

# Kernel

December 8, 2021

# Road map

- 1 Intuition and Motivation
- 2 Theoretical framework for kernels
  - Summary
  - Kernel on vectors
  - Reproducing kernel Hilbert space
- 3 Building kernels
  - Kernel algebra
  - Kernels on generic data
- 4 Key tools for ML
  - Kernel trick
  - The representer theorem

# Important References

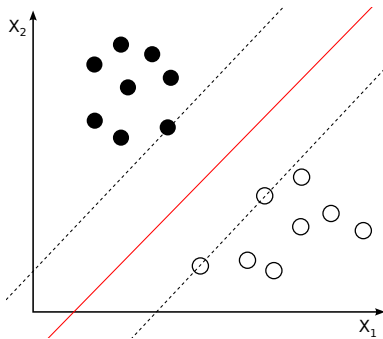
## Books

- Learning with kernel  
<https://mitpress.mit.edu/books/learning-kernels>
- Kernel methods for Pattern Analysis. J. Shawe-Taylor et al.
- Reproducing kernel Hilbert spaces in Probability and Statistics. Berlinet et al.

## Lecture notes

- J. Mairal and JP. Vert, MVA
- S. Canu, G. Gasso, B. Gauzere, INSA de Rouen

# Linear SVM



- Inputs are vectors of  $\mathbb{R}^d$
- Linear SVM seeks linear classification function

# Limitations

- Data are **not always vectors**: (string, time series, graphs, images ...)
- The decision function **can not always be linear** (text categorization; email filtering; gene detection; protein image classification; handwriting recognition; prediction of loan defaulting)

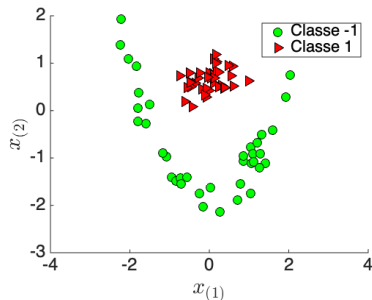
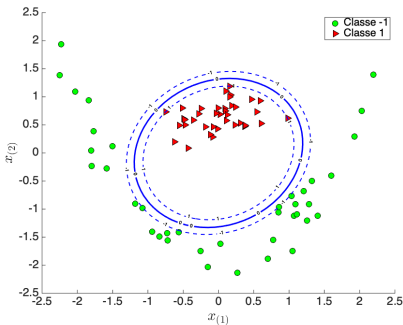


Figure: How do you classify these data?

# From linear to non-linear decision function

Data might be separable with a non-linear function



- **Non-linear embedding** of  $x = \begin{pmatrix} x_{(1)} \\ x_{(2)} \end{pmatrix}$

$$\mathbb{R}^2 \rightarrow \mathcal{H}$$

$$x \mapsto \Phi(x) = \begin{pmatrix} x_{(1)}^2 \\ x_{(2)}^2 \\ \sqrt{2}x_{(1)}x_{(2)} \end{pmatrix}$$

- Train a linear SVM with samples  $\{(\Phi(x_i), y_i)\}$

Resulting SVM model

$$f(x) = b + \sum_{i \in SV} \alpha_i y_i \Phi(x_i)^\top \Phi(x)$$

# Non-linear decision function

## Decision function

$$f(\mathbf{x}) = b + \sum_{i \in SV} \alpha_i y_i \underbrace{\Phi(\mathbf{x}_i)^\top \Phi(\mathbf{x})}_{\text{Kernel } k(\mathbf{x}_i, \mathbf{x})}$$

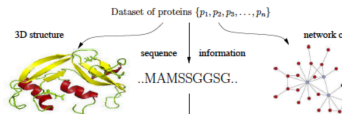
## Kernel function: the trick

- No explicit knowledge of  $\Phi(\mathbf{x})$
- We only need to define a function  $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

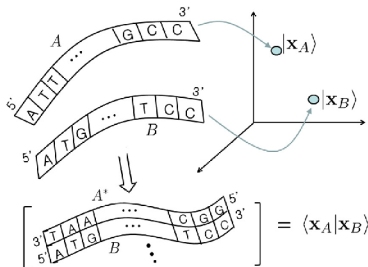
Problem linearly non-separable in the original space  $\mathcal{X}$  but linear separable in the space  $\mathcal{H}$  induced by the kernel  $k$

Also...

