

# Kernel Machines

Alain Celisse

SAMM

Paris 1-Panthéon Sorbonne University

`alain.celisse@univ-paris1.fr`

## *Lecture 3: Kernelizing classical strategies*

---

Master 2 Data Science – Centrale Lille, Lille University  
Fall 2022

## Successive topics of the coming lectures:

1. Introduction to Kernel methods
2. Support vector classifiers and Kernel methods
3. Extending classical strategies to high dimension  
(Today!)
  - ▶ KPCA
  - ▶ KRR
4. Duality gap and KKT conditions
5. Designing reproducing kernels
6. Maximum Mean Discrepancy (MMD)
7. Change-point detection, KCP

# Outline of the lecture

Kernel Machines

Alain Celisse

Introduction

Kernel PCA

Kernel Ridge  
Regression

- ▶ KPCA
- ▶ KRR

Introduction

Motivation

Extensions of classical  
strategies

Kernel PCA

Kernel Ridge  
Regression

# Introduction

# Limitations of classical learning approaches

- ▶ Cannot easily deal with “complex structured” data (metric to be defined)
  - ▶ Combining real measures with histograms, categorical data, and graphs is a hard problem
  - ▶ Extracting relevant information requires particular metric
- ▶ Easy to understand (linear classifier), but not versatile tools (see the SVM lecture)
- ▶ Measuring the dependence is not an easy task
  - ▶ Covariance and correlation measure linear dependence between variables  
( $\text{Cov}(X, Y) = 0$  does not imply  $X \perp Y$ )
  - ▶ General measures of dependence require comparison between marginal and joint distributions

## Combining kernels

- ▶  $k_1, k_2$ : reproducing kernels  $\Rightarrow \alpha k_1 + \beta k_2$ : reproducing kernel ( $\alpha, \beta \geq 0$ )
- ▶  $(k_1)^\alpha$  ( $\alpha > 0$ ): reproducing kernel

## Kernels for structured objects

- ▶ The  $\chi^2$ -kernel deals with histograms ( $G$  bins)

$$k(x, y) = \exp \left[ - \sum_{g=1}^G \frac{(x_g - y_g)^2}{x_g + y_g} \right]$$

- ▶ Helpful in video streams analysis, dealing with texts (bag of words), ...

# Versatile tool improving the performance

Kernel Machines

Alain Celisse

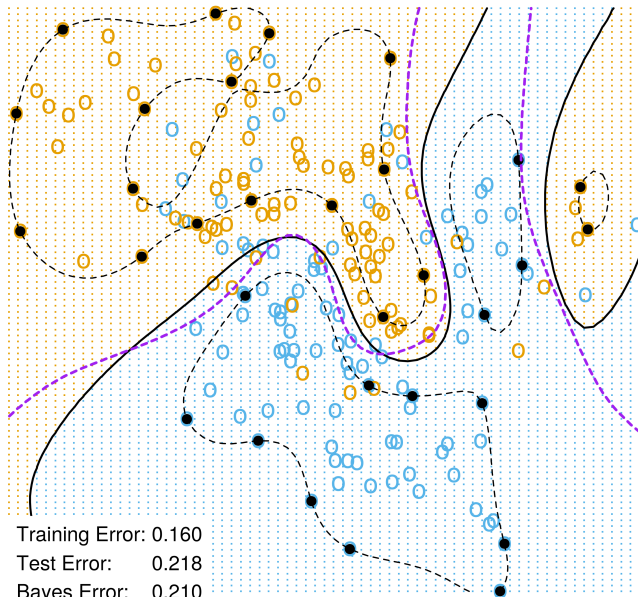
## Introduction

### Motivation

Extensions of classical strategies

### Kernel PCA

### Kernel Ridge Regression



# Measuring the dependence between variables with kernels

►  $X_1, \dots, X_n \in \mathbb{R}^d$

► Empirical covariance matrix:

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X_i X_i^\top$$

► Empirical covariance operator (from  $\mathcal{H}$  to  $\mathcal{H}$ ):

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\phi(X_i) \otimes \phi(X_i))$$

$\phi(X_i) \in \mathcal{H}$ : Extended feature vector (centered)

$\otimes$ : tensor product defined by

$$(a \otimes b)g = \langle b, g \rangle_{\mathcal{H}} a \in \mathcal{H}, \quad \forall g \in \mathcal{H}$$

► Captures non-linear dependencies between the  $d$  features of the  $X_i$ s



There are two classical lines of thought for extending a learning strategy to a kernelized version:

- ▶ Strategies relying on scalar products can be easily extended to the "kernel world"
  - ▶ SV classifier, Ridge regression, exponential families,...
- ▶ Strategies relying on similarity measures between "individuals" can be combined with any psd kernel
  - ▶  $K$ -means,  $k$ -Nearest Neighbors, Spectral clustering,...

- ▶ Kernelized Tikhonov regularization (LS-SVM, or Kernel Ridge Regression)
- ▶ Kernel-principal component Analysis
- ▶ Kernel-Canonical Correlation Analysis
- ▶ MMD and HSIC criteria (coming lecture. . . )
  
- ▶ Kernelized Support Vector Classifier (SVM)
- ▶ Multi-task learning
- ▶  $K$ -means, Kernelized HAC, Kernelized spectral clustering
- ▶ ...

