

MLSS 2019 London

Kernels

Part II: Reproducing Kernel Hilbert Spaces (RKHS)

Lorenzo Rosasco

MaLGa- Machine learning Genova Center

Università di Genova

MIT

IIT

Linear models \mapsto **features** \mapsto **kernels**

$$\mathbf{x}^\top \mathbf{x}' \mapsto \Phi(\mathbf{x})^\top \Phi(\mathbf{x}') \mapsto \mathbf{k}(\mathbf{x}, \mathbf{x}')$$

“The kernel trick”

Outline

PD kernels

RKHS

A bit more than a trick...

- ▶ Feature maps
- ▶ PD kernels
- ▶ RKHS

Can we start from a kernel?

Can we start from a kernel?

What is a kernel?

$$k(\mathbf{x}, \mathbf{x}')$$

Can we start from a kernel?

What is a kernel?

$$k(\mathbf{x}, \mathbf{x}')$$

Two words:

Positive definite

Positive definite kernels

$k : X \times X \rightarrow \mathbb{R}$ is **positive definite** if:

1. for all $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$, the $N \times N$ matrix \hat{K} with entries

$$\hat{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

is positive semidefinite (non negative eigenvalues)

Positive definite kernels

$k : X \times X \rightarrow \mathbb{R}$ is **positive definite** if:

1. for all $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$, the $N \times N$ matrix \hat{K} with entries

$$\hat{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

is positive semidefinite (non negative eigenvalues)

2. Equivalently for all $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$, and $\mathbf{a} \in \mathbb{R}^N$

$$\mathbf{a}^\top \hat{K} \mathbf{a} \geq 0, \quad \forall \mathbf{a} \in \mathbb{R}^N.$$

Positive definite kernels

$k : X \times X \rightarrow \mathbb{R}$ is **positive definite** if:

1. for all $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$, the $N \times N$ matrix \hat{K} with entries

$$\hat{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

is positive semidefinite (non negative eigenvalues)

2. Equivalently for all $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$, and $\mathbf{a} \in \mathbb{R}^N$

$$\mathbf{a}^\top \hat{K} \mathbf{a} \geq 0, \quad \forall \mathbf{a} \in \mathbb{R}^N.$$

3. Equivalently, for all $a_1, \dots, a_N \in \mathbb{R}$, $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$,

$$\sum_{i,j=1}^N k(\mathbf{x}_i, \mathbf{x}_j) a_i a_j \geq 0.$$

Positive definite kernels

$k : X \times X \rightarrow \mathbb{R}$ is **positive definite** if:

1. for all $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$, the $N \times N$ matrix \hat{K} with entries

$$\hat{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

is positive semidefinite (non negative eigenvalues)

2. Equivalently for all $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$, and $\mathbf{a} \in \mathbb{R}^N$

$$\mathbf{a}^\top \hat{K} \mathbf{a} \geq 0, \quad \forall \mathbf{a} \in \mathbb{R}^N.$$

3. Equivalently, for all $a_1, \dots, a_N \in \mathbb{R}$, $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$,

$$\sum_{i,j=1}^N k(\mathbf{x}_i, \mathbf{x}_j) a_i a_j \geq 0.$$

Symmetry is also always assumed.

Inner product kernels are PD

For $\Phi : \mathbf{X} \rightarrow \mathbb{R}^p$, $p \leq \infty$ let

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

Inner product kernels are PD

For $\Phi : X \rightarrow \mathbb{R}^p$, $p \leq \infty$ let

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

► Then

$$\sum_{i,j=1}^N k(\mathbf{x}_i, \mathbf{x}_j) \mathbf{a}_i \mathbf{a}_j = \sum_{i,j=1}^N \Phi(\mathbf{x}_i)^\top \Phi(\mathbf{x}_j) \mathbf{a}_i \mathbf{a}_j = \left\| \sum_{i=1}^N \Phi(\mathbf{x}_i) \mathbf{a}_i \right\|^2 \geq 0.$$

Inner product kernels are PD

For $\Phi : \mathbf{X} \rightarrow \mathbb{R}^p$, $p \leq \infty$ let

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

► Then

$$\sum_{i,j=1}^N k(\mathbf{x}_i, \mathbf{x}_j) \mathbf{a}_i \mathbf{a}_j = \sum_{i,j=1}^N \Phi(\mathbf{x}_i)^\top \Phi(\mathbf{x}_j) \mathbf{a}_i \mathbf{a}_j = \left\| \sum_{i=1}^N \Phi(\mathbf{x}_i) \mathbf{a}_i \right\|^2 \geq 0.$$

► Clearly k is symmetric.