How do we classify dataset composed of proteins?



Use the kernel trick:  $f(\text{protein}) = b + \sum_{i \in SV} \alpha_i y_i k(\text{protein}_i, \text{protein})$





# Remark

## Intuition

By simply modifying the dot product, the algorithm works in another space

## Kernel Trick

- 1 Linear SVM relies on inner product between the SV and the sample to predict
- 2  $\rightarrow$  Replaces the inner product between the sample in the ambient space by a **kernel**  $k(\cdot, \cdot)$
- 3  $\rightarrow$  Leads to a non-linear version of the SVM

# Motivation of Kernel methods

- From
  - linear techniques
  - operating on vector spaces
- to
  - non linear prediction models
  - operating on various, structured, high-dimensional data
- Using a:  
well known mathematical framework
- leading to:  
efficient and powerful algorithm and tools

→ Kernel method

# One slide summary

## Motivation

- develop generic algorithms for analyzing and learning from data
- ... without making any assumptions on the type of the data

## The approach

- Introducing framework based on similarities between pair of examples
- By appropriately defining the notion of similarity, we define a framework named as **learning in RKHS**.

# Prerequisites

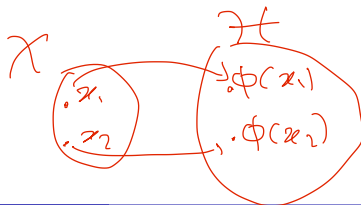
decision function  $f(x) = \sum_i \alpha_i \phi(x_i) \phi(x) = \sum_{i=1}^N \alpha_i k(x_i, x)$

## Some definitions and notations

- $\mathcal{X}$ : non empty input space ( $\mathbb{R}^N$ , graphs, objects, ...)
- $x \in \mathcal{X}$ ,
- $\mathcal{H}$ : feature space endowed with a dot product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$

→ •  $\Phi : \mathcal{X} \rightarrow \mathcal{H}$ : embedding function from  $\mathcal{X}$  to  $\mathcal{H}$

→  $k(\cdot, \cdot)$



# Kernel as a similarity

## Kernel

A kernel  $k$  is a function  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  such that in some way  $k(x, z)$  captures the "similarity" between  $x$  and  $z$ .

## Ideas

- the matrix  $K$  where each  $K_{i,j}$  represents all pairwise similarities
- For  $n$  data,  $K$  is of size  $n \times n$  matrix of  $\mathbb{R}$  regardless of the data type.
- Modularity between the data representation (through  $K$  and the choice of the algorithm (that uses  $K$ ))

# Positive Definite Kernels (1)

## Positive definite kernel

A kernel  $k(x, z)$  on  $\mathcal{X} \times \mathcal{X}$  is said to be positive definite

- if it is symmetric:  $k(x, z) = k(z, x)$
- and if for any finite positive integer  $n$ :

$$\forall \{\alpha_i\}_{i=1,n} \in \mathbb{R}, \forall \{x_i\}_{i=1,n} \in \mathcal{X}, \quad \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(x_i, x_j) \geq 0$$

it is strictly positive definite if for  $\alpha_i \neq 0$

$$\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(x_i, x_j) > 0$$

let  $K_{ij} = k(x_i, x_j)$ ,  $K \in \mathbb{R}^{N \times N}$ ,

equivalently, the gram matrices  $K$  are p.s.d.

$$\alpha^T K \alpha \geq 0$$

$$\alpha^T K \alpha > 0$$

## Positive Definite Kernels (2)

### Gram Matrix

Given a kernel  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , and samples  $\{x_1, \dots, x_n\}$ , the **Gram Matrix**  $K$  is a  $n \times n$  matrix with entries  $K_{i,j} := k(x_i, x_j)$

### Another way to characterize positive definite kernel

For any set of  $n \in \mathbb{N}$  samples  $\{x_1, \dots, x_n\}$  the associated Gram Matrix  $K \in \mathbb{R}^{n \times n}$  is positive definite, iff  $k$  is a **positive definite kernel** on  $\mathcal{X}$ .

### Kernel method

We define kernel methods as algorithms taking positive definite matrix as input.

