

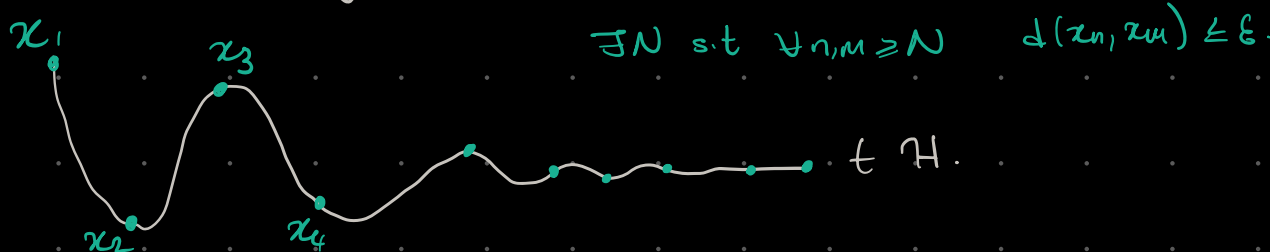
11- Kernel Methods: Basics

- Hilbert Space:

A Hilbert space \mathcal{H} is a real (or complex) inner product space that is also a complete metric space wrt the norm induced by its inner product.

Remark: Two key properties: 1- inner product 2- completeness

- Complete: Every Cauchy sequence in \mathcal{H} has a limit in \mathcal{H} .



Def (Inner prod): An inner product is a fnc $\langle \cdot, \cdot \rangle: \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$ satisfying for $f, g, h \in \mathcal{H}$

1- Symmetry: $\langle f, g \rangle = \langle g, f \rangle$

2- Linearity: $a, b \in \mathbb{R} \quad \langle af + bg, h \rangle = a \langle f, h \rangle + b \langle g, h \rangle$

3- Non-negativity: i) $\langle f, f \rangle \geq 0$

ii) $\langle f, f \rangle = 0 \iff f = 0$.

Norm induced by the inner product $\|f\|_{\mathcal{H}} = \sqrt{\langle f, f \rangle}$.

Ex (Euclidean space): $\mathcal{H} = \mathbb{R}^d$ and standard inner prod

$u, v \in \mathbb{R}^d, \quad \langle u, v \rangle = \sum_{i=1}^d u_i v_i$ which defines a norm $\|u\|_2 = \sqrt{\sum_{i=1}^d u_i^2}$.

Ex (Square integrable fncs on $[0,1]$):

$$L^2([0,1]) = \left\{ f: [0,1] \rightarrow \mathbb{R} \text{ and } \int_0^1 f(x)^2 dx < \infty \right\}$$

with inner prod $\langle f, g \rangle = \int_0^1 f(x)g(x) dx$.

Def (Dual space): Dual space \mathcal{H}^* of a Hilbert space \mathcal{H} is the space of all cnts and linear fncs from \mathcal{H} to \mathbb{R} . It carries a norm, $F \in \mathcal{H}^*$ $\|F\|_* = \sup_{\|x\|_{\mathcal{H}}=1} |F(x)|$.

Def (Linear fnc): A fnc $f: X \rightarrow \mathbb{R}$ is linear if for $x, x' \in X$ and any $c \in \mathbb{R}$, it satisfies $f(x+x') = f(x) + f(x')$.

Remark: $f(x) = ax$ is linear, but $f(x) = ax + b$ is not for $b \neq 0$.

Ex (Euclidean space): $\mathcal{H} = \mathbb{R}^d$, then its dual \mathcal{H}^*

$$\mathcal{H}^* = \left\{ F: \mathbb{R}^d \rightarrow \mathbb{R} \text{ where } F \text{ is linear (and cnts)} \right\}$$

Fncs of the form $F(x) = \langle x, u \rangle$ for some $u \in \mathbb{R}^d$ satisfy the condition in \mathcal{H}^* . Are there any more? — *