Kernel properties

Let $K_1 : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, $K_2 : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, $K_3 : \mathbb{R}^t \times \mathbb{R}^t \rightarrow \mathbb{R}$ be kernels, $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$, $\mathbf{z}, \mathbf{z}' \in \mathbb{R}^t$ and $\alpha, \beta > 0$ then the following are also kernels

1. $\alpha K_1(\mathbf{x}, \mathbf{x}') + \beta K_2(\mathbf{x}, \mathbf{x}')$
2. $K_1(\mathbf{x}, \mathbf{x}') K_2(\mathbf{x}, \mathbf{x}')$

Kernel properties

Let $K_1 : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, $K_2 : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, $K_3 : \mathbb{R}^t \times \mathbb{R}^t \rightarrow \mathbb{R}$ be kernels, $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$, $\mathbf{z}, \mathbf{z}' \in \mathbb{R}^t$ and $\alpha, \beta > 0$ then the following are also kernels

1. $\alpha K_1(\mathbf{x}, \mathbf{x}') + \beta K_2(\mathbf{x}, \mathbf{x}')$
2. $K_1(\mathbf{x}, \mathbf{x}')K_2(\mathbf{x}, \mathbf{x}')$
3. $p(K_1(\mathbf{x}, \mathbf{x}'))$ for any p a function whose polynomial expansion has only non-negative coefficients
4. $f(\mathbf{x})K_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$ for any $f : \mathbb{R}^d \rightarrow \mathbb{R}$
5. $\frac{K_1(\mathbf{x}, \mathbf{x}')}{\sqrt{K_1(\mathbf{x}, \mathbf{x})K_1(\mathbf{x}', \mathbf{x}')}}}$
6. $K_3(\psi(\mathbf{x}), \psi(\mathbf{x}'))$ for any $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^t$
7. $\alpha K_1(\mathbf{x}, \mathbf{x}') + \beta K_3(\mathbf{z}, \mathbf{z}')$
8. $K_1(\mathbf{x}, \mathbf{x}')K_3(\mathbf{z}, \mathbf{z}')$

There are many PD kernels

All the examples seen so far are based on inner products.

- ▶ linear $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$
- ▶ polynomial $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + 1)^s$
- ▶ Gaussian $k(\mathbf{x}, \mathbf{x}') = e^{-\|\mathbf{x} - \mathbf{x}'\|^2 \gamma}$

\mathbf{X} need not be \mathbb{R}^d , one can consider

- ▶ kernels on probability distributions
- ▶ kernels on strings
- ▶ kernels on functions
- ▶ kernels on groups
- ▶ kernels graphs
- ▶ ...

Convolution kernels

- ▶ Sets X, X_1, \dots, X_D and $X^D = X_1 \times X_2 \times \dots \times X_D$.
- ▶ Pd kernels $k_1 : X_1 \times X_1 \rightarrow \mathbb{R}, \dots, k_D : X_D \times X_D \rightarrow \mathbb{R}$
- ▶ Relation $R : X^D \times X \rightarrow \{0, 1\}$ and

$$R^{-1}(\mathbf{x}) = \{\mathbf{x}^D \in X^D \mid R(\mathbf{x}^D, \mathbf{x}) = 1\}$$

Interpretation

- ▶ $\mathbf{x}^D = (x_1, \dots, x_D) \in X^D$ are the possible **parts** of \mathbf{x} .
- ▶ $R(\mathbf{x}^D, \mathbf{x}) = 1$ if \mathbf{x}^D are the parts of \mathbf{x}

Example $X = X_1 = X_2$ strings over a finite alphabet.