# Linear Kernel

$$k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^\top \mathbf{z}$$

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

- $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$
- symmetric:  $\mathbf{x}^\top \mathbf{z} = \mathbf{z}^\top \mathbf{x}$
- positive:

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) &= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \mathbf{x}_i^\top \mathbf{x}_j \\ &= \left( \sum_{i=1}^n \alpha_i \mathbf{x}_i \right)^\top \left( \sum_{j=1}^n \alpha_j \mathbf{x}_j \right) \\ &= \left\| \sum_{i=1}^n \alpha_i \mathbf{x}_i \right\|^2 \end{aligned}$$



# Product kernel

$$k(x, z) = g(x)g(z)$$

- $x, z \in \mathcal{X}$
- for some  $g : \mathcal{X} \rightarrow \mathbb{R}$
- symmetric: by construction
- positive:

*x, z are graphs.  
g(x) = nb of  
connected  
components in x*

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(x_i, x_j) &= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j g(x_i) g(x_j) \\ &= \left( \sum_{i=1}^n \alpha_i g(x_i) \right) \left( \sum_{j=1}^n \alpha_j g(x_j) \right) \\ &= \left( \sum_{i=1}^n \alpha_i g(x_i) \right)^2 \end{aligned}$$

# Finite kernel

let  $\phi_j, j = 1, p$  be a finite dictionary of functions from  $\mathcal{X}$  to  $\mathbb{R}$  (polynomials, wavelets...)

the feature map and linear kernel

feature map:  $\Phi : \mathcal{X} \rightarrow \mathbb{R}^p$   
 $x \mapsto \Phi = (\phi_1(x), \dots, \phi_p(x))$

kernel in the feature space is a positive definite kernel:

$$k(x, z) = (\phi_1(x), \dots, \phi_p(x))^T (\phi_1(z), \dots, \phi_p(z))$$

$$\sum_{i,j} \alpha_i \alpha_j k(x_i, x_j) \geq 0$$

# The quadratic kernel

For  $x, z \in \mathbb{R}^d$ , we define the feature map as

Feature map

$$\begin{aligned}\Phi : \mathbb{R}^d &\rightarrow \mathbb{R}^{p=\frac{d(d+1)}{2}} \\ x &\mapsto \Phi = (x_1^2, \dots, x_j^2, \dots, x_d^2, \dots, \sqrt{2}x_i x_j, \dots)\end{aligned}$$

Here the  $x_j$  represent the variables of  $x \in \mathbb{R}^d$

The kernel  $k(x, z)$  can be shown to be

$$k(x, z) = (x^\top z)^2$$

# Kernels as inner products

## Theorem (Aronszajn, 1950)

$k$  is a positive definite kernel on  $\mathcal{X}$  if and only if, there exists a Hilbert space  $\mathcal{H}$  and a mapping

$$\Phi : \mathcal{X} \mapsto \mathcal{H}$$

such that for any  $x$  and  $z$  in  $\mathcal{X}$  :

$$k(x, z) = \langle \Phi(x), \Phi(z) \rangle_{\mathcal{H}}$$

# Reproducing Kernel Map

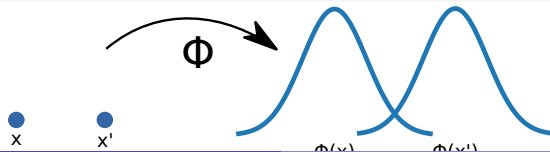
## Preliminaries

- We need a particular Hilbert space as defined by Aronszjan's theorem
- $\Phi$  maps each point  $x \in \mathcal{X}$  to a function in  $\mathcal{H}$
- $\mathcal{H}$  is a space of function from  $\mathcal{X} \rightarrow \mathbb{R}$
- Build  $\Phi$  from the psd kernel  $k(\cdot, \cdot)$

$$\begin{aligned}\Phi : \mathcal{X} &\rightarrow \mathbb{R}^{\mathcal{X}} \\ x &\mapsto k(\cdot, x)\end{aligned}$$

gaussian kernel  
or RBF kernel  
 $k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$

- Example :  $\mathcal{X} = \mathbb{R}$  and  $k(\cdot, x) = y \mapsto e^{-\frac{(x-y)^2}{2\sigma^2}}$



# RKHS Definition

## Definition

Let  $\mathcal{X}$  be a set and  $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$  be a class of function forming a real Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ . The function  $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$  is called a reproducing kernel (r.k.) of  $\mathcal{H}$  if