# Kernel PCA (KPCA)

# Displaying the data in high dimension

Kernel Machines

Alain Celisse

Introduction

Kernel PCA

Motivation

Covariance matrix

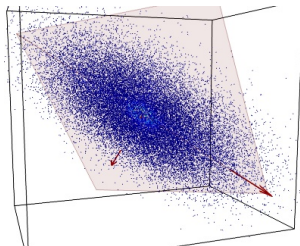
Eigenvalue  
decomposition

principal components

Kernelized PCA

Kernel covariance  
operator

Kernel Ridge  
Regression



- ▶ A large part of collected data is high-dimensional
- ▶ Displaying data in  $\mathbb{R}^d$  is challenging if  $d \geq 4$

## Strategy

- ▶ Looking for directions along which the data exhibit the largest variability
- ▶ Use these directions as a new vector basis for drawing graphs
- ▶ These privileged directions allow for visually identifying hidden structures (clusters)

## Definition (Covariance matrix)

- ▶  $X = (X^1, \dots, X^d)^\top \in \mathbb{R}^d$ : column vector
- ▶  $X^\top$ : the transpose of  $X$ , row vector
- ▶ The covariance matrix is defined by

$$\mathbb{V}(X) = \mathbb{E} \left[ (X - \mathbb{E}(X)) \cdot (X - \mathbb{E}(X))^\top \right]$$

Rks:

- ▶ "Covariance Matrix":

$$\begin{aligned} [\mathbb{V}(X)]_{i,j} &= \mathbb{E} \left[ (X^i - \mathbb{E}(X^i)) \cdot (X^j - \mathbb{E}(X^j)) \right] \\ &= \text{Cov}(X^i, X^j) \end{aligned}$$

- ▶ The covariance matrix captures some dependence between the variables  $X^j$ s
- ▶ If  $X \in \mathbb{R}^d$  is Gaussian,  $\text{Cov}(X^i, X^j) = 0 \Leftrightarrow X^i \perp X^j$

# Estimating the covariance matrix

With  $\mathbb{E}(X) = 0$ ,

$$\mathbb{V}(X) = \mathbb{E} \left[ X \cdot X^\top \right] \in \mathcal{M}_d(\mathbb{R})$$

- ▶  $X \in \mathbb{R}^d$ : column vector
- ▶  $\mathbb{V}(X)$ :  $d \times d$  matrix, positive semidefinite (PSD)

## Definition (Empirical covariance matrix)

If  $\mathbb{E}(X) = 0$ ,

$$\widehat{\mathbb{V}(X)} = \frac{1}{n} \sum_{i=1}^n X_i \cdot X_i^\top = \frac{1}{n} X^\top X$$

Rk:

- ▶ In practice, start by centering the  $X_i$ s!
- ▶ If  $d = 1$ ,  $\widehat{\mathbb{V}(X)} = 1/n \sum_{i=1}^n X_i^2 \approx \mathbb{V}(X) = \mathbb{E} [X^2]$
- ▶  $O(nd^2)$  "elementary operations" to be computed
- ▶ With  $d \gg 1$  (large), becomes computationally heavy

# Classical PCA

## Goal

- Find unit vector  $v_1 \in \mathbb{R}^d$  such that

$$\begin{aligned} v_1 \in \operatorname{Arg} \max_{\|u\|=1} \left\{ \operatorname{Var}(X^\top u) \right\} &= \operatorname{Arg} \max_{\|u\|=1} \left\{ \mathbb{E} \left[ \left( X^\top u \right)^2 \right] \right\} \\ &= \operatorname{Arg} \max_{\|u\|=1} \left\{ u^\top \mathbb{V}(X) u \right\} \end{aligned}$$

- Repeat with  $v_{j+1} \in \operatorname{Vect}(v_1, \dots, v_j)^\perp$ , for  $1 \leq j \leq d$

## Eigenvalues and eigenvectors

- Amounts to find the largest eigenvalue (and the corresponding eigenvector) of the PSD matrix  $\mathbb{V}(X)$
- Amounts to find successive basis vectors (eigenvectors) maximizing the variance at each iteration



# PCA: Eigenvalue decomposition

$\mathbb{V}(X)$ :  $d \times d$  psd matrix

## Theorem

For any psd  $d \times d$  matrix  $\Sigma$ , there exist

- ▶  $O$ :  $d \times d$  orthogonal matrix ( $O^\top \cdot O = O \cdot O^\top = I_d$ )
- ▶  $\Lambda$ :  $d \times d$  diagonal matrix ( $\Lambda_{i,i} = \lambda_i$ ,  $\Lambda_{i,j} = 0$ ,  $i \neq j$ )
- ▶  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$

such that

$$\Sigma = O \cdot \Lambda \cdot O^\top = \sum_{j=1}^d \lambda_j O_{\cdot j} O_{\cdot j}^\top$$

## Comments:

- ▶  $O$ : rotation matrix from canonical to the new basis
- ▶ Each column of  $O_{\cdot j}$  is an eigenvector of  $\Sigma$
- ▶ Each diagonal element  $\lambda_i$  is an eigenvalue of  $\Sigma$
- ▶ Each  $O_{\cdot j}$  is the eigenvector of  $\lambda_j$

# PCA: Diagonalizing $\mathbb{V}(X)$

## Application to $\mathbb{V}(X)$

With  $\mathbb{V}(X) = V\Lambda V^\top$ , ( $V$ : orthogonal)

$$w \in \operatorname{Arg} \max_{\|u\|=1} \left\{ u^\top \mathbb{V}(X) u \right\}$$

$$\Leftrightarrow w \in \operatorname{Arg} \max_{\|u\|=1} \left\{ u^\top \sum_{j=1}^d \lambda_j V_{\cdot j} V_{\cdot j}^\top u \right\}$$

$$\Leftrightarrow w \in \operatorname{Arg} \max_{\|u\|=1} \left\{ \sum_{j=1}^d \lambda_j \left( V_{\cdot j}^\top u \right)^2 \right\}$$

$$\Leftrightarrow w \in \left\{ u \in \mathbb{R}^d \mid \mathbb{V}(X)u = \lambda_1 u \right\} \cap B_{\|\cdot\|}(1)$$

## Key ingredient

( $v_j = V_{\cdot j}$ )

