

# KERNEL METHODS

- Good framework for understanding linear models
- Key object: kernel functions  

$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle$$
- Precise study of statistical & approximation properties  
 in many cases
- Links to neural networks with random weights

$$\varphi(x) = (\sigma(w_1^T x), \dots, \sigma(w_m^T x)) \in \mathbb{R}^m$$

$$f(x) = \sum_{i=1}^m r_i \sigma(w_i^T x) = \underbrace{\langle r, \varphi(x) \rangle}_{\mathcal{N}(0, 1)}$$

## 1. REPRESENTER THEOREM

$\theta \in \mathbb{R}^d$  d very large, or  $d \rightarrow \infty$   
 Q: How can we learn such models in a tractable way

ERM: 
$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(y_i, \langle \theta, \varphi(x_i) \rangle) + \frac{\lambda}{2} \|\theta\|_2^2$$

Theorem (Representer theorem, (Kimeldorf & Wahba, 1971))

Let  $G(\theta) = \Psi(\langle \theta, \varphi(x_1) \rangle, \dots, \langle \theta, \varphi(x_m) \rangle, \|\theta\|_{\mathcal{H}})$

$\Psi$  is strictly increasing w.r.t. last variable.

Then, any minimizer of  $G$  takes the form:

$$\boxed{\theta = \sum_{i=1}^m \alpha_i \varphi(x_i)}$$

Remark: Let  $K = [\langle \varphi(x_i), \varphi(x_j) \rangle]_{ij} \in \mathbb{R}^{n \times n}$

We have  $\langle \theta, \varphi(x_i) \rangle = [K\alpha]_i$

$$\|\theta\|_{\mathcal{H}}^2 = \sum_{i,j} \alpha_i \alpha_j K_{ij} = \alpha^\top K \alpha$$

$\Rightarrow$  ERM becomes an optimization problem over  $\alpha \in \mathbb{R}^n$

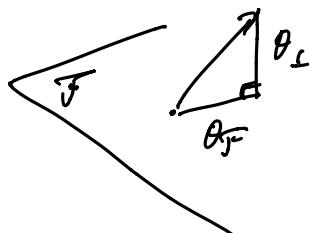
Proof: Let  $\mathcal{F} := \text{span}\{\varphi(x_i)\}_{i=1 \dots m}$

we want to show that minimizers belong to  $\overline{\mathcal{F}}$

Let  $\theta$  be a minimizer, write

$$\theta = \theta_{\mathcal{F}} + \theta_{\perp} \quad \text{with } \theta_{\mathcal{F}} \in \mathcal{F}$$

$$\theta_{\perp} \in \mathcal{F}^{\perp} \quad (\text{i.e. } \langle \theta_{\perp}, \varphi(x_i) \rangle = 0 \forall i)$$



want to show  $\theta_{\perp} = 0$

Assume, by contradiction,  $\theta_{\perp} \neq 0$

We have  $\langle \theta, \varphi(x_i) \rangle = \langle \theta_F, \varphi(x_i) \rangle$

and  $\|\theta\|_{\mathcal{H}}^2 = \|\theta_F\|^2 + \|\theta_\perp\|^2$  (Pythagorean theorem)  
 $\geq 0$   
 $> \|\theta_F\|^2$

$$\begin{aligned} G(\theta) &= \Psi\left(\langle \theta, \varphi(x_1) \rangle, \dots, \langle \theta, \varphi(x_m) \rangle, \|\theta\|_{\mathcal{H}}\right) \\ &= \Psi\left(\langle \theta_F, \varphi(x_1) \rangle, \dots, \langle \theta_F, \varphi(x_m) \rangle, \|\theta\|_{\mathcal{H}}\right) \\ &> G(\underline{\theta_F}) \end{aligned}$$

this contradicts  $\theta$  being a minimizer.

$$\Rightarrow \theta \in \mathcal{F}$$



## 2. Kernels and RKHS

Def (positive-definite kernel)

A function  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a p.d. kernel on  $\mathcal{X}$  if  
for any  $x_1, \dots, x_n \in \mathcal{X}$  and  $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ ,

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

(i.e.  $K = [k(x_i, x_j)]_{ij}$  is p.s.d.)

