

Q

[Bonus: Hilbert spaces

$(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$

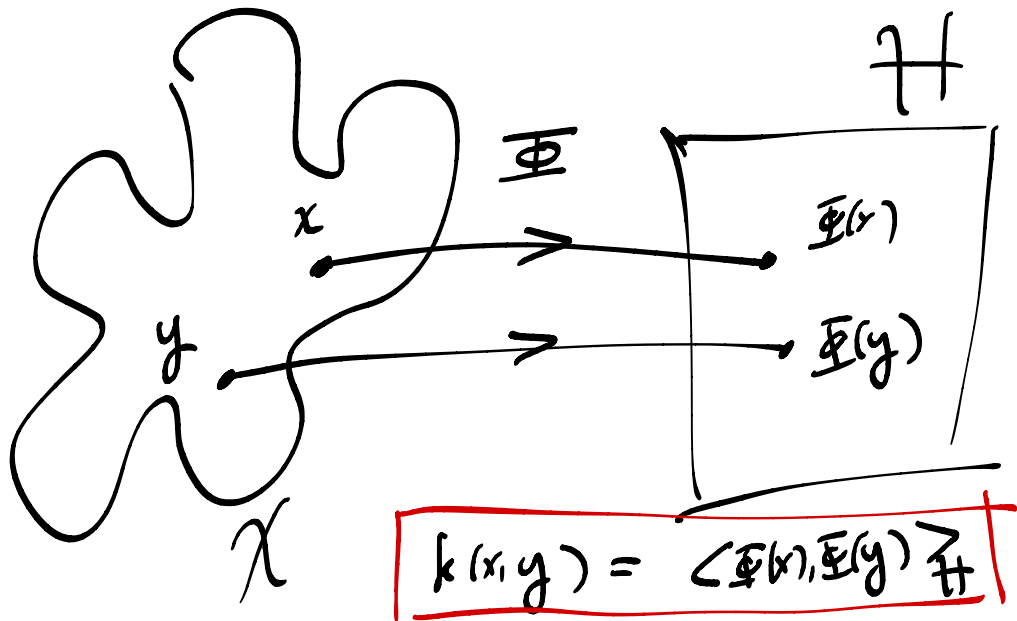
Definition: A Hilbert space is a real or complex vector space which is also a complete metric space with respect to the distance function induced by the inner product.

- ▶ **Remark:** recall the linear kernel, all we used were properties of inner product
- ▶ let $\Phi : \mathcal{X} \rightarrow \mathcal{H}$ be some mapping, \mathcal{H} a Hilbert space with scalar product $\langle \cdot, \cdot \rangle$
- ▶ then $k(x, y) = \langle \Phi(x), \Phi(y) \rangle$ is positive definite:

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j \langle \Phi(x_i), \Phi(x_j) \rangle = \left\| \sum_{i=1}^n c_i \Phi(x_i) \right\|^2 \geq 0,$$

by linearity of the inner product.

→ large class of kernels: "nice" vector space \mathcal{H} + embedding of \mathcal{X} in \mathcal{H}



$k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, sym., Gram matrix K is psd.
Kernel as inner products

all $k(x_i, x_j)$
($1 \leq i, j \leq n$)

- ▶ Remarkable fact: the converse statement is true!

Theorem:⁹ For any kernel k on \mathcal{X} , there exists a Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle)$ and a mapping $\Phi: \mathcal{X} \rightarrow \mathcal{H}$ such that

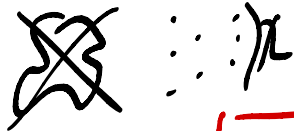
$$\forall x, y \in \mathcal{X},$$

$$k(x, y) = \langle \Phi(x), \Phi(y) \rangle.$$

$$(\Phi = \varphi)$$

- ▶ **Reminder:** Hilbert space = inner product + *complete* for the associated norm (Cauchy sequences converge in \mathcal{H})
- ▶ **Consequence:** we can think of any kernel as a dot product in the feature space
- ▶ **Main idea:** forget about Φ and work only with kernel evaluations (more on that later)