Finding the  $\lambda_i$ s and eigenvectors  $v_i$ s such that

$$\mathbb{V}(X) = \sum_{j=1}^d \lambda_j v_j \cdot v_j^\top$$

### Introduction

### Kernel PCA

Motivation

Covariance matrix

Eigenvalue  
decomposition

principal components

Kernelized PCA

Kernel covariance  
operator

### Kernel Ridge Regression

# PCA: Diagonalizing $\widehat{\mathbb{V}}(\mathbf{X})$

## Application to $\widehat{\mathbb{V}}(\mathbf{X})$ (empirical covariance matrix)

With  $\widehat{\mathbb{V}}(\mathbf{X}) = \widehat{\mathbf{V}} \widehat{\Lambda} \widehat{\mathbf{V}}^\top$ , ( $\widehat{\mathbf{V}}$ : orthogonal)

- Find the  $\widehat{\lambda}_j$ s and  $\widehat{\mathbf{v}}_j$ s such that

$$\widehat{\mathbb{V}}(\mathbf{X}) = \sum_{j=1}^d \widehat{\lambda}_j \widehat{\mathbf{v}}_j \cdot \widehat{\mathbf{v}}_j^\top \quad (\text{SVD of } \widehat{\mathbb{V}}(\mathbf{X}))$$

- The “principal component” of individual  $i$  on the component  $j$  is the score

$$\mathbf{X}_i^\top \cdot \widehat{\mathbf{v}}_j \in \mathbb{R} \quad (\mathbf{X}_i \text{ centered})$$

- The “ $j$ th principal component” of all individuals is

$$\tilde{\mathbf{X}}_{\cdot j} = \mathbf{X} \cdot \widehat{\mathbf{v}}_j \in \mathbb{R}^n \quad (\text{projection})$$

- principal component decomposition:

$$\tilde{\mathbf{X}} = \mathbf{X} \cdot \widehat{\mathbf{V}} \in \mathcal{M}_{n,d}(\mathbb{R}) \quad (\text{Score matrix})$$

### Introduction

### Kernel PCA

Motivation

Covariance matrix

Eigenvalue  
decomposition

principal components

Kernelized PCA

Kernel covariance  
operator

### Kernel Ridge Regression

## Summarizing the information

From

$$\tilde{X} = X \cdot \hat{V} \in \mathcal{M}_{n,d}(\mathbb{R})$$

- ▶ The row  $i$  of  $\tilde{X}$  describes individual  $i$  in the new basis  $\hat{V}$
- ▶ Displaying component  $\ell$  versus component  $k$  consists in displaying the point  $(\tilde{X}_{i,k}, \tilde{X}_{i,\ell})$  for each individual  $i$

## Dimension reduction

- ▶ From  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_d \geq 0$ , the relevant information about the cloud of points can be summarized by the  $r$  largest components
- ▶ To ease the storage of the data when  $d \gg 1$ , "build a smaller summary of the data" by means of

$$\tilde{X}^{(r)} = \left[ \tilde{X}_{\cdot,1}, \tilde{X}_{\cdot,2}, \dots, \tilde{X}_{\cdot,r} \right]$$

Introduction

Kernel PCA

Motivation

Covariance matrix

Eigenvalue  
decomposition

principal components

Kernelized PCA

Kernel covariance  
operatorKernel Ridge  
Regression

Introduction

Kernel PCA

Motivation

Covariance matrix

Eigenvalue  
decomposition

principal components

Kernelized PCA

Kernel covariance  
operator

Kernel Ridge  
Regression

# K-PCA

## Projection of extended features

- ▶ For  $i = 1, \dots, n$ ,  $\phi(x_i) = k_{x_i} \in \mathcal{H}$ : Extended features
- ▶ **Warning:** The  $k_{x_i}$ s assumed to be centered in what follows
- ▶ Orthogonal projection of  $k_{x_i}$  onto  $g \in \mathcal{H}$ :

$$\left\langle k_{x_i}, \frac{g}{\|g\|_{\mathcal{H}}} \right\rangle_{\mathcal{H}} \cdot \frac{g}{\|g\|_{\mathcal{H}}}$$

## Population variance of the projection

$$\text{Var} \left( \frac{\langle k_X, g \rangle_{\mathcal{H}}}{\|g\|_{\mathcal{H}}} \right) = \mathbb{E} \left[ \frac{(g(X))^2}{\|g\|_{\mathcal{H}}^2} \right]$$

## Empirical variance of the projection

$$\widehat{\text{Var}} \left( \frac{\langle k_X, g \rangle_{\mathcal{H}}}{\|g\|_{\mathcal{H}}} \right) = \widehat{\mathbb{E}} \left[ \frac{(g(X))^2}{\|g\|_{\mathcal{H}}^2} \right] = \frac{1}{n} \sum_{i=1}^n \frac{(g(x_i))^2}{\|g\|_{\mathcal{H}}^2}$$

### Introduction

### Kernel PCA

Motivation

Covariance matrix

Eigenvalue  
decomposition

principal components

### Kernelized PCA

Kernel covariance  
operator

### Kernel Ridge Regression

## General Kernel PCA algorithm

- Find a vector  $g_1 \in \mathcal{H}$  such that

$$\begin{aligned} g_1 &\in \operatorname{Arg\,max}_{g \in \mathcal{H}} \left\{ \widehat{\operatorname{Var}} \left( \frac{\langle k_X, g \rangle_{\mathcal{H}}}{\|g\|_{\mathcal{H}}} \right) \right\} \\ &= \operatorname{Arg\,max}_{g \in \mathcal{H}} \underbrace{\left\{ \frac{1}{n} \sum_{i=1}^n \frac{(g(x_i))^2}{\|g\|_{\mathcal{H}}^2} \right\}}_{= -\Psi(g(x_1), \dots, g(x_n), \|g\|_{\mathcal{H}})} \end{aligned}$$