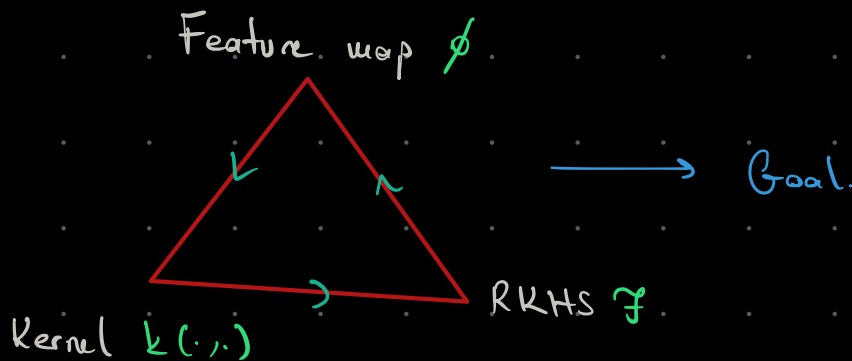
Thm (Riesz - Fréchet Representation): For every $f \in \mathcal{H}$, $\exists F_f \in \mathcal{H}^*$ unique s.t. $F_f(g) = \langle f, g \rangle$. Also, for every $F \in \mathcal{H}^*$ $\exists f_F \in \mathcal{H}$ unique s.t. $F(g) = \langle f_F, g \rangle$.

→ No! $\mathcal{H}^* = \{ F_u(x) = \langle u, x \rangle \mid u \in \mathbb{R}^d \}$

$\downarrow f$ $\downarrow \mathcal{H}$

Dual norm $\|F\|_* = \sup_{\|x\|_2=1} |\langle u, x \rangle| = \|u\|_2.$

— Kernels: Formal Definition



Def (Feature map) A fnc $\phi: \mathcal{X} \rightarrow \mathcal{H}$ where \mathcal{X} is the input space and \mathcal{H} is a Hilbert space.

Def (Kernel): A kernel is a fnc $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ s.t. for any $x_1, \dots, x_n \in \mathcal{X}$, the matrix $K_{ij} = k(x_i, x_j)$ is PSD.

Ex: Linear kernel $k(x, x') = \langle x, x' \rangle$ $\mathcal{X} = \mathbb{R}^d$

Any $x_1, \dots, x_n \in \mathbb{R}^d$, $K_{ij} = \langle x_i, x_j \rangle$

$$X = \begin{bmatrix} -x_1^T \\ \vdots \\ -x_n^T \end{bmatrix} \quad K = X X^T \quad \forall u \quad \langle u, Ku \rangle = \|X^T u\|_2^2 \geq 0.$$

Thm ($\phi \rightarrow k$): A feature map $\phi: X \rightarrow H$ defines a kernel.

proof: - $k(x, x') = \langle \phi(x), \phi(x') \rangle$

- For any x_1, \dots, x_n $K_{ij} = k(x_i, x_j)$ is PSD. \square

Thm ($k \rightarrow \phi$): For every kernel $k: X \times X \rightarrow \mathbb{R}$,

$\exists H$ a Hilbert space and a feature map $\phi: X \rightarrow H$ s.t.

$$k(x, x') = \langle \phi(x), \phi(x') \rangle.$$

proof (for finite X): Let $X = \{x_1, \dots, x_n\}$ and

$$K_{ij} = k(x_i, x_j) \Rightarrow K \text{ is PSD} \Rightarrow K = U \Lambda U^T = \Phi \Phi^T$$

$\phi(x_i) = \Lambda^{\frac{1}{2}} u_i$ defines a feature map. \square

Remark: The choice of ϕ is not unique. $\phi'(x) = Q \phi(x)$ s.t.

Q is orthogonal $Q^T Q = I$.

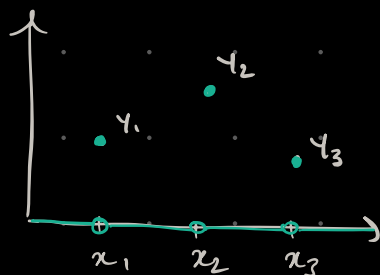
- Hilbert Space defined by the 'Reproducing Kernel'

For a dataset (y_i, x_i) $i=1 \dots n$ and $\mathcal{F} = L^2([0, 1])$.