Examples: • linear kernel  $k(x, x') = x^T x'$

$$\begin{aligned} \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j) &= \left( \sum_i \alpha_i x_i \right)^T \left( \sum_j \alpha_j x_j \right) \\ &= \left\| \sum_i \alpha_i x_i \right\|^2 \geq 0 \end{aligned}$$

- feature map kernel  $k(x, x') = \langle \varphi(x), \varphi(x') \rangle$   
 $\varphi: X \rightarrow H$  fixed

$$\sum_{ij} \alpha_i \alpha_j k(x_i, x_j) = \left\langle \sum_i \alpha_i \varphi(x_i), \sum_j \alpha_j \varphi(x_j) \right\rangle \\ = \left\| \sum_i \alpha_i \varphi(x_i) \right\|_H^2 \geq 0$$

- quadratic kernel  $k(x, x') = (x^T x')^2 \stackrel{?}{=} \langle \phi(x), \phi(x') \rangle$

$$= x^T x' x'^T x$$

$$= Tr \left( \underbrace{x^T x'}_{x x^T} \underbrace{x'^T x}_{\cancel{x' x'^T}} \right)$$

$$= Tr(x x^T \cdot x' x'^T)$$

$$\phi(x)_{i_1 i_2 \dots i_n} = x_{i_1} x_{i_2} \dots x_{i_n} \quad = \langle x x^T, x' x'^T \rangle \\ = \langle \phi(x), \phi(x') \rangle \text{ w/ } \phi(x) = x x^T$$

- polynomial kernel  $k(x, x') = (x^T x')^n$  (homogeneous poly.)

$$k(x, x') = (1 + x^T x')^n \text{ (arbitrary poly.)}$$

Exercise: the sum of p.d. kernels is p.d.

- product

$$k(x, x') = e^{x^T x'}$$

$$k(x, x') = e^{-\alpha \|x - y\|^2}$$

Theorem [Aronszajn '50]

$k: X \times X \rightarrow \mathbb{R}$  is p.d. if and only if there exists a Hilbert space  $H$  and  $\varphi: X \rightarrow H$  such that

$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle_H$$

## Def/Theorem (RKHS)

Let  $k$  be a p.d. kernel on  $X$

Then, there exist a unique Hilbert space  $\mathcal{H}$  such that:

$$(1) \forall x \in X, k(x, \cdot) \in \mathcal{H} \quad (k(x, \cdot) \Leftrightarrow x \mapsto k(x, x'))$$

(2) [Reproducing property]

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

This is called the reproducing Kernel Hilbert Space (RKHS)  
of  $k$ .

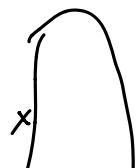
Remark:

