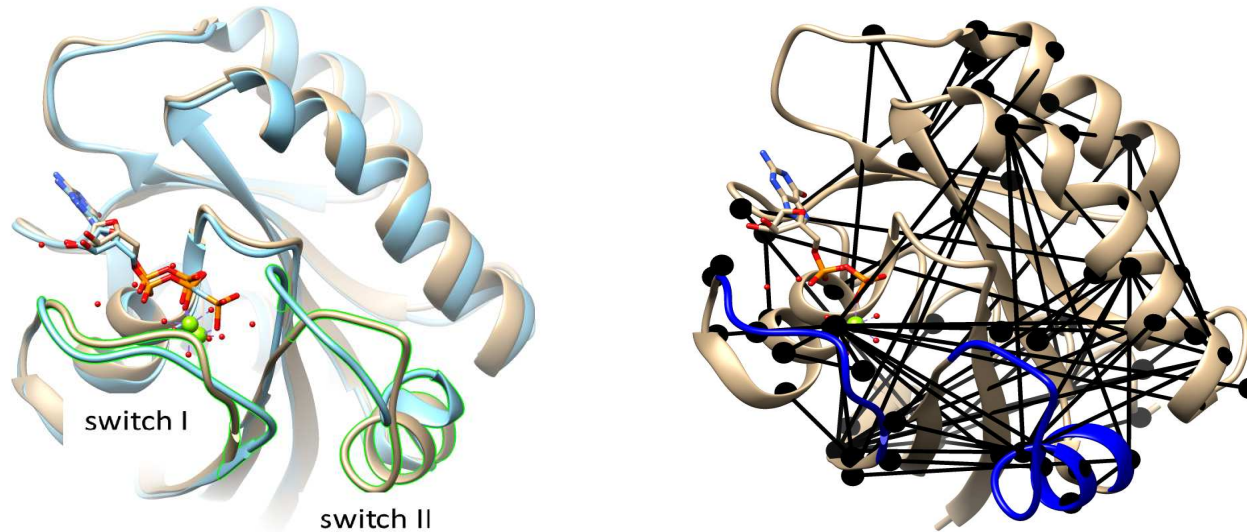
↓ replace  $\mathbf{X} \mathbf{X}^\top$  with  $\mathbf{K}_X$  and  $\mathbf{Y} \mathbf{Y}^\top$  with  $\mathbf{K}_Y + \underline{\text{regularize}}$  ↓

$$\hat{\rho}_{\mathbf{K}}(x, y) = \max_{\substack{\boldsymbol{\alpha}^\top (\mathbf{K}_X^2 + \lambda \mathbf{I}) \boldsymbol{\alpha} = 1 \\ \boldsymbol{\theta}^\top (\mathbf{K}_Y^2 + \lambda \mathbf{I}) \boldsymbol{\theta} = 1}} \boldsymbol{\alpha}^\top \mathbf{K}_X \mathbf{K}_Y \boldsymbol{\theta}$$

# Discussion

- Easier
  - to specify a good feature space,  $\varphi(\cdot)$ , directly, or
  - to specify a good similarity measure,  $K(\cdot; \cdot)$  instead?
- Many “fancy” kernels formed directly by  $[\varphi(\cdot)]^T[\varphi(\cdot)]$ , e.g.,
  - for text,
  - for protein sequences,
  - and so on.

# Application: Proteins



Want to quantify the dependence between a pair of **residues**, say  $a$  and  $b$ , in a protein. (Ongoing work by L. Soltan-Ghoraie, F. Burkowski & M. Zhu)

# Application: Proteins

- we do so by computing  $\hat{\rho}_K(\mathbf{x}_{[-z]}, \mathbf{y}_{[-z]})$ , where
  - $\mathbf{x}$  = vector of  $p$  dihedral angles for  $a$  ( $0 \leq p \leq 4$ ),
  - $\mathbf{y}$  = vector of  $q$  dihedral angles for  $b$  ( $0 \leq q \leq 4$ ),
  - $\mathbf{z}$  = vector of  $d$  dihedral angles for the rest ( $d$  large)
- need a sample  $\mathbb{D} = \{(\mathbf{x}_i, \mathbf{y}_i, \mathbf{z}_i)\}_{i=1}^n$  — i.e., different conformations of the same protein
- need a kernel function  $K(\cdot, \cdot)$ , for measuring the similarity of two conformations (expressed in terms of dihedral angles)

**Remark** The “ $\mathbf{x}_{[-z]}, \mathbf{y}_{[-z]}$ ” part is important for our application, but not so important for this course — can ignore it.