$$R((x_1, x_2), x) = 1 \quad \text{iff} \quad x_1 \circ x_2 = x.$$

Convolution kernels

- ▶ Sets X, X_1, \dots, X_D and $X^D = X_1 \times X_2 \times \dots \times X_D$.
- ▶ Pd kernels $k_1 : X_1 \times X_1 \rightarrow \mathbb{R}, \dots, k_D : X_D \times X_D \rightarrow \mathbb{R}$
- ▶ Relation $R : X^D \times X \rightarrow \{0, 1\}$ and

$$R^{-1}(\mathbf{x}) = \{\mathbf{x}^D \in X^D \mid R(\mathbf{x}^D, \mathbf{x}) = 1\}$$

Interpretation

- ▶ $\mathbf{x}^D = (x_1, \dots, x_D) \in X^D$ are the possible **parts** of \mathbf{x} .
- ▶ $R(\mathbf{x}^D, \mathbf{x}) = 1$ if \mathbf{x}^D are the parts of \mathbf{x}

Convolution kernel

$$k(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{x}^D \in R^{-1}(\mathbf{x}), \mathbf{z}^D \in R^{-1}(\mathbf{z})} \prod_{j=1}^D k_j(x_j, z_j)$$

A bit more than a trick...

- ▶ Feature maps
- ▶ PD kernels
- ▶ RKHS

Outline

PD kernels

RKHS

Function spaces

$$\Phi \mapsto \mathbf{f}(\mathbf{x}) = \mathbf{w}^\top \Phi(\mathbf{x})$$

$$\mathbf{k} \mapsto \mathbf{f}(\mathbf{x}) = ?$$

Function spaces

$$\Phi \mapsto \mathcal{H}_{\Phi} = \{\mathbf{f} : \mathbf{X} \rightarrow \mathbb{R} \mid \exists! \mathbf{w} \in \mathcal{F} \text{ s.t. } \mathbf{f}(\mathbf{x}) = \mathbf{w}^{\top} \Phi(\mathbf{x}), \forall \mathbf{x} \in \mathcal{X}\}$$

$$\mathbf{k} \mapsto \mathcal{H}_{\mathbf{k}} = ?$$

Hilbert spaces

Hilbert space H

- ▶ Linear space (closed under sum/multiplication with reals)

$$\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}, \mathbf{a}, \mathbf{b} \in \mathbb{R} \quad \Rightarrow \quad \mathbf{a}\mathbf{h}_1 + \mathbf{b}\mathbf{h}_2 \in H$$

Hilbert spaces

Hilbert space H

- ▶ Linear space (closed under sum/multiplication with reals)

$$\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}, \mathbf{a}, \mathbf{b} \in \mathbb{R} \quad \Rightarrow \quad \mathbf{a}\mathbf{h}_1 + \mathbf{b}\mathbf{h}_2 \in H$$

- ▶ With an inner product (positive and symmetric bilinear form)

$$\mathbf{h}_1^\top \mathbf{h}_2$$

Hilbert spaces

Hilbert space H

- ▶ Linear space (closed under sum/multiplication with reals)

$$\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}, \mathbf{a}, \mathbf{b} \in \mathbb{R} \quad \Rightarrow \quad \mathbf{a}\mathbf{h}_1 + \mathbf{b}\mathbf{h}_2 \in H$$

- ▶ With an inner product (positive and symmetric bilinear form)

$$\mathbf{h}_1^\top \mathbf{h}_2$$

- ▶ complete (Cauchy sequences converge)¹

$$(\mathbf{h}_j)_j \quad \mathbf{i}, \mathbf{j} > \mathbf{N}, \quad \|\mathbf{h}_j - \mathbf{h}_i\| \leq \epsilon \quad \Rightarrow \quad \lim_{j \rightarrow \infty} \mathbf{h}_j = \mathbf{h} \in H$$

Function spaces

$$\Phi \mapsto \mathcal{H}_\Phi = \{\mathbf{f} : \mathbf{X} \rightarrow \mathbb{R} \mid \exists! \mathbf{w} \in \mathcal{F} \text{ s.t. } \mathbf{f}(\mathbf{x}) = \mathbf{w}^\top \Phi(\mathbf{x}), \forall \mathbf{x} \in \mathcal{X}\}$$

It's a Hilbert space since \mathcal{F} is!

$$\mathbf{k} \mapsto \mathcal{H}_{\mathbf{k}} = ?$$

From PD kernels to function spaces

For a set X and PD kernel k :

- define the linear the space of functions²

$$f(\mathbf{x}) = \sum_{i=1}^N k(\mathbf{x}, \mathbf{x}_i) a_i$$

for any $a_1, \dots, a_N \in \mathbb{R}$, $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$ and any $N \in \mathbb{N}$.