- 1  $\mathcal{H}$  contains all functions of the form

$$\forall x \in \mathcal{X}, \quad k(x, \cdot) : z \mapsto k(x, z)$$

*functions in  $\mathcal{H}$  can be written as  $f(x) = \sum \alpha_i k(x_i, x)$*

*$\|f\|_{\mathcal{H}}^2 = \langle \sum \alpha_i k(x_i, \cdot), \sum \alpha_j k(x_j, \cdot) \rangle_{\mathcal{H}} = \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j)$*

- 2 For every  $x \in \mathcal{X}$  and  $f \in \mathcal{H}$ , the reproducing property holds :

$$f(x) = \langle f, k(x, \cdot) \rangle_{\mathcal{H}}$$

*$\Rightarrow \|k(x, \cdot)\|_{\mathcal{H}}$   
 $= \langle k(x, \cdot), k(x, \cdot) \rangle_{\mathcal{H}}$   
 $= k(x, x)$*

if a r.k exists, then  $\mathcal{H}$  is called a reproducing kernel Hilbert space (RKHS).

# RKHS and Machine Learning

Using RKHS as a hypothesis space for machine learning problem leads to a simple recipe for non-linear models

- 1 maps data  $x \in \mathcal{X}$  to a high-dimensional r.k. Hilbert space  $\mathcal{H}$  through the mapping  $\Phi : \mathcal{X} \mapsto \mathcal{H}$  with  $\Phi(x) = k(x, \cdot)$
- 2 in  $\mathcal{H}$ , consider linear model with  $f(x) = \langle f, k(x, \cdot) \rangle_{\mathcal{H}}$
- 3 use this linear model in your learning framework. For supervised learning, we would have

$$\min_{f \in \mathcal{H}} \sum_{i=1}^n L(y_i, f(x_i)) + \lambda \|f\|_{\mathcal{H}}^2$$

# Uniqueness of r.k and RKHS

## Theorem

- if  $\mathcal{H}$  is a RKHS then it has an unique r.k
- A function  $k(\cdot, \cdot)$  can be the r.k of at most one RKHS



# RKHS, reproducing kernel and positive definite kernel

## Theorem

A function  $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$  is positive definite if and only if it is a reproducing kernel.

## Proof : an r.k is pd

- a r.k is symmetric
- a r.k leads to a pd matrix for any  $n$  subset of  $\{x_i\}_{i=1}^n$

# A positive definite kernel is a reproducing kernel I

We aim to encode the image of  $\Phi$  into a vector space

## The vector space

Linear combinations of  $k(\cdot, x)$ :

$$f(\cdot) = \sum_{i=1}^m \alpha_i k(\cdot, x_i) \quad \text{with any } x_1, \dots, x_m \in \mathcal{X}$$

## The dot product

- given  $g(\cdot) = \sum_{j=1}^{m'} \beta_j k(\cdot, x'_j)$  with  $x'_1, \dots, x'_{m'} \in \mathcal{X}$  we define the dot product as

$$\langle f, g \rangle := \sum_{i=1}^m \sum_{j=1}^{m'} \alpha_i \beta_j k(x_i, x'_j)$$

## A positive definite kernel is a reproducing kernel II

We have:

$$f(x'_j) = \sum \alpha_i k(x'_j, x_i) = \sum \alpha_i k(x_i, x'_j)$$

Hence:

$$\langle f, g \rangle = \sum \beta_j f(x'_j)$$

Note that it does not depend on the expansion of  $f$ . Similarly we have

$$\langle f, g \rangle = \sum \alpha_i g(x_i)$$

It is easy to show that our dot product is:

- bilinear
- symmetric
- positive definite

and thus constitutes a valid dot product on the vector space  $\mathbb{R}^{\mathcal{X}}$ .

# A positive definite kernel is a reproducing kernel III

## Reproducing kernel

$$\langle k(\cdot, \mathbf{x}), f \rangle = \sum \alpha_i k(\mathbf{x}_i, \mathbf{x}) = f(\mathbf{x})$$

Considering  $f(\cdot) = k(\cdot, \mathbf{x}')$ , we have:

$$\langle k(\cdot, \mathbf{x}), k(\cdot, \mathbf{x}') \rangle = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle = k(\mathbf{x}, \mathbf{x}')$$

$k$  is thus the **reproducing kernel** of  $\mathcal{H}$  and corresponds to a dot product in the vector space of functions  $\mathcal{H}$ .

