

Lecture 4: Reproducing Kernel Hilbert Spaces

→ Last week: we constructed very expressive function classes

$$H = \text{span } \{ \sigma(\langle \cdot, w \rangle + b) ; (w, b) \in \mathbb{R}^d \}$$

→ We verified (VAT) that these are very expressive
↳ From Fourier perspective.

Today: (i) study the framework for linear approximation
in the high-dimensional limit \rightarrow Kernel methods
(ii) relate these spaces with (shallow) NNs.

— — — — — — features.

→ We will focus on linear models, ie functions

$$f_\theta : X \rightarrow \mathbb{R} \text{ of the form } f_\theta(x) = \langle \theta, \varphi(x) \rangle_H$$

where $\varphi : X \rightarrow H$ H a Hilbert space of very high
(even infinite) dimension. $\theta \in H$.

→ The key object to understand these models is the
kernel function $k(x, x') = \langle \varphi(x), \varphi(x') \rangle_H$.

→ Linear / Kernel-based models are interesting for two reasons

(i) allows for a "complete" analysis of learning
(approx, statistical + optim errors can be analyzed).

(ii) enables us to better understand Neural Nets.

→ Running motivating example: $\Phi(x) = (\Gamma(w_1^T x), \dots, \Gamma(w_m^T x))$
shallow Neural Net with random $w_i \sim N(0, I)$.
weights provides random features.

Q: How can we do anything practical out of an infinite-dimensional space?

Representer Theorem

• Consider the optimization problem: given data $\{(x_i, y_i)\}_{i=1..n}$

$$\min_{\theta \in H} \frac{1}{n} \sum_{i=1}^n l(y_i, \langle \theta, \Phi(x_i) \rangle) + \frac{\lambda}{2} \|\theta\|_H^2 \quad X \times Y.$$

→ Even though the model exists in an infinite-dim space,
its effective dimension once it interacts with data is finite.

Theorem (Representer Thm, Kindeldorf, Wahba '71) Let
 $\Phi: X \rightarrow H$ and $(x_1 \dots x_n) \in X^n$. Assume that
the functional $\underline{\Psi}: \mathbb{R}^{n+1} \rightarrow \mathbb{R}$ is strictly increasing
along its last variable. Then

$$\inf_{\theta \in H} \Psi(\langle \theta, \Phi(x_1) \rangle, \dots, \langle \theta, \Phi(x_n) \rangle, \|\theta\|^2)$$

can be obtained by restricting θ to

$$H_D := \text{span} \left\{ \varphi(x_i) \right\}_{i=1 \dots n}$$

Proof: Let $\theta \in H$. We use the fact that H is Hilbert + decompose $\theta = \theta_D + \theta_{\perp}$ with $\theta_D \in H_D$ $\theta_{\perp} \in (H_D)^{\perp}$

key observations: $\forall i \quad \langle \theta, \varphi(x_i) \rangle = \langle \theta_D, \varphi(x_i) \rangle$

$$\|\theta\|^2 = \|\theta_D\|^2 + \|\theta_{\perp}\|^2$$

$$\Psi(\langle \theta, \varphi(x_1) \rangle, \dots, \langle \theta, \varphi(x_n) \rangle, \|\theta\|^2) =$$

$$= \Psi(\langle \theta_D, \varphi(x_1) \rangle, \dots, \langle \theta_D, \varphi(x_n) \rangle, \|\theta_D\|^2 + \|\theta_{\perp}\|^2)$$

Ψ is non-decreasing $\hookrightarrow (\geq) \quad \Psi(\langle \theta_D, \varphi(x_1) \rangle, \dots, \langle \theta_D, \varphi(x_n) \rangle, \|\theta_D\|^2) \quad \square$

\rightarrow Remark: no convexity assumption is needed!

\rightarrow Now we can write $\theta_D \in H_D$ as

$$\theta = \sum_{i=1}^n \alpha_i \varphi(x_i) \quad \alpha \in \mathbb{R}^n$$

\rightarrow Let $K \in \mathbb{R}^{n \times n}$ $K_{ij} = K(x_i, x_j) = \langle \varphi(x_i), \varphi(x_j) \rangle_H$.

Then $\forall j, \quad \langle \theta, \varphi(x_j) \rangle = \sum \alpha_i K(x_i, x_j) = (K\alpha)_j$

$$0 \leq \|\theta\|^2 = \left\langle \sum_i \alpha_i \varphi(x_i), \sum_i \alpha_i P(\varphi_i) \right\rangle = \alpha^T K \alpha \geq 0$$

L) In particular K psd + symmetric.