From PD kernels to function spaces

For a set X and PD kernel k :

- define the linear the space of functions²

$$f(\mathbf{x}) = \sum_{i=1}^N k(\mathbf{x}, \mathbf{x}_i) a_i$$

for any $a_1, \dots, a_N \in \mathbb{R}$, $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$ and any $N \in \mathbb{N}$.

- define the inner product

$$\langle f, f' \rangle = \sum_{i=1}^N \sum_{j=1}^{N'} k(\mathbf{x}_i, \mathbf{x}'_j) a_i a'_j.$$

From PD kernels to function spaces

For a set X and PD kernel k :

- ▶ define the linear the space of functions²

$$f(\mathbf{x}) = \sum_{i=1}^N k(\mathbf{x}, \mathbf{x}_i) a_i$$

for any $a_1, \dots, a_N \in \mathbb{R}$, $\mathbf{x}_1, \dots, \mathbf{x}_N \in X$ and any $N \in \mathbb{N}$.

- ▶ define the inner product

$$\langle f, f' \rangle = \sum_{i=1}^N \sum_{j=1}^{N'} k(\mathbf{x}_i, \mathbf{x}'_j) a_i a'_j.$$

\mathcal{H}_k is the completion to a Hilbert space.

A key result

Theorem

Given a PD kernel k there exists Φ s.t. $k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle_{\mathcal{H}_k}$ and

$$\mathcal{H}_\Phi \simeq \mathcal{H}_k$$

A key result

Theorem

Given a PD kernel k there exists Φ s.t. $k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle_{\mathcal{H}_k}$ and

$$\mathcal{H}_\Phi \simeq \mathcal{H}_k$$

Proof: Not so easy...

A key result

Theorem

Given a PD kernel k there exists Φ s.t. $k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle_{\mathcal{H}_k}$ and

$$\mathcal{H}_\Phi \simeq \mathcal{H}_k$$

Proof: Not so easy...

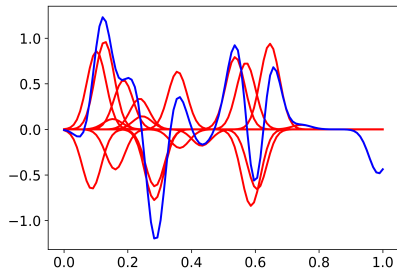
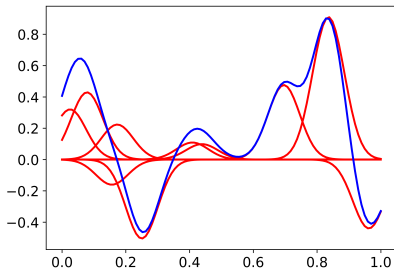
Roughly speaking

$$\mathbf{f}(\mathbf{x}) = \mathbf{w}^\top \Phi(\mathbf{x}) \quad \simeq \quad \mathbf{f}(\mathbf{x}) = \sum_{i=1}^N k(\mathbf{x}, \mathbf{x}_i) \mathbf{a}_i$$

An illustration

Functions defined by Gaussian kernels with large and small widths.

$$f(\mathbf{x}) = \sum_{i=1}^N k(\mathbf{x}, \mathbf{x}_i) \mathbf{a}_i = \sum_{i=1}^N e^{-\frac{\|\mathbf{x}-\mathbf{x}_i\|^2}{2\sigma^2}} \mathbf{a}_i$$



RKHS and representer

Functions in \mathcal{H}_k are (limit of) combinations

$$f(\mathbf{x}) = \sum_{i=1}^N k(\mathbf{x}, \mathbf{x}'_i) c_i$$

for a sequence $\mathbf{x}'_1, \dots, \mathbf{x}'_N$.