# The RKHS associated to a p.d kernel $k$ .

## Hilbert Space

- $\mathcal{H}$  is Hilbert Space obtained from
  - pre-Hilbert space of functions defined as above endowed with a inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  inducing a norm  $\|f\| := \sqrt{\langle f, f \rangle}$
  - and the space completed with all limits of Cauchy sequences

$\mathcal{H}$  is called a **reproducing kernel Hilbert space (RKHS)** associated to kernel  $k$

# Smoothness functional

## A simple inequality

- By Cauchy-Schwarz, we have for any function  $f \in \mathcal{H}$  and two points  $x, z \in \mathcal{X}$  :

$$\begin{aligned} |f(x) - f(z)| &= \langle f, k(x, \cdot) - k(z, \cdot) \rangle_{\mathcal{H}} \\ &\leq \|f\|_{\mathcal{H}} \times \|k(x, \cdot) - k(z, \cdot)\|_{\mathcal{H}} \\ &= \|f\|_{\mathcal{H}} \times d_k(x, z) \end{aligned}$$

- the norm of a function  $f$  in the RKHS controls the variation of  $f$  over  $\mathcal{X}$  with respect to the geometry defined by the kernel.

## Take-home message

Small norm  $\Rightarrow$  Small variations

# RKHS

## Positive kernel $\Leftrightarrow$ RKHS

- it defines the inner product
- it defines **regularity** (smoothness) of  $\mathcal{H}$
- there exists a "mapping" function  $\Phi : \mathcal{X} \rightarrow \mathcal{H}$  such that  $\forall x, z \in \mathcal{X}$  the inner product  $\langle \Phi(x), \Phi(z) \rangle_{\mathcal{H}} = k(x, z)$  and thus it defines the function space

# Let's summarize

## From kernel to feature space

Given a valid kernel  $k$ , we can associate a RKHS  $\mathcal{H}$  which corresponds to the feature space of  $k$ .

## From feature space to kernel

Now consider that you have  $\Phi : \mathcal{X} \rightarrow \mathcal{H}$  a mapping function.  
A positive kernel  $k$  is defined by:

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}}$$



# Kernel algebra

## Convex cone:

The set of kernels forms a convex cone, closed under pointwise convergence.

- **Linear combination:**

- if  $k_1$  and  $k_2$  are kernels,  $a_1, a_2 \geq 0$ , then  $a_1 k_1 + a_2 k_2$  is a kernel
- if  $k_1, k_2, \dots$  are kernels, and  $k(x, x') := \lim_{n \rightarrow \infty} k_n(x, x')$  exists for all  $x, x'$ , then  $k$  is a kernel

- **Product kernel:**

if  $k_1$  and  $k_2$  are kernels, then  $k_1 k_2(x, x') := k_1(x, x') k_2(x, x')$  is a kernel.

## Proofs

- by linearity:

$$\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j (a_1 k_1(i, j) + k_2(i, j)) = a_1 \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k_1(i, j) + \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k_2(i, j)$$

- assuming  $\exists \psi_\ell$  s.t.  $k_1(s, t) = \sum_{\ell} \psi_\ell(s) \psi_\ell(t)$

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k_1(x_i, x_j) k_2(x_i, x_j) &= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \left( \sum_{\ell} \psi_\ell(x_i) \psi_\ell(x_j) k_2(x_i, x_j) \right) \\ &= \sum_{\ell} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i \psi_\ell(x_i)) (\alpha_j \psi_\ell(x_j)) k_2(x_i, x_j) \end{aligned}$$

# Kernel engineering: building PDK

- for any polynomial with positive coef.  $\phi$  from  $\mathbb{R}$  to  $\mathbb{R}$

$$\phi(k(s, t))$$

- if  $\Psi$  is a function from  $\mathbb{R}^d$  to  $\mathbb{R}^d$

$$k(\Psi(s), \Psi(t))$$

Example : the Gaussian kernel is a PDK

$$\begin{aligned}\exp(-\|s - t\|^2) &= \exp(-\|s\|^2 - \|t\|^2 - 2s^\top t) \\ &= \exp(-\|s\|^2) \exp(-\|t\|^2) \exp(2s^\top t)\end{aligned}$$