→ So our learning problem becomes

$$\begin{cases} \inf_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n \ell(y_i, (K\alpha)_i) + \frac{\lambda}{2} \alpha^T K \alpha \\ f(x) = \sum_{i=1}^n \alpha_i K(x, x_i) \end{cases}$$

Neither during training
nor testing we need
features $\varphi(x)$

only Kernel $K(x, x')$. "Kernel Trick".

Q: How to build kernels? What are some natural examples? (+ How do the resulting spaces relate to what we saw last week?)

- Kernels and RKHS (Reproducing Kernel Hilbert Space).

→ We just saw that a kernel, when it intersects with data, defines a symmetric psd kernel matrix. In fact, this is precisely what defines a pd kernel:

df: a function $K: X \times X \rightarrow \mathbb{R}$ is a pd kernel
iff all associated kernel matrices are symmetric psd.

examples: • linear kernel: $K(x, x') = x^T x'$

$$\sum_{ij} \alpha_i \alpha_j K(x_i, x_j) = \left(\sum_i \alpha_i x_i \right)^T \left(\sum_{i' > i} \alpha_{i'} x_{i'} \right)$$

$$= \left\| \sum_i \alpha_i x_i \right\|^2 > 0.$$

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

$$\sum_{i,j} \alpha_i \alpha_j K(x_i, x_j) = \left\| \sum_i \alpha_i \varphi(x_i) \right\|_H^2 > 0.$$

- quadratic kernel: $K(x, x') = (x^T x')^2 = \langle \varphi(x), \varphi(x') \rangle$

$$K(x, x') = x^T x' (x')^T x = \text{Tr}(\underbrace{x^T x'}_{\mathbb{R}^{n \times n}} (x')^T x) = \text{Tr}(\underbrace{x x^T}_{\mathbb{R}^{n \times n}} \underbrace{x' (x')^T}_{\mathbb{R}^{n \times n}})$$

$$= \langle x x^T, x x^T \rangle \quad (\varphi(x) = \underline{\underline{x x^T}}).$$

→ Does the representation $K(x, x') = \langle \varphi(x), \varphi(x') \rangle_H$ always exist? A: Yes!

Theorem (Aronszajn' 50) K is a pd kernel iff
there exists a Hilbert space H and a function
 $\varphi: X \rightarrow H$ such that $\forall x, x' \quad K(x, x') = \langle \varphi(x), \varphi(x') \rangle_H$.

Prof. [Chapter 7 of Francis Bach's new book MCFP].

Def/Thm. [RKHS] Given a pd kernel, $K: X \times X \rightarrow \mathbb{R}$,
there exists a unique Hilbert space H with the following
 properties:

(i) $K(x, \cdot) : X \rightarrow \mathbb{R}$ belongs to H $\forall x \in X$.

(ii) (reproducing property) function evaluation are in H :

$$\forall x \in X, \forall f \in H \quad \langle f, K(x, \cdot) \rangle_H = f(x)$$