## General Kernel PCA algorithm

- Find a vector  $g_1 \in \mathcal{H}$  such that

$$\begin{aligned} g_1 &\in \operatorname{Arg} \max_{g \in \mathcal{H}} \left\{ \widehat{\operatorname{Var}} \left( \frac{\langle k_X, g \rangle_{\mathcal{H}}}{\|g\|_{\mathcal{H}}} \right) \right\} \\ &= \operatorname{Arg} \max_{g \in \mathcal{H}} \underbrace{\left\{ \frac{1}{n} \sum_{i=1}^n \frac{(g(x_i))^2}{\|g\|_{\mathcal{H}}^2} \right\}}_{= -\Psi(g(x_1), \dots, g(x_n), \|g\|_{\mathcal{H}})} \end{aligned}$$

(Representer theorem)



## General Kernel PCA algorithm

- Find a vector  $g_1 \in \mathcal{H}$  such that

$$\begin{aligned} g_1 &\in \operatorname{Arg} \max_{g \in \mathcal{H}} \left\{ \widehat{\operatorname{Var}} \left( \frac{\langle k_X, g \rangle_{\mathcal{H}}}{\|g\|_{\mathcal{H}}} \right) \right\} \\ &= \operatorname{Arg} \max_{g \in \mathcal{H}} \underbrace{\left\{ \frac{1}{n} \sum_{i=1}^n \frac{(g(x_i))^2}{\|g\|_{\mathcal{H}}^2} \right\}}_{= -\Psi(g(x_1), \dots, g(x_n), \|g\|_{\mathcal{H}})} \end{aligned}$$

(Representer theorem)

- Repeat with  $g_{j+1} \in \operatorname{Vect}(g_1, \dots, g_j)^\perp$ , for  $1 \leq j \leq n$

# Simplified Kernel PCA algorithm

Kernel Machines

Alain Celisse

## Representer theorem

- ▶ The solution  $\hat{g}_1 = \sum_{i=1}^n \hat{\alpha}_i k_{x_i} \in \mathcal{H}$
- ▶  $\|\hat{g}_1\|_{\mathcal{H}}^2 = \hat{\alpha}^\top K \hat{\alpha}$
- ▶  $\sum_{i=1}^n \hat{g}_1(x_i)^2 = \hat{\alpha}^\top K^2 \hat{\alpha}$

Solving the problem amounts to compute:

## Practical Kernel PCA algorithm

- ▶ Compute

$$\hat{\alpha}_1 \in \underset{\alpha \in \mathbb{R}^n}{\text{Arg max}} \left\{ \frac{\alpha^\top K^2 \alpha}{\alpha^\top K \alpha} \right\}$$

- ▶ Repeat with  $\hat{\alpha}_{j+1} \perp_K \{\hat{\alpha}_1, \dots, \hat{\alpha}_j\}$  where

$$\alpha \perp_K \beta \quad \Leftrightarrow \quad \alpha^\top K \beta = 0$$

## Conclusion

Solving KPCA means diagonalizing the **centered** Gram matrix

Introduction

Kernel PCA

Motivation

Covariance matrix

Eigenvalue  
decomposition

principal components

Kernelized PCA

Kernel covariance  
operator

Kernel Ridge  
Regression

- ▶ The eigenvectors are now functions in  $\mathcal{H}$

$$\hat{g}_j = \sum_{i=1}^n \hat{\alpha}_{j,i} k_{x_i}, \quad \forall j \geq 1$$

- ▶ The “principal component of individual  $i$  on the component  $j$  is the score

$$\frac{\langle k_{x_i}, \hat{g}_j \rangle_{\mathcal{H}}}{\|\hat{g}_j\|_{\mathcal{H}}}$$

- ▶ The “ $j$ th principal component” of all individuals is

$$(\hat{g}_j(x_1), \dots, \hat{g}_j(x_n))^{\top} \times \frac{1}{\|\hat{g}_j\|_{\mathcal{H}}}$$

- ▶ Any individual  $x$  can be displayed through its  $(i, j)$ th coordinates in the new basis by

$$\left( \frac{\langle k_x, \hat{g}_i \rangle_{\mathcal{H}}}{\|\hat{g}_i\|_{\mathcal{H}}}, \frac{\langle k_x, \hat{g}_j \rangle_{\mathcal{H}}}{\|\hat{g}_j\|_{\mathcal{H}}} \right)$$

# Towards the covariance operator

For all  $g \in H$ , ( $k_X$ : centered)

$$\text{Var} \left( \left\langle k_X, \frac{g}{\|g\|_{\mathcal{H}}} \right\rangle_{\mathcal{H}} \right) = \mathbb{E} \left[ \left\langle k_X, \frac{g}{\|g\|_{\mathcal{H}}} \right\rangle_{\mathcal{H}}^2 \right]$$

## Tensor product

For all  $a, b \in \mathcal{H}$

$$(a \otimes b) : f \mapsto (a \otimes b)f = \langle b, f \rangle_{\mathcal{H}} a$$

## Proposition

For all  $f, g \in \mathcal{H}$

$$\langle (a \otimes b)f, g \rangle_{\mathcal{H}} = \langle a, f \rangle_{\mathcal{H}} \cdot \langle b, g \rangle_{\mathcal{H}}$$

Then

$$\langle (a \otimes a)f, f \rangle_{\mathcal{H}} = \langle a, f \rangle_{\mathcal{H}}^2$$