- $s^\top t$  is a PDK and function  $\exp$  as the limit of positive series expansion, so  $\exp(2s^\top t)$  is a PDK
- $\exp(-\|s\|^2) \exp(-\|t\|^2)$  is a PDK as a product kernel
- the product of two PDK is a PDK

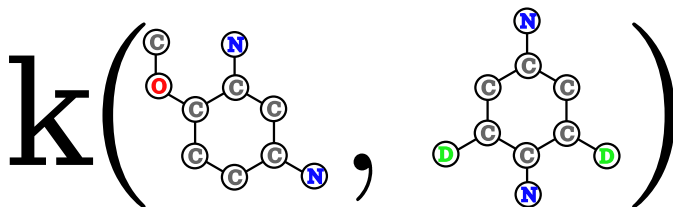
# Positive definite kernels: some common examples

Type	Name	$k(x, z)$
radial	Gaussian	$\exp\left(-\frac{\ x-z\ ^2}{2\sigma^2}\right)$
radial	Laplacian	$\exp(-\ x-z\ /\sigma)$
non stat.	$\chi^2$	$\exp(-r/\sigma), r = \sum_k \frac{(x_k - z_k)^2}{x_k + z_k}$
projectif	polynomial	$(x^\top z + \sigma)^p$
projectif	cosinus	$x^\top z / \ x\  \ z\ $
projectif	correlation	$\exp\left(\frac{x^\top z}{\ x\  \ z\ } - \sigma\right)$

- The kernel may involve hyper-parameter(s) to tune (polynom order  $p$ , bandwidth  $\sigma$ )
- Their value has to be set by cross-validation

# Kernels on structures

- $\mathcal{X}$  may not be a vector space.
- we can define kernels on any kind of data :
  - Strings
  - Time series
  - Graphs
  - Images
  - ...



# RKHS, kernel and machine learning : Kernel trick

## Proposition

Any algorithm to process finite-dimensional vectors that can be expressed only in terms of pairwise inner-products can be applied to potentially infinite-dimensional vectors in the space of a p.d kernel by replacing each inner product evaluation by a kernel evaluation.

## Applications

- replace inner product by kernel evaluation
- replace any inner-product induced norm by kernel evaluation

# RKHS and machine learning : representer theorem

## The representer theorem

- let  $\mathcal{H}$  be a RKHS with  $k$  as associated kernel and  $\mathcal{S} = \{x_1, \dots, x_n\}$  be a finite set of points in  $\mathcal{X}$
- let  $\Psi : \mathbb{R}^{n+1} \mapsto \mathbb{R}$  be a function of  $n + 1$  variables, strictly increasing with respect to the last variable.
- Then, any solution to the optimization problem :

$$\min_{f \in \mathcal{H}} \Psi(f(x_1), \dots, f(x_n), \|f\|_{\mathcal{H}})$$

admits a representation of the form:

$$f(x) = \sum_{i=1}^n \alpha_i K(x_i, x)$$

# Using the Representer theorem

## In practice

- When the representer theorem holds for our learning problem, then we can look for solution  $f$  of the form :

$$f(\cdot) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \cdot)$$

- Pointwise evaluation : for any  $j = 1, \dots, n$ , we have

$$f(\mathbf{x}_j) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x}_j) = [\mathbf{K}\alpha]_j$$

- The norm

$$\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \sum_{i,j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) = \alpha^\top \mathbf{K} \alpha$$



# Using the Representer theorem

## In practice

- A problem of the form

$$\min_{f \in \mathcal{H}} \Psi(f(x_1), \dots, f(x_n), \|f\|_{\mathcal{H}})$$

is equivalent to the n-dimensional optimization problem

$$\min_{\alpha \in \mathbb{R}^n} \Psi([K\alpha]_1, \dots, [K\alpha]_n, \alpha^\top K \alpha)$$

- The resulting problem can be usually solved analytically or by numerical methods

# Conclusion

## What we seen

- Kernels corresponds to scalar product in some Hilbert space:
  - value corresponds to high dimensional scalar product,
  - on non linear embedding
  - without explicit representations of  $\Phi$
- Can be defined on any kind of data
- Kernels define the functional space (and its smoothness) in which we are looking for the solution

## Key tools for ML

- Representer theorem
- Kernel trick