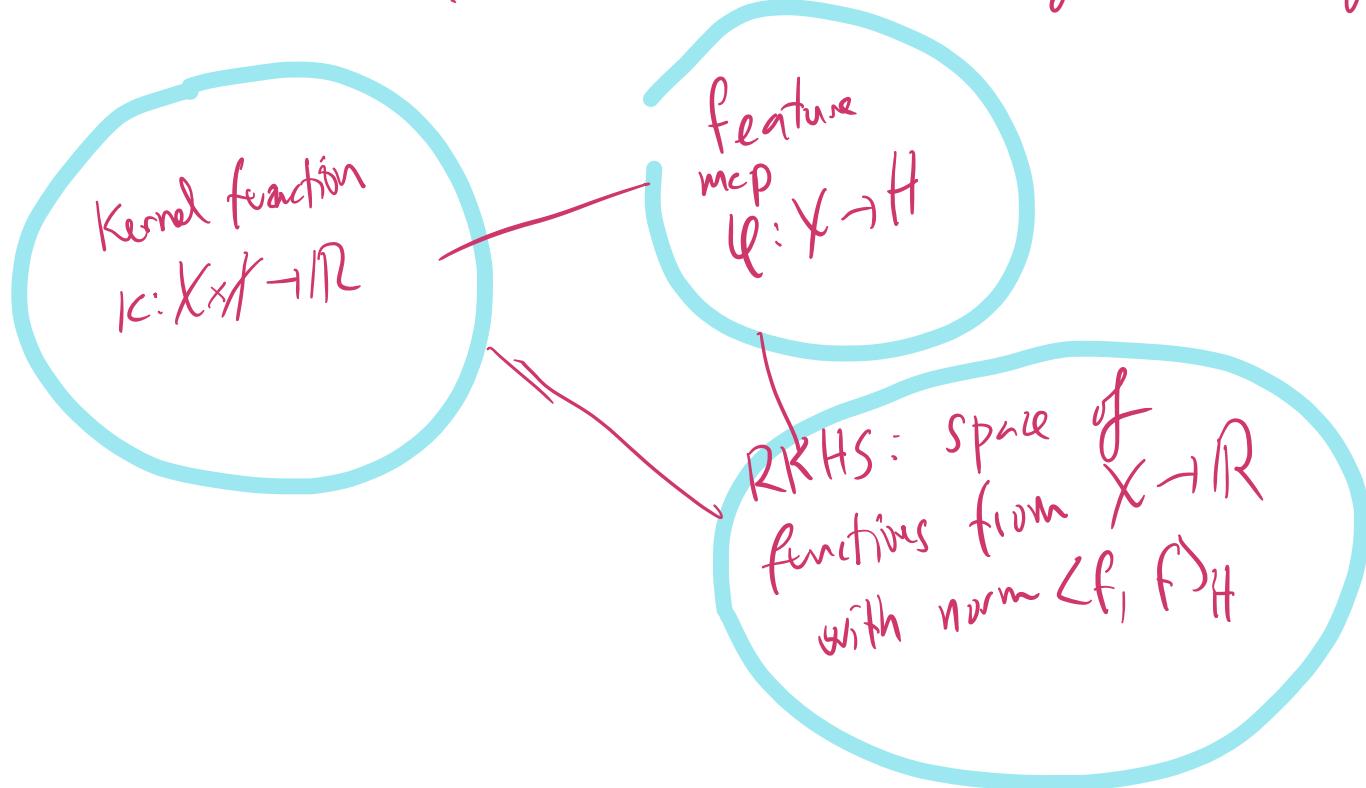
Q: Is $L_2(\mathbb{R}^d)$ an RKHS? Hell no!

Suppose it is: then there would exist a function

$K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\forall f \in L_2 \quad \forall x \in \mathbb{R}^d$

$$\langle f, K(x, \cdot) \rangle = \int_{\mathbb{R}^d} f(y) [K(x, y)] dy = f(x)$$

That would imply that $K(x, \cdot) dy = \delta_x$ Dirac
at x !!. Dirac has no density wrt Lebesg.!



Kernels give us a device to build hypothesis space via
Haar associated RKHS., keeping nice hilbertian structure.
What are the tradeoffs?

Choice of Kernel \leftrightarrow Choice of smoothness
"inductive bias"

- Assume compact input space X .

- Let $\{\phi_i\}_{i=1\dots}$ orthonormal basis of $L_2(X)$.

Ex: $X = [0, 1] \setminus \{\phi_i\}$: Fourier basis ($\phi_i = e^{j2\pi i t}$)

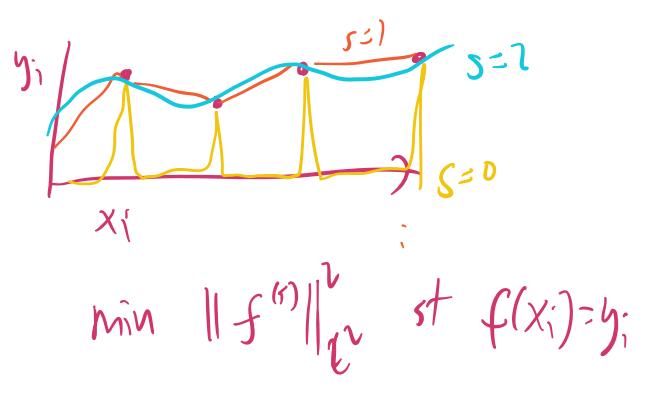
$X = S^{d-1}$ $\setminus \{x \in \mathbb{R}^d ; \|x\|_2 = 1\}$ $\setminus \{\phi_i\}$ of spherical harmonics
(basically it is the natural equivalent of Fourier w.)