We consider

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 + \frac{\lambda}{2} \|f\|_{\mathcal{F}}^2$$

- Choose f



$$f(x) = \sum_{i=1}^n y_i \mathbb{1}_{\{x=x_i\}}$$

$$\|f\|_{\mathcal{F}} = 0$$

\Rightarrow 0 training error

\Rightarrow overfitting

- Clearly, a Hilbert space is too complex (even has indicators).

Def (Evaluation functional): For a Hilbert space of fncs $h: X \rightarrow \mathbb{R}$, $\forall x \in X$, the evaluation functional $\mathcal{E}_x: \mathcal{H} \rightarrow \mathbb{R}$ is defined as $\mathcal{E}_x(h) = h(x)$.

Remark: Evaluation functionals are linear.

$$h, h' \in \mathcal{H} \quad \mathcal{E}_x(h+h') = (h+h')(x) = h(x) + h'(x) = \mathcal{E}_x(h) + \mathcal{E}_x(h')$$

$$c \in \mathbb{R} \quad \mathcal{E}_x(c \cdot h) = (c \cdot h)(x) = c \cdot h(x) = c \cdot \mathcal{E}_x(h)$$

Ex (Euclidean input space): $X = \mathbb{R}^d$, $\mathcal{H} = \{h_\theta(x) = \langle \theta, x \rangle \mid \theta \in \mathbb{R}^d\}$

$$\mathcal{E}_x(h_\theta) = h_\theta(x) = \langle \theta, x \rangle$$

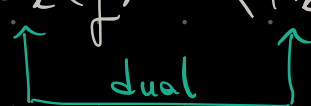
Def (RKHS): An RKHS \mathcal{F} is a Hilbert space over fncs $f: X \rightarrow \mathbb{R}$ s.t. evaluation fnc'ls are Lipschitz.

Remarks: - The constraint on the eval fnc'ls restricts \mathcal{F} .
For example, indicators no longer belong to \mathcal{F} .

- Eval. fnc'ls are cnts and linear $\Rightarrow \mathcal{E}_x \in \mathcal{F}^*$

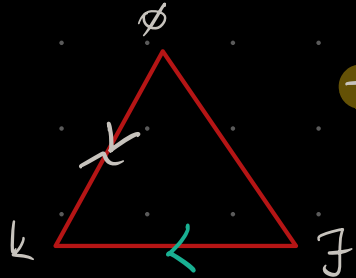
- By Riesz-Fréchet, $\exists f_{\mathcal{E}_x} \triangleq \Psi_x \in \mathcal{F}$ s.t.

$$\forall g \quad g(x) = \mathcal{E}_x(g) = \langle \Psi_x, g \rangle$$



! Ψ_x is called the representer

- Fnc evaluations can be written as inner products!



Thm ($\mathcal{F} \rightarrow k$): Every RKHS \mathcal{F} defines a unique kernel.

proof: - RKHS \Rightarrow eval fnc'ls E_x are Lipschitz cnts.

$$\Rightarrow E_x \in \mathcal{F}^*$$

- By Riesz-Fréchet, $\exists \Psi_x \in \mathcal{F}$ unique s.t.

$$\forall f \in \mathcal{F}, E_x(f) = \langle f, \Psi_x \rangle = f(x). \quad (*)$$

(Ψ_x is the representer, $(*)$ is the reproducing property)

- Since $\Psi_x \in \mathcal{F}$, for $x' \in X$, by $(*)$

$$\Psi_x(x') = E_{x'}(\Psi_x) = \langle \Psi_x, \Psi_{x'} \rangle$$

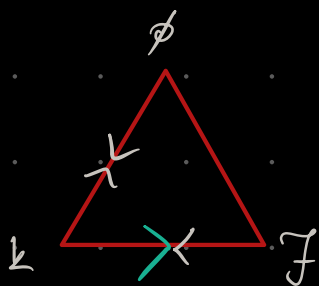
$$\triangleq k(x, x').$$

- $\Psi_x \in \mathcal{F}$, $\Psi_x = \phi(x)$ defines a feature map, which defines a kernel. \square

Remark: RKHS \mathcal{F} defines a unique kernel k called the reproducing kernel.

$$f(x) = E_x(f) = \langle f, \Psi_x \rangle = \langle f, k(x, \cdot) \rangle$$

\uparrow representer \uparrow



Thm ($k \rightarrow \mathcal{F}$; Moore-Aronszajn):

Every kernel corresponds to a unique RKHS.

proof: - Basic idea: Use $k(x, \cdot)$ as basis for RKHS.