For all  $g \in H$ , ( $k_X$ : centered)

$$\begin{aligned}\mathbb{E} \left[ \left\langle k_X, \frac{g}{\|g\|_{\mathcal{H}}} \right\rangle_{\mathcal{H}}^2 \right] &= \mathbb{E} \left[ \left\langle (k_X \otimes k_X) \frac{g}{\|g\|_{\mathcal{H}}}, \frac{g}{\|g\|_{\mathcal{H}}} \right\rangle_{\mathcal{H}} \right] \\ &= \left\langle \mathbb{E} [k_X \otimes k_X] \frac{g}{\|g\|_{\mathcal{H}}}, \frac{g}{\|g\|_{\mathcal{H}}} \right\rangle_{\mathcal{H}} \\ &= \left\langle \Sigma \frac{g}{\|g\|_{\mathcal{H}}}, \frac{g}{\|g\|_{\mathcal{H}}} \right\rangle_{\mathcal{H}}\end{aligned}$$

## Definition

Covariance operator: Unique operator  $\Sigma$  from  $\mathcal{H}$  to  $\mathcal{H}$  such that, for all  $f, g \in \mathcal{H}$ ,

$$\begin{aligned}\langle \Sigma f, g \rangle_{\mathcal{H}} &= \mathbb{E} [\langle k_X, f \rangle_{\mathcal{H}} \langle k_X, g \rangle_{\mathcal{H}}] \\ &= \langle \mathbb{E} [k_X \otimes k_X] f, g \rangle_{\mathcal{H}}\end{aligned}$$

## Introduction

### Kernel PCA

Motivation

Covariance matrix

Eigenvalue  
decomposition

principal components

Kernelized PCA

Kernel covariance  
operator

### Kernel Ridge Regression

## Definition

Empirical covariance operator: Unique operator  $\hat{\Sigma}$  from  $\mathcal{H}$  to  $\mathcal{H}$  such that, for all  $f, g \in \mathcal{H}$ ,

$$\begin{aligned}\langle \hat{\Sigma}f, g \rangle_{\mathcal{H}} &= \frac{1}{n} \sum_{i=1}^n \langle k_{X_i}, f \rangle_{\mathcal{H}} \langle k_{X_i}, g \rangle_{\mathcal{H}} \\ &= \left\langle \frac{1}{n} \sum_{i=1}^n (k_{X_i} \otimes k_{X_i}) f, g \right\rangle_{\mathcal{H}}\end{aligned}$$

Remark:

$$\langle \hat{\Sigma}f, g \rangle_{\mathcal{H}} = \frac{1}{n} \sum_{i=1}^n f(X_i)g(X_i) \xrightarrow[n \rightarrow +\infty]{P} \langle \Sigma f, g \rangle_{\mathcal{H}} = \mathbb{E}[f(X)g(X)]$$

Introduction

Kernel PCA

Motivation

Covariance matrix

Eigenvalue  
decomposition

principal components

Kernelized PCA

Kernel covariance  
operatorKernel Ridge  
Regression

# Kernel Ridge Regression (KRR)

## Linear regression model

$$Y = X \cdot \beta^* + \epsilon \in \mathbb{R}^n$$

with  $\beta^* \in \mathbb{R}^d$ ,  $Y = (Y_1, \dots, Y_n)^\top$ ,  $\epsilon = (\epsilon_1, \dots, \epsilon_n)^\top$ , and  $X = [X^1, \dots, X^d]$ :  $n \times d$  matrix

### Assumptions:

- ▶  $\mathbb{E}_\epsilon[Y] = X \cdot \beta^*$
- ▶  $\text{Var}_\epsilon(\epsilon) = \sigma^2 I_n$ ,  $\sigma^2 > 0$

### Question

How to efficiently predict  $Y$  at a new  $X \in \mathbb{R}^d$ ?



# $\ell_0$ -norm and lack of convexity

## Ideal optimization problem to solve

For every  $k \geq 1$ , solve

$$\hat{\beta}_k = \underset{\|\beta\|_0 \leq k}{\operatorname{Arg\,min}} \left\{ \frac{1}{n} \sum_{i=1}^n \left( Y_i - X_i^\top \beta \right)^2 \right\}$$

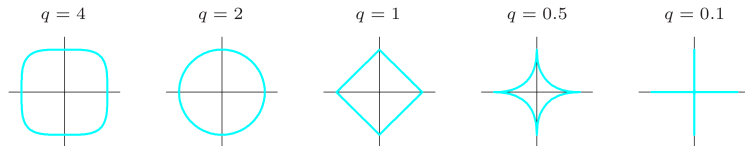
which is equivalent to:

For every  $\lambda > 0$ , solve

$$(\|\beta\|_0 = \sum_{j=1}^d \mathbb{1}_{\beta_j \neq 0})$$

$$\hat{\beta}_{NP,\lambda} = \underset{\beta \in \mathbb{R}^d}{\operatorname{Arg\,min}} \left\{ \frac{1}{n} \sum_{i=1}^n \left( Y_i - X_i^\top \beta \right)^2 + \lambda \|\beta\|_0 \right\}$$

**Remark:**  $\beta \mapsto \|\beta\|_0$  is non-convex!  $\Rightarrow$  NP-hard to solve



# Convex relaxation of the $\ell_0$ -norm

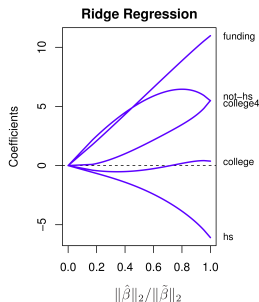
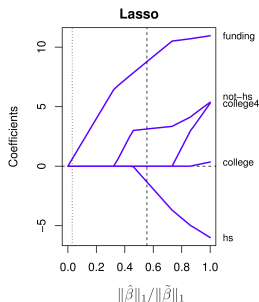
## Convex relaxation with $\ell_2$ -norm: Ridge