$f \in L^2(X) ; f = \sum_{j=1}^{\infty} \alpha_j \phi_j \quad \|f\|_{L^2(X)}^2 = \sum_j |\alpha_j|^2$

We already saw that regularity of f can be measured
in the F domain by penalising this L_2 norm; eg

$$\|f'\|_{L^2}^2 = \sum_j j^2 |\alpha_j|^2$$

$$\left(\|f^{(s)}\|_{L^2}^2 \right) = \sum_j j^{2s} |\alpha_j|^2$$



Penalty \rightarrow Kernel ?

$$f = \sum_j a_j \phi_j \quad \text{Consider a penalty}$$

$$\sum_j b_j^2 |a_j|^2 = \|f\|_H^2$$

$$\langle f, g \rangle_H = \left(\sum_j b_j^2 a_j(f) \cdot a_j(g) \right)$$

We define $\varphi(x) = (b_j^{-1} \phi_j(x))_{j \in \mathbb{N}}$ in $\ell_2(\mathbb{N})$

$$K(x, x') := \sum_j \frac{1}{b_j^2} \phi_j(x) \phi_j(x') \quad (= \langle \varphi(x), \varphi(x') \rangle_{\ell_2(\mathbb{N})})$$

$$K(x, x) = \sum_j b_j^{-2} < +\infty \quad \sum_j \frac{1}{b_j} < +\infty \text{ is necessary.}$$

$$f = \sum_j a_j \phi_j$$

$$K(x, \cdot) = \sum_j \frac{1}{b_j^2} \phi_j(x) \phi_j(\cdot)$$

$$\langle f, K(x, \cdot) \rangle_H = \sum_j b_j^{-2} a_j \cdot \frac{1}{b_j^2} \phi_j(x) = \sum_j a_j \phi_j(x) = f(x)$$

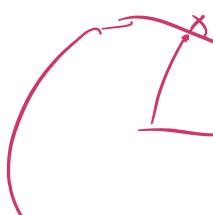
In the Fourier case,

$$K(x, x') = h(x - x')$$

in particular it is
translation invariant.

in the spherical case

$$K(x, x') = h(\langle x, x' \rangle)$$

 the kernel only depends on the angle between

* Kernel \rightarrow Penalty on $L_2(X)$?

Define $T_K : L_2(X) \rightarrow L_2(X)$ the integral op.

$$(T_K f)(x) = \int_X K(x, y) f(y) dy$$

Theorem [Mercer] If K is a pd kernel, then there exists an orthonormal basis $\{\phi_j\}$ st T_K is diagonal, $T_K \phi_j = \lambda_j \phi_j$ $\forall j$, $\lambda_j > 0$.

$$\rightarrow \text{As a result, } K(x, x') = \sum_j \lambda_j \phi_j(x) \phi_j(x')$$

\rightarrow The RKHS can be embedded inside $L_2(X)$

$$\|f\|_H^2 = \sum_j (\lambda_j^{-1}) \langle f, \phi_j \rangle^2$$

Ex: $K(x, x') = h(\underbrace{\langle x, x' \rangle}_S)$, $x, x' \in S^{d-1}$

? $\{\phi_j\}$? $\{\lambda_j\}$?

spherical harmonics.

$$K = S^{d-1} \\ L(S^{d-1})$$

$$h: [-1, 1] \rightarrow \mathbb{R}$$

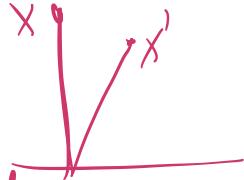
$\lambda_j \leftrightarrow$ Legendre polynomial decomposition of h .

$$\rightarrow \text{If } K(x, x') = \mathbb{E}_{w \sim N(0, I)} [\nabla(w^T x) \nabla(w^T x')]$$

$$K(Ux, Ux') = K(x, x') \quad \text{if } U^T U = \text{Id},$$

\hookrightarrow since $N(0, I)$ is Rot. inv.

$$\exists \quad K(x, x') \leq h(\langle x, x' \rangle)$$



For $\nabla(t) = \max(0, t)$, the associated h is

$$h(u) = T^{-1} \left(u(\pi - \arccos(u)) + \sqrt{1-u^2} \right)$$

[Cho & Saul, 08].