# Application: Proteins

- $\mathbb{D}$ : randomly fixed positions for 20% of the residues, and asked existing prediction algorithms to locate the rest
- $K(\cdot, \cdot)$ : for  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^p$ , used the von-Mises kernel,

$$K(\mathbf{x}_i, \mathbf{x}_j) = w_i w_j \prod_{t=1}^p \exp[\kappa_t \cos(x_{it} - x_{jt})],$$

and likewise for  $\mathbf{y}_i, \mathbf{y}_j \in \mathbb{R}^q$

- used weights  $w_i, w_j$  inversely proportional to the energy of samples  $i, j$



# Application: Proteins

$$\cos(x) \approx 1 - \frac{x^2}{2}$$

↓

von-Mises kernel

$$\begin{aligned} & \exp[\kappa_t \cos(x_{it} - x_{jt})] \\ & \approx (e^{\kappa_t}) \exp\left[-\frac{(x_{it} - x_{jt})^2}{2/\kappa_t}\right] \end{aligned}$$

like a Gaussian kernel with “standard deviation unit”

$$\sigma_t = \sqrt{\frac{1}{\kappa_t}} \times \frac{360^\circ}{2\pi}$$

$\kappa_t$  for different dihedral angles

$t$	$\kappa_t$	$\sigma_t$
1st	8	20°
2nd	8	20°
3rd	4	30°
4th	2	40°

$\sigma_t$  rounded to nearest 10-th

[current setting; still fine-tuning]

# Eigenvalue Problem

$$\max_{\|\mathbf{a}\|=1} \mathbf{a}^T \mathbf{S} \mathbf{a}$$

$$\mathbf{S} \text{ symmetric} \Rightarrow \mathbf{S} = \mathbf{U} \mathbf{D} \mathbf{U}^T \quad [\mathbf{U} \text{ orthonormal, } \mathbf{D} \text{ diagonal}]$$

$$\mathbf{a}^T \mathbf{S} \mathbf{a} = \underbrace{\mathbf{a}^T \mathbf{U}}_{\mathbf{b}^T} \mathbf{D} \underbrace{\mathbf{U}^T \mathbf{a}}_{\mathbf{b}} = \mathbf{b}^T \mathbf{D} \mathbf{b} = \sum_i d_i b_i^2$$

$$\|\mathbf{b}\|^2 = \mathbf{b}^T \mathbf{b} = \mathbf{a}^T \underbrace{\mathbf{U} \mathbf{U}^T}_{\mathbf{I}} \mathbf{a} = \mathbf{a}^T \mathbf{a} = \|\mathbf{a}\|^2$$

$$w_i = b_i^2 \Rightarrow \max_{w_1, w_2, \dots} \sum_i d_i w_i \quad \text{s.t.} \quad w_i \geq 0 \quad \text{and} \quad \sum_i w_i = 1$$

# Eigenvalue Problem

$$d_1 \geq d_2 \geq \dots \Rightarrow \hat{w}_1 = 1, \hat{w}_2 = \hat{w}_3 = \dots = 0$$

$$\hat{\mathbf{b}} = (1, 0, \dots, 0)^T$$

$$\begin{aligned} \hat{\mathbf{a}} = \mathbf{U}\hat{\mathbf{b}} &= \text{first column of } \mathbf{U} \\ &= \text{first eigenvector of } \mathbf{S} \end{aligned}$$

$$\begin{aligned} \hat{\mathbf{a}}^T \mathbf{S} \hat{\mathbf{a}} &= d_1 \\ &= \text{largest eigenvalue of } \mathbf{S} \end{aligned}$$

# Generalized Eigenvalue Problem

$$\underbrace{\max_{\substack{\mathbf{a}^T \mathbf{M} \mathbf{a} = 1 \\ \mathbf{b}^T \mathbf{N} \mathbf{b} = 1}} \mathbf{a}^T \mathbf{S} \mathbf{b}} \Leftrightarrow \underbrace{\max_{\|\mathbf{c}\| = \|\mathbf{d}\| = 1} \mathbf{c}^T \mathbf{M}^{-1/2} \mathbf{S} \mathbf{N}^{-1/2} \mathbf{d}} \\ \mathbf{c} = \mathbf{M}^{1/2} \mathbf{a}, \quad \mathbf{d} = \mathbf{N}^{1/2} \mathbf{b}$$

solve RHS problem for  $\hat{\mathbf{c}}$  and  $\hat{\mathbf{d}}$ , and then obtain

$$\hat{\mathbf{a}} = \mathbf{M}^{-1/2} \hat{\mathbf{c}} \quad \text{and} \quad \hat{\mathbf{b}} = \mathbf{N}^{-1/2} \hat{\mathbf{d}}$$

## Exercise

(a) Show that the solution to the RHS problem above is:

$$\begin{aligned} \hat{\mathbf{c}} &= \text{first left singular vector of } \mathbf{M}^{-1/2} \mathbf{S} \mathbf{N}^{-1/2}; \\ \hat{\mathbf{d}} &= \text{first right singular vector of } \mathbf{M}^{-1/2} \mathbf{S} \mathbf{N}^{-1/2}. \end{aligned}$$

(b) Why do we have to regularize kernel CCA?

# Summary

- key ideas:
  - margin; generalization error
  - constrained optimization; KKT conditions
  - regularized functional estimation in  $\mathcal{H}_K$
  - kernelization (of algorithms driven by inner products)
- specific methods:
  - SVM
  - kernel ridge regression; kernel CCA
- applications:
  - long-range dependences in a protein molecule