For every  $\lambda > 0$ , solve

$$(\|\beta\|_2^2 = \sum_{j=1}^d |\beta_j|^2)$$

$$\begin{aligned}\hat{\beta}_\lambda &= \underset{\beta \in \mathbb{R}^d}{\text{Arg min}} \left\{ \frac{1}{n} \sum_{i=1}^n \left( Y_i - X_i^\top \beta \right)^2 + \lambda \|\beta\|_2^2 \right\} \\ &= \underset{\beta \in \mathbb{R}^d}{\text{Arg min}} \left\{ \|Y - X\beta\|_n^2 + \lambda \|\beta\|_2^2 \right\}\end{aligned}$$

where  $\|u\|_n^2 = \sum_{i=1}^n u_i^2 / n$



- ▶ Unlike LASSO, Ridge leads to a closed-form expression for the estimator
- ▶ The  $\ell_2$ -norm differentiable everywhere (unlike the  $\ell_1$ -norm)

Ridge estimator: Closed-form expression

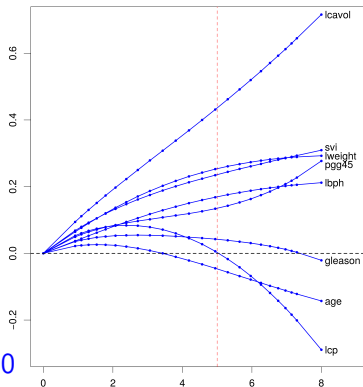
$$\hat{\beta}_\lambda = \left( \frac{\mathbf{X}^\top \mathbf{X}}{n} + \lambda \mathbf{I}_d \right)^{-1} \frac{\mathbf{X}^\top \mathbf{Y}}{n}, \quad \forall \lambda > 0$$

Remark:

- ▶  $\frac{\mathbf{X}^\top \mathbf{X}}{n} + \lambda \mathbf{I}_d$ : always invertible
- ▶  $\frac{\mathbf{X}^\top \mathbf{X}}{n}$ :  $d \times d$  matrix  $\Rightarrow$  Inversion is time consuming
- ▶ SVD of  $\mathbf{X}\mathbf{X}^\top$  can be more efficient than that of  $\mathbf{X}^\top \mathbf{X}$
- ▶ Ridge shrinks the coefficients (as LASSO does)

# Ridge: No sparsity constraint

Continuous shrinking towards 0



$$\hat{\beta}_{\lambda} = \left( \frac{\mathbf{X}^T \mathbf{X}}{n} + \lambda \mathbf{I}_d \right)^{-1} \frac{\mathbf{X}^T \mathbf{Y}}{n}, \quad \forall \lambda > 0$$

- ▶ Coefficient of  $\hat{\beta}_{\lambda}$  are shrunk continuously as  $\lambda$  grows (continuous function)
- ▶ Contrasts with the LASSO estimator

## Model (Reminder)

$$Y_i = \underbrace{X_i^\top \beta}_{=f_\beta(X_i)} + \epsilon_i$$

## Ridge estimator (Reminder)

$$\hat{\beta}_\lambda = \left( \frac{X^\top X}{n} + \lambda I_d \right)^{-1} \frac{X^\top Y}{n}, \quad \forall \lambda > 0$$

## Ridge predictor

$$f_{\hat{\beta}_\lambda}(x) = x^\top \hat{\beta}_\lambda = x^\top \left( \frac{X^\top X}{n} + \lambda I_d \right)^{-1} \frac{X^\top Y}{n}, \quad \forall x \in \mathbb{R}^d$$

## Ridge predictor evaluated at the design points (the $X_i$ s)

$$\begin{aligned} F_{\hat{\beta}_\lambda} &= X \hat{\beta}_\lambda = X \left( \frac{X^\top X}{n} + \lambda I_d \right)^{-1} \frac{X^\top Y}{n} \\ &= \left( \frac{XX^\top}{n} + \lambda I_n \right)^{-1} \frac{XX^\top}{n} Y \end{aligned}$$

Introduction

Kernel PCA

Kernel Ridge  
Regression

Ridge regression

Reproducing Kernel  
Hilbert SpaceKernelized optimization  
problem

## From linear to Nonparametric regression model

- ▶ **Linear regression Model:** For all  $1 \leq i \leq n$ ,

$$Y_i = \underbrace{X_i^\top \beta^\star}_{=f_{\beta^\star}(X_i)} + \epsilon_i = \langle X_i, \beta \rangle_{\mathbb{R}^d} + \epsilon_i$$

with  $\mathbb{E}_{\epsilon_i} [\epsilon_i] = 0$ , and  $\text{Var}_{\epsilon_i}(\epsilon_i) = \sigma^2 > 0$

- ▶ **Nonparametric regression model:** For all  $1 \leq i \leq n$ ,

$$Y_i = f^\star(X_i) + \epsilon_i,$$

with  $\mathbb{E}_{\epsilon_i} [\epsilon_i] = 0$ , and  $\text{Var}_{\epsilon_i}(\epsilon_i) = \sigma^2 > 0$

### Remark:

- ▶ No necessary Gaussian assumption!
- ▶  $f^\star$ : not necessarily linear
- ▶  $f^\star(x) = \mathbb{E}[Y \mid X = x]$ : regression function

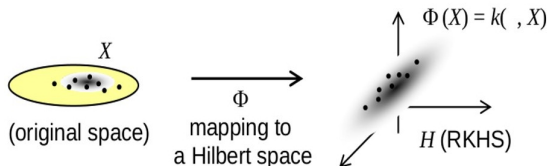
$$Y_i = f^*(X_i) + \epsilon_i,$$

## Estimating $f^*$ and Smoothness

- ▶ Estimating  $f^*$  requires defining candidate functions (Statistical Model)
- ▶ Classical assumptions are that:
  - ▶  $f^*$  bounded a.s. (because  $Y_i$  is so!)
  - ▶  $f^* \in L^2(P_X)$ , where  $P_X$ : Proba. distrib. of  $X$
  - ▶ ...
- ▶ Here, we assume that there exists a reproducing kernel  $k(\cdot, \cdot)$  and an RKHS  $\mathcal{H}$  (uniquely defined from the kernel) such that  $f^*$  is "not too far" from  $\mathcal{H}$