\rightarrow For more general activations [Le Roux & Bengio]

[Francis Bach '17].

Bottom line: the smoother h is, the "smaller" the associated RKHS.

Kernels from Feature Maps

Given a feature map $\phi: X \rightarrow \mathcal{H}$ Hilbert space

$$K(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}} \quad H \text{ is not necessarily RKHS.}$$

Q: What is the RKHS associated with K ?

Ex Infinite width shallow NN: \mathcal{E} : finite measure over H^*

$$\begin{aligned}
 K(x, x') &= \mathbb{E}_{w \sim \mathcal{T}} \left[\sigma(w^\top x) \sigma(w^\top x') \right] \\
 &= \int_{\mathbb{R}^d} \sigma(w^\top x) \sigma(w^\top x') d\mathcal{T}(w) \\
 &= \langle \phi(x), \phi(x') \rangle_{L_2(\mathbb{R}^d, d\mathcal{T})}
 \end{aligned}$$

$\phi(x) = \sigma(\langle \cdot, x \rangle) \in L_2(d\mathcal{T})$

→ There is a closely related kernel,

$$\begin{aligned}
 \hat{K}(x, x') &= \langle \phi(x), \phi(x') \rangle_{L_2(\mathbb{R}^d, d\mathcal{T})} \\
 &= \frac{1}{m} \sum_{j=1}^m \sigma(\langle x, w_j \rangle) \sigma(\langle x', w_j \rangle) \\
 &\quad \text{where } w_1, \dots, w_m \stackrel{\text{iid}}{\sim} \mathcal{T} \\
 &= \langle \hat{\phi}(x), \hat{\phi}(x') \rangle_{\mathbb{R}^m}
 \end{aligned}$$

$K \rightsquigarrow (\mathcal{H})$? RICHS over $X \rightarrow \mathbb{R}$ [Rahimi Recht]

$\hat{K} \rightsquigarrow \hat{\mathcal{H}}$ a different from \mathcal{H} .

$$\begin{aligned}
 \text{Good news: as } m \rightarrow \infty \quad \hat{K} &\rightsquigarrow K \\
 &\Rightarrow \hat{\mathcal{H}} \rightarrow \mathcal{H}
 \end{aligned}$$

Theorem: The RKHS of $K(x, x) = E(\delta(x, w) \otimes w, x))$ is given by

$$H = \left\{ f(x) = \int \Gamma(\langle x, w \rangle) g(w) d\pi(w); g \in L_2(d\pi) \right\}$$

$$\text{with } \|f\|_H = \inf \left\{ \|g\|_{L_2(d\pi)} : f(x) = \int \Gamma(\langle x, w \rangle) g(w) d\pi(w) \right\}$$

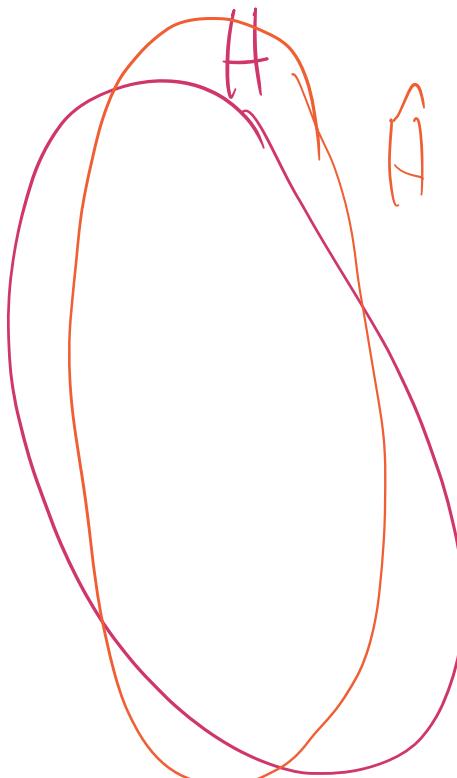
Prof: [Bach '17 App A].

Remark: $f(x) = \Gamma(\langle w_0, x \rangle)$ belongs to the RKHS?

Does

No! in general (Γ) we would need $g(w) = \int_{w_0}$
which is not in $L_2(d\pi)$!

$$H \subseteq L^2(X)$$



Pick $f \in \{f' ; \|f'\|_H \leq 1\}$

$$\sup_{f \in B_1(H)} \|f\|_H \quad B_1(H)$$