- Let $n \in \mathbb{N}$, $x_1, \dots, x_n \in X$, $x_i, \theta_i \in \mathbb{R}$.

- Consider $f(x) = \sum_{i=1}^n x_i k(x, x_i)$ and $g(x) = \sum_{i=1}^n \theta_i k(x, x_i)$

- $\mathcal{F} = \left\{ f(x) = \sum_{i=1}^n x_i k(x, x_i) : n \in \mathbb{N}, x_1, \dots, x_n \in X, x_i \in \mathbb{R} \right\}$

(this is a vector space, but not necessarily complete)

- Define the fnc $\langle \cdot, \cdot \rangle : \mathcal{F} \times \mathcal{F} \rightarrow \mathbb{R}$ as


$$\langle f, g \rangle \triangleq \sum_{i,j} x_i \theta_j k(x_i, x_j)$$

This defines an inner product:

1. Symmetry ✓

2. Linearity ✓

3. Non-negativity: i) $\langle f, f \rangle = \sum_{i,j} x_i x_j k(x_i, x_j) = x^T K x \geq 0$.

ii) $\langle f, f \rangle = 0 \iff f = 0$.


→ Define $c(x) = [k(x, x_1) \dots k(x, x_n)]^T \in \mathbb{R}^n$

Augmented kernel matrix for $\{x_1, \dots, x_n, x\}$

$$K' = \begin{bmatrix} K & c(x) \\ c(x)^T & k(x, x) \end{bmatrix} \geq 0$$

Assume $\langle f, f \rangle = \underline{x^T K x} = 0$ but $f \neq 0 \Rightarrow x \neq 0$

For any scalar $b \in \mathbb{R}$, let $u = \begin{bmatrix} x \\ b \end{bmatrix} \in \mathbb{R}^{n+1}$

$$\begin{aligned} u^T K' u &= \begin{bmatrix} x \\ b \end{bmatrix}^T \begin{bmatrix} K & c(x) \\ c(x)^T & k(x, x) \end{bmatrix} \begin{bmatrix} x \\ b \end{bmatrix} = \underbrace{x^T K x}_{=0} + 2b x^T c(x) + b^2 k(x, x) \\ &= 2b x^T c(x) + b^2 k(x, x) \geq 0. \end{aligned}$$

But for any ξ_1 and $\xi_2 \in \mathbb{R}$, $\exists b$, $b \xi_1 + b^2 \xi_2 < 0$, so contradiction.

- Need to show all eval. functionals are Lipschitz.

$$f(x) = \sum_i x_i k(x, x_i) = \langle f, k(x, \cdot) \rangle$$

$$\text{Notice } k(x, \cdot) = \sum_{j=1}^n \theta_j k(x_j, \cdot) = 1 \cdot k(x, \cdot) \quad \left(\begin{array}{l} n=1 \\ \theta_1=1 \\ x_1=x \end{array} \right)$$

$$\langle f, k(x, \cdot) \rangle = \sum_{ij} x_i \theta_j k(x_i, x) = \sum_i x_i k(x_i, x)$$

- E_x be an evaluation functional. $\forall f, g \in \mathcal{F}$

$$\begin{aligned} |E_x(f-g)| &= |\langle f-g, k(x, \cdot) \rangle| \\ &\leq \|f-g\|_{\mathcal{F}} \|k(x, \cdot)\|_{\mathcal{F}} \quad (\text{by Cauchy-Schwarz}) \\ &= \|f-g\|_{\mathcal{F}} \sqrt{k(x, x)} \end{aligned}$$

$$\text{since } \|k(x, \cdot)\|_{\mathcal{F}}^2 = \langle k(x, \cdot), k(x, \cdot) \rangle = k(x, x).$$

- To complete the proof, one needs to consider the completion of \mathcal{F} including all limit points. This is skipped. \square

Remark: The main take away property of RKHS:

$f \in \text{RKHS}$ then $f(x) = \sum_{i=1}^n x_i k(z_i, x)$ for some $x_i \in X$ and $x_i \in \mathbb{R}$.