## Mapping data from $\mathcal{X}$ to $\mathcal{H}$

- ▶  $k(\cdot, \cdot)$ : reproducing kernel
- ▶  $x \in \mathbb{R}^d \mapsto k_x \in \mathcal{H}$ : canonical feature map from  $\mathbb{R}^d$  to  $\mathcal{H}$
- ▶  $k_x = \phi(x) \in \mathcal{H}$ : new “observation” in the RKHS



## Reminder

### Definition (Hilbert space)

Vector space endowed with a scalar product (pre-Hilbertian space), which is complete for the induced norm



## Classical Ridge

For every  $\lambda > 0$ , solve

$$(\|\beta\|_2^2 = \sum_{j=1}^d |\beta_j|^2)$$

$$\hat{\beta}_\lambda = \underset{\beta \in \mathbb{R}^d}{\text{Arg min}} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - f_\beta(X_i))^2 + \lambda \|\beta\|_2^2 \right\}$$

where  $f_\beta(x) = \langle x, \beta \rangle_{\mathbb{R}^d}$  and  $\|u\|_n^2 = \sum_{i=1}^n u_i^2 / n$   
(Important normalization by  $1/n!$ )

Introduction

Kernel PCA

Kernel Ridge  
Regression

Ridge regression

Reproducing Kernel  
Hilbert SpaceKernelized optimization  
problem

## Classical Ridge

For every  $\lambda > 0$ , solve

$$(\|\beta\|_2^2 = \sum_{j=1}^d |\beta_j|^2)$$

$$\hat{\beta}_\lambda = \mathop{\text{Arg min}}_{\beta \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - f_\beta(X_i))^2 + \lambda \|\beta\|_2^2 \right\}$$

where  $f_\beta(x) = \langle x, \beta \rangle_{\mathbb{R}^d}$  and  $\|u\|_n^2 = \sum_{i=1}^n u_i^2 / n$   
(Important normalization by  $1/n!$ )

## Kernel Ridge Regression (KRR)

For every  $\lambda > 0$ , solve

$$\hat{f}_\lambda = \mathop{\text{Arg min}}_{f \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\}$$

## Remark:

Optimization over an infinite dimensional space  $\mathcal{H}$

$$\hat{f}_\lambda = \operatorname{Arg} \min_{f \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\}$$

## Theorem (Representer theorem)

$\Psi : \mathbb{R}^n \times \mathbb{R}_+ \rightarrow \mathbb{R}$ , non-decreasing w.r.t. its  $n+1$ th argument

$$\operatorname{Arg} \min_{g \in \mathcal{H}} \{ \Psi [g(x_1), \dots, g(x_n), \|g\|_{\mathcal{H}}] \}$$

Any solution  $\hat{g}$  to the above optimization problem can be written as

$$\hat{g}(x) = \sum_{i=1}^n \hat{\alpha}_i k(x_i, x), \quad \forall x \in \mathcal{X}$$

where  $\hat{\alpha}_i \in \mathbb{R}$ , for all  $1 \leq i \leq n$

# Representer theorem (Proof)

►  $S = \text{Vect}(\{k_{x_1}, \dots, k_{x_n}\}) \subset \mathcal{H}$

► For all  $g \in \mathcal{H}$ ,

$$g = p_S^\perp(g) + \overbrace{(g - p_S^\perp(g))}^{\in S^\perp}$$

Then,  $\|g\|_{\mathcal{H}} \geq \|p_S^\perp(g)\|_{\mathcal{H}}$  (orthog. proj.)

► Besides,

$$g(x_i) = \langle g, k_{x_i} \rangle_{\mathcal{H}} \quad (\text{reproducing property})$$

$$= \langle p_S^\perp(g), k_{x_i} \rangle_{\mathcal{H}} \quad (\text{since } k_{x_i} \in S)$$

Then  $(\Psi \text{ non-decreasing})$

$$\begin{aligned} & \Psi[g(x_1), \dots, g(x_n), \|g\|_{\mathcal{H}}] \\ &= \Psi[\langle g, k_{x_1} \rangle_{\mathcal{H}}, \dots, \langle g, k_{x_n} \rangle_{\mathcal{H}}, \|g\|_{\mathcal{H}}] \\ &= \Psi[\langle p_S^\perp(g), k_{x_1} \rangle_{\mathcal{H}}, \dots, \langle p_S^\perp(g), k_{x_n} \rangle_{\mathcal{H}}, \|g\|_{\mathcal{H}}] \\ &\geq \Psi[\langle p_S^\perp(g), k_{x_1} \rangle_{\mathcal{H}}, \dots, \langle p_S^\perp(g), k_{x_n} \rangle_{\mathcal{H}}, \|p_S^\perp(g)\|_{\mathcal{H}}] \end{aligned}$$

Hence, any minimizer belongs to  $S$



## Conclusion for the KRR estimator

$$\hat{f}_\lambda = \sum_{i=1}^n \hat{\alpha}_i k_{x_i} \in \mathcal{H}, \quad \forall \lambda > 0$$

where  $\{\hat{\alpha}_i\}_i \subset \mathbb{R}^n$  are to be calculated

## Estimateur evaluated at the design points

$$\hat{F}_\lambda = (\hat{f}_\lambda(x_1), \dots, \hat{f}_\lambda(x_n))^T = K \hat{\alpha}$$

Then

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left( Y_i - \hat{f}_\lambda(X_i) \right)^2 + \lambda \|f\|_{\mathcal{H}}^2 &= \left\| Y - \hat{F}_\lambda \right\|_n^2 + \lambda \|f\|_{\mathcal{H}}^2 \\ &= \|Y - K \hat{\alpha}\|_n^2 + \lambda \hat{\alpha}^T K \hat{\alpha} \\ &= \varphi(\hat{\alpha}) \end{aligned}$$