The representer theorem ensures we can consider kernel combinations on data

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

From features and kernels to RKHS and beyond

\mathcal{H}_k has many properties, characterizations, connections:

From features and kernels to RKHS and beyond

\mathcal{H}_k has many properties, characterizations, connections:

- ▶ **Reproducing property**
- ▶ **Reproducing kernel Hilbert spaces (RKHS)**
- ▶ Mercer theorem (Karhunen Loève expansion)
- ▶ Gaussian processes
- ▶ Stochastic calculus Cameron-Martin spaces
- ▶ Physics: POVM
- ▶ Harmonic analysis: group representation and wavelet transforms

Reproducing property

By definition of \mathcal{H}_k

- ▶ $k_x = k(\mathbf{x}, \cdot) \in \mathcal{H}_k$
- ▶ Reproducing property

$$f(\mathbf{x}) = \langle f, k_x \rangle$$

for all $f \in \mathcal{H}_k$, $\mathbf{x} \in \mathbf{X}$.

- ▶ Note that

$$|f(\mathbf{x}) - f'(\mathbf{x})| \leq \|k_x\| \|f - f'\|, \quad \forall \mathbf{x} \in \mathbf{X}.$$

The above observations have a converse.

RKHS

Definition

A RKHS \mathcal{H} is a Hilbert space of functions where $\exists k : X \times X \rightarrow \mathbb{R}$ s.t.

- ▶ $k_x = k(x, \cdot) \in \mathcal{H}_k$,
- ▶ and

$$f(x) = \langle f, k_x \rangle .$$

Theorem

If \mathcal{H} is a RKHS then k is pos. def.

Proof hint: Let $\Phi(x) = k_x \dots$

Equivalent RKHS definition by evaluation functionals

If \mathcal{H} is a RKHS then the evaluation functionals

$$\mathbf{e}_x(\mathbf{f}) = \mathbf{f}(\mathbf{x})$$

are continuous. i.e.

$$|\mathbf{e}_x(\mathbf{f}) - \mathbf{e}_x(\mathbf{f}')| \lesssim \|\mathbf{f} - \mathbf{f}'\|, \quad \forall \mathbf{x} \in \mathbf{X}$$

since

$$\mathbf{e}_x(\mathbf{f}) = \langle \mathbf{f}, \mathbf{k}_x \rangle .$$

Note that $L^2(\mathbb{R}^d)$ or $C(\mathbb{R}^d)$ don't have this property³!

Alternative RKHS definition

Turns out the previous property also characterizes a RKHS.

Theorem

A Hilbert space of functions with continuous evaluation functionals is a RKHS.

Proof hint: direct application of Riesz lemma.

Feature maps & kernel/RKHS

- ▶ RKHS and PD kernels are in one to one relation.
- ▶ Feature map and kernel/RKHS are not.

Examples

- ▶ $\Phi(\mathbf{x}) = \mathbf{k}_x$ is a feature map (Aronzajn).
- ▶ If $(v_i)_i$ is a o.n.b. then $\Phi(\mathbf{x}) = (v_1(\mathbf{x}), v_2(\mathbf{x}), \dots)$ is a feature map.
- ▶ ...

Mercer theorem

A famous feature map.

Let k bounded kernel and ρ probability distribution, and

$$L_k f(\mathbf{x}) = \int k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}') d\rho(\mathbf{x}').$$

Theorem (Mercer theorem)

If $(\lambda_j, \psi_j)_j$ is the eigensystem ⁴ of L_k , then

$$k(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^{\infty} \lambda_j \psi_j(\mathbf{x}) \psi_j(\mathbf{x}')$$

where the series converges absolutely and pointwise.

► $\Phi(\mathbf{x}) = (\sqrt{\lambda_1} \psi_1(\mathbf{x}), \sqrt{\lambda_2} \psi_2(\mathbf{x}), \dots)$ is a feature map.

End of the tour

- ▶ Feature maps
- ▶ PD kernels
- ▶ RKHS

My view

"It's hard to find a useful function space which is not a RKHS".

L. Rosasco

The one exception are neural nets. But even there RKHS are useful for understanding (something...).