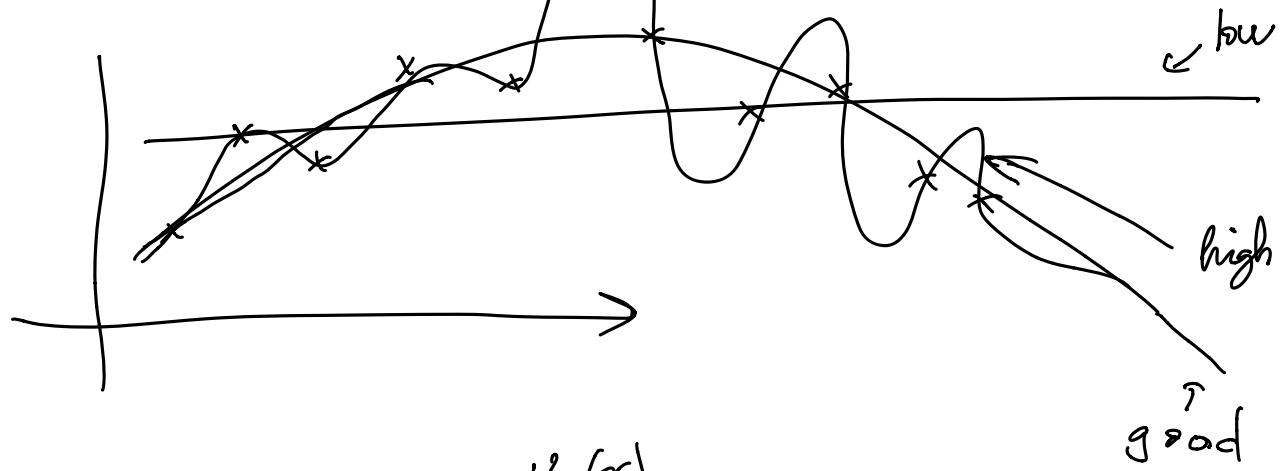
- $\langle \underbrace{k(\cdot, \cdot)}_f, k(x, \cdot) \rangle_{\mathcal{H}} = \underbrace{k(x, \cdot)}_g(x) = k(x, x') \Leftarrow$
- $\phi(x) = \underbrace{k(x, \cdot)}_{\mathcal{H}} \in \mathcal{H} \Rightarrow k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$
- $\langle f, \varphi(x) \rangle_{\mathcal{H}} = f(x)$  vs.  $\langle f, \varphi(x) \rangle_{\mathbb{R}}$

## 3. Regularization and Examples of RKHSs

Choice of kernel  $\hookrightarrow$  choice of norm  $\|f\|_{\mathcal{H}}^2$   
regularization penalty

In ML: "prior", "inductive bias", "simplicity"





### 3.1 Weighted penalties in $L^2(X)$

Assume  $X$  compact

$$\{\phi_i\}_i: \text{orthonormal basis of } L^2(X) \rightarrow \langle \phi_i, \phi_j \rangle_{L^2(X)} \\ = \int_X \phi_i(x) \phi_j(x) dx \\ = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{o/w.} \end{cases}$$

Ex:  $\cdot X = [0,1]$ ,  $\phi_i$ : Fourier basis

$$\cdot X = S^{d-1} = \{x \in \mathbb{R}^d, \|x\|=1\}$$

if  $d=1$ , wide  $\cong [0,1]$

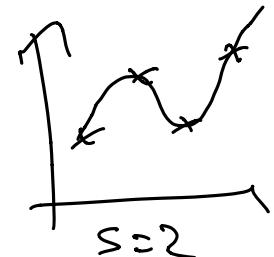
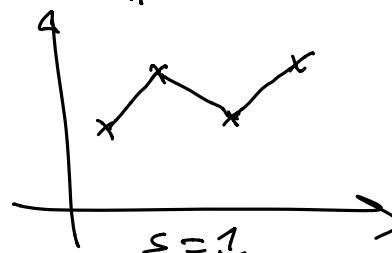
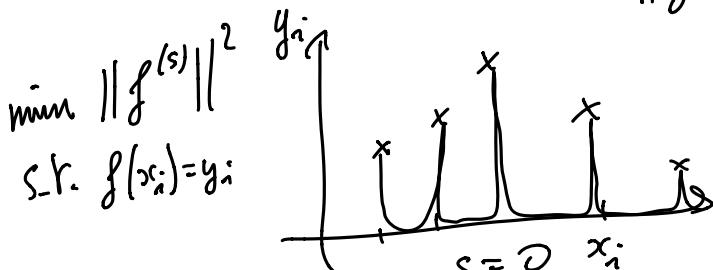
$d>1$   $\phi_i$  can be basis of Spherical Harmonics

$$L^2 \text{ norm: } f = \sum_i a_i \phi_i$$

$$\|f\|_{L^2(X)}^2 = \sum_i a_i^2 \quad (\text{Parseval})$$

$$\text{other norms? } \|f'\|_{L^2(X)}^2 \approx \sum_i (ia_i)^2 = \sum_i i^2 a_i^2$$

$$\|f^{(s)}\|^2 \approx \sum_i i^{2s} a_i^2$$



• Penalty  $\rightarrow$  Kernel?