Introduction

Kernel PCA

Kernel Ridge  
Regression

Ridge regression

Reproducing Kernel  
Hilbert SpaceKernelized optimization  
problem

# Closed-form expression (Cont'd)

$$\begin{aligned}\blacktriangleright \partial\varphi(\hat{\alpha}) = 0 &\Leftrightarrow \frac{1}{n}K(Y - K\hat{\alpha}) = \lambda K\hat{\alpha} \\ &\Leftrightarrow (K + n\lambda I_n)^{-1} KY = K\hat{\alpha} = \hat{F}_\lambda\end{aligned}$$

Kernel Machines

Alain Celisse

Introduction

Kernel PCA

Kernel Ridge  
Regression

Ridge regression

Reproducing Kernel  
Hilbert Space

Kernelized optimization  
problem

# Closed-form expression (Cont'd)

- ▶  $\partial\varphi(\hat{\alpha}) = 0 \Leftrightarrow \frac{1}{n}K(Y - K\hat{\alpha}) = \lambda K\hat{\alpha}$   
 $\Leftrightarrow (K + n\lambda I_n)^{-1}KY = K\hat{\alpha} = \hat{F}_\lambda$
- ▶  $\hat{F}_\lambda = \left( \left\langle \hat{f}_\lambda, k_{x_1} \right\rangle_{\mathcal{H}}, \dots, \left\langle \hat{f}_\lambda, k_{x_n} \right\rangle_{\mathcal{H}} \right)^\top$  (reprod. prop.)

# Closed-form expression (Cont'd)

- ▶  $\partial\varphi(\hat{\alpha}) = 0 \Leftrightarrow \frac{1}{n}K(Y - K\hat{\alpha}) = \lambda K\hat{\alpha}$   
 $\Leftrightarrow (K + n\lambda I_n)^{-1} KY = K\hat{\alpha} = \hat{F}_\lambda$
- ▶  $\hat{F}_\lambda = \left( \left\langle \hat{f}_\lambda, k_{x_1} \right\rangle_{\mathcal{H}}, \dots, \left\langle \hat{f}_\lambda, k_{x_n} \right\rangle_{\mathcal{H}} \right)^\top$  (reprod. prop.)
- ▶

$$\begin{aligned}(K + n\lambda I_n)^{-1} KY &= K (K + n\lambda I_n)^{-1} Y \\ &= \sum_{j=1}^n K_{\cdot,j} \left[ (K + n\lambda I_n)^{-1} Y \right]_j\end{aligned}$$



# Closed-form expression (Cont'd)

- ▶  $\partial\varphi(\hat{\alpha}) = 0 \Leftrightarrow \frac{1}{n}K(Y - K\hat{\alpha}) = \lambda K\hat{\alpha}$   
 $\Leftrightarrow (K + n\lambda I_n)^{-1} KY = K\hat{\alpha} = \hat{F}_\lambda$
- ▶  $\hat{F}_\lambda = \left( \langle \hat{f}_\lambda, k_{x_1} \rangle_{\mathcal{H}}, \dots, \langle \hat{f}_\lambda, k_{x_n} \rangle_{\mathcal{H}} \right)^\top$  (reprod. prop.)
- ▶

$$\begin{aligned}(K + n\lambda I_n)^{-1} KY &= K(K + n\lambda I_n)^{-1} Y \\ &= \sum_{j=1}^n K_{\cdot j} \left[ (K + n\lambda I_n)^{-1} Y \right]_j\end{aligned}$$

Since  $K_{\cdot j} = \left( \langle k_{x_1}, k_{x_j} \rangle_{\mathcal{H}}, \dots, \langle k_{x_n}, k_{x_j} \rangle_{\mathcal{H}} \right)^\top$ , we get

$$\hat{f}_\lambda = \sum_{j=1}^n k_{x_j} \underbrace{\left[ (K + n\lambda I_n)^{-1} Y \right]_j}_{=\hat{\alpha}_j} \in \mathcal{H}$$

No *a priori* unicity: one possible solution at this stage!

► Assume  $K\hat{\alpha} = K\hat{\beta} = \hat{F}_\lambda \in \mathbb{R}^n$

- ▶ Assume  $K\hat{\alpha} = K\hat{\beta} = \hat{F}_\lambda \in \mathbb{R}^n$
- ▶  $K(\hat{\alpha} - \hat{\beta}) = 0 \Leftrightarrow \hat{\alpha} - \hat{\beta} \in \text{Null}(K)$

- ▶ Assume  $K\hat{\alpha} = K\hat{\beta} = \hat{F}_\lambda \in \mathbb{R}^n$
- ▶  $K(\hat{\alpha} - \hat{\beta}) = 0 \Leftrightarrow \hat{\alpha} - \hat{\beta} \in \text{Null}(K)$
- ▶ Moreover

$$\begin{aligned} & \left\| \sum_{i=1}^n \hat{\alpha}_i k_{x_i} - \sum_{i=1}^n \hat{\beta}_i k_{x_i} \right\|_{\mathcal{H}} \\ &= (\hat{\alpha} - \hat{\beta})^\top K (\hat{\alpha} - \hat{\beta}) \\ &= 0 \end{aligned}$$

## Solution path

$$\lambda > 0 \mapsto \sum_{j=1}^n k_{x_j} \left[ (K + n\lambda I_n)^{-1} Y \right]_j$$

## Efficient computations

- ▶ Computing the whole path is possible by means of only one (careful) SVD
- ▶ When  $K = XX^\top$  (linear kernel  $\leftrightarrow$  classical Ridge regression), both  $XX^\top$  ( $n \times n$ ) and  $X^\top X$  ( $d \times d$ ) share the same non-zero singular values