$$\langle f, g \rangle_{\mathcal{H}} = \sum_i b_i^2 a_i(f) \cdot a_i(g)$$

$$f = \sum_i a_i \phi_i \quad \|f\|_{\mathcal{H}}^2 = \sum_i b_i^2 a_i^2 \quad b_i > 0$$

$$\text{ex: } \|f^{(k)}\|_L^2: \text{Sobolev spaces.}$$

$$b_i^2 \sim i^{-2\gamma}$$

- approach:
- . design feature map  $\phi$
  - . define corresponding kernel  $k(x, x') = \langle \phi(x), \phi(x') \rangle$
  - . check that it yields the RKHS

$$\text{Define } \phi(x) = \left(\frac{1}{b_i} \phi_i\right)_{i \in \mathbb{N}} \in \ell^2(\mathbb{N})$$

$$h(x, x') = \sum_i \frac{1}{b_i^2} \phi_i(x) \phi_i(x') \leftarrow \text{need: } \sum_i \frac{1}{b_i^2} < \infty$$

$$f = \sum a_i \phi_i$$

$$h(x, \cdot) = \sum_i \underbrace{\left| \frac{1}{b_i^2} \phi_i(x) \right|}_{\text{---}} \underbrace{\phi_i(\cdot)}$$

$$\langle f, h(x, \cdot) \rangle_{\mathcal{H}} = \sum_i b_i^2 \cdot a_i \cdot \frac{1}{b_i^2} \phi_i(x) = \sum_i a_i \phi_i(x) = f(x)$$

$\Rightarrow \mathcal{H}$  is the RKHS w/ reproducing kernel  $k$ .

⚠  $k$  is not always known in closed form.

$$\text{Fourier} \Rightarrow h(x, x') = K(x - x')$$

$$\text{Spherical H} \Rightarrow h(x, x') = K(\langle x, x' \rangle) \quad (\text{e.g. shallow NNs})$$

Kernel  $\rightarrow$  penalty / RKHS

Define  $T_h: L^2(X) \rightarrow L^2(X)$  (integral operator of  $k$ )

$$T_h f(x) = \int_X k(x, y) f(y) dy$$

Theorem (Mercer's theorem)

If  $k$  is p.d. then there exists an orthonormal basis  $\{\phi_i\}_i$  s.t.  $T_h$  is diagonal in this basis, i.e.

$$T_h \phi_i = \mu_i \phi_i, \mu_i \geq 0$$

Then,  $k(x, x') = \sum_i \mu_i \phi_i(x) \phi_i(x')$

The RKHS of  $k$  takes a similar form, with  $b_i^2 = \frac{1}{\mu_i}$

$$\|f\|_{\mathcal{H}}^2 = \sum_i \frac{a_i^2}{\mu_i}$$

Ex: •  $k(x, x') = \kappa(x - x')$   $\kappa$  is 1-periodic,  $X = [0, 1]$

$T_h$  is diagonalized in the Fourier basis

$\|f\|_X^2$  penalizes high frequencies more if  $\mu_i$  decays quickly

$$\mu_i \leftarrow \hat{\kappa}_i \text{ (Fourier coeff of } \kappa)$$

$$\bullet \quad h(x, x') = K(\langle x, x' \rangle), \quad x \in \mathbb{S}^{d-1}$$

$\mu_i \leftrightarrow$  Legendre coefs of  $K: [-1, 1] \rightarrow \mathbb{R}$   
 Gegenbauer

e.g. if  $h(x, x') = \mathbb{E}_{w \sim \mathcal{N}(0, I)} [\sigma(w^T x) \sigma(w^T x')]$

thus  $h(Vx, Vx') = h(x, x')$  if  $V^T V = I \Rightarrow h(x, x') = K(\langle x, x' \rangle)$   
 (rotation-invariant)

→ for  $\sigma(u) = \max(u, 0)$ :

then  $K(u) = \frac{1}{\pi} (u \cdot (\pi - \arccos(u)) + \sqrt{1-u^2})$

(arc-cosine kernel)

[Cho & Saul 2009]

[Bach 2017]

Convolutional  
Kernels:

$$h(x, x') = \sum_p \underbrace{h(x_p, x'_p)}_{K(x_p, x'_p)} \quad \begin{matrix} \text{patches} \\ \text{(Mairal 2016)} \end{matrix}$$

(Bietti-Mairal 2019)

$$= \langle \phi(x), \phi(x') \rangle_H \quad H \neq \text{RKHS}$$

$$\phi(x) = (\bar{\varphi}(x_1), \dots, \bar{\varphi}(x_p))$$

]

### 3.2. Kernels from feature maps

Feature map  $\phi: X \rightarrow H$ ,  $H$  Hilbert space (not RKHS)

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_H$$

Q: What is the RKHS for such kernels?

Ex: sum kernel  $k = k_1 + k_2$

$\varphi_1: X \rightarrow \mathcal{H}_1$  RKHS of  $k_1$

$\varphi_2: X \rightarrow \mathcal{H}_2$  RKHS of  $k_2$

$$h(x, x') = \langle \phi(x), \phi(x') \rangle_H$$

$$\phi(x) = (\varphi_1(x), \varphi_2(x)) \in H = \mathcal{H}_1 \oplus \mathcal{H}_2$$

Infinite-width NN:

$$h(x, x') = \mathbb{E}_{w \sim \mathcal{G}} \{ \sigma(w^T x) \sigma(w^T x') \}$$

$$= \int \sigma(w^T x) \sigma(w^T x') d\tau(w)$$

$$= \langle \phi(x), \phi(x') \rangle_{L^2(d\tau)}$$

$$\phi(x) = \sigma(\langle \cdot, x \rangle) \in L^2(d\tau)$$

$$h(x, x') = \frac{1}{m} \sum_{i=1}^m \sigma(w_i^T x) \sigma(w_i^T x') = \langle \phi(x), \phi(x') \rangle_{\mathbb{R}^m}$$

(Random Features: [Bach 2017b]  
[Rudi-Rosasco 2017])

Theorem The RHS of  $k(x, x') = \langle \phi(x), \phi(x') \rangle$   
is given by

$$\mathcal{H} = \{ \langle \theta, \phi(\cdot) \rangle, \theta \in H \}$$

$$\|f\|_{\mathcal{H}}^2 = \min_{\theta \in H} \|\theta\|_H^2$$

s.t.  $f = \langle \theta, \phi(\cdot) \rangle$

(e.g. [Bach, 2017, App.A] for proof)

Ex: (sum Kernel)

$$\|f\|_{\mathcal{H}}^2 = \min_{f_1, f_2} \|f_1\|_{\mathcal{H}_1}^2 + \|f_2\|_{\mathcal{H}_2}^2$$

s.t.  $f = f_1 + f_2$

• (Net)

$$f(x) = \int p(\omega) \sigma(w^T x) d\tau(\omega)$$

with  $p \in L^2(d\tau)$

$$\|f\|_{\mathcal{H}}^2 = \min_p \|p\|_{L^2(d\tau)}^2$$

s.t.  $f(x) = \int p(\omega) \sigma(w^T x) d\tau(\omega)$

finite width:  $f(x) = \sum_i p_i \sigma(w_i^T x)$

⚠  $f(x) = \sigma(w^T x) \Rightarrow$  not always in RKHS !  
(need  $p(\omega) \rightarrow$  Dirac on  $w^*$ )