⁹Aronszajn, *Theory of reproducing kernels*, Transactions of the American Mathematical Society, 1950



$$h(x, y) = k(x_i, x_j)$$

Proof in the finite case

- ▶ assume that $\mathcal{X} = \{x_1, \dots, x_N\}$ is finite of size N
- ▶ any kernel k is entirely defined by the $N \times N$ positive semi-definite matrix $K := (k(x_i, x_j))_{i,j=1}^N$
- ▶ we can diagonalize K in an orthonormal basis (u_1, \dots, u_N) with associated (non-negative) eigenvalues $\lambda_1, \dots, \lambda_N$: $K = U\Lambda U^\top$, with $U_{:,i} = u_i$, $\Lambda = \text{diag}(\lambda)$, $UU^\top = U^\top U = I$
- ▶ then we write

$$\begin{aligned} k(x_i, x_j) &= \left(\sum_{\ell=1}^N \lambda_\ell u_\ell u_\ell^\top \right)_{i,j} \\ &= \sum_{\ell=1}^N \lambda_\ell (u_\ell)_i (u_\ell)_j = \langle \Phi(x_i), \Phi(x_j) \rangle, \end{aligned}$$

with

$$\Phi(x_i) := \left(\sqrt{\lambda_1} (u_1)_i, \dots, \sqrt{\lambda_n} (u_N)_i \right)^\top.$$



Reminder: $K = U \Lambda U^T \in \mathbb{R}^{N \times N}$ real
✓

with $U \in \mathbb{R}^{N \times N}$ $U U^T = U^T U = I_N$ (identity) since K sym.

additionally, $K \succeq 0$

\Downarrow
 $\forall i \in [N], \lambda_i \geq 0$

$$\Lambda = \text{Diag}(\lambda_1, \dots, \lambda_N) \\ = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_N \end{pmatrix}$$

$$\Rightarrow K = \sum_{\ell=1}^N \lambda_{\ell} \mu_{\ell} \mu_{\ell}^T \quad \text{with } \mu_{\ell} = \text{colth column at } \ell.$$

\uparrow
 eigenvalues of K

$$\begin{aligned}
 U \Lambda &= \left(\mu_1 \mid \mu_2 \mid \dots \mid \mu_N \right) \begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \ddots \\ & & & \lambda_N \end{pmatrix} \\
 &= \left(\lambda_1 \mu_1 \mid \lambda_2 \mu_2 \mid \dots \mid \lambda_N \mu_N \right)
 \end{aligned}$$

$$K = U \Lambda U^T = \begin{pmatrix} \lambda_1 \mu_1 & \dots & \lambda_N \mu_N \end{pmatrix} \begin{pmatrix} \frac{\mu_1^T}{\mu_N^T} \\ \vdots \\ \frac{\mu_N^T}{\mu_N^T} \end{pmatrix} = \sum_{\ell=1}^N \lambda_\ell \mu_\ell \mu_\ell^T \quad \text{pos.}$$

$\mathbb{R}^{N \times N} \ni \mu_i, \mu_j^T$? $U U^T = I_N$

$$\Rightarrow \begin{pmatrix} \mu_1 & \dots & \mu_N \end{pmatrix} \begin{pmatrix} \frac{\mu_1^T}{\mu_N^T} \\ \vdots \\ \frac{\mu_N^T}{\mu_N^T} \end{pmatrix} = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}$$

$$\mu_i \mu_j^T = \begin{cases} 1 & i=j \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned}
 k(x_i, x_j) &= K_{ij} = \left(\sum_{l=1}^N \lambda_l \mu_l \mu_l^T \right)_{ij} \\
 &= \sum_{l=1}^N \lambda_l (\mu_l \mu_l^T)_{ij} \\
 &= \sum_{l=1}^N \lambda_l \mu_{l,i} \mu_{l,j}
 \end{aligned}$$

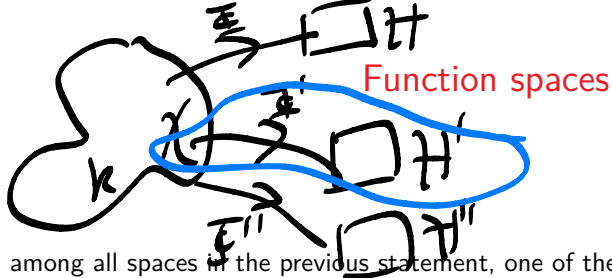
Diagram illustrating the matrix structure and element-wise multiplication:

$$\begin{pmatrix} \mu_{1,1} \\ \vdots \\ \mu_{1,N} \end{pmatrix} \begin{pmatrix} \mu_{1,1} & \mu_{1,2} & \dots & \mu_{1,N} \end{pmatrix}$$

Diagram illustrating the element-wise multiplication of two vectors:

$$\begin{pmatrix} \mu_{1,i} \\ \vdots \\ \mu_{N,i} \end{pmatrix} \times \begin{pmatrix} \mu_{1,j} \\ \vdots \\ \mu_{N,j} \end{pmatrix} = \begin{pmatrix} \mu_{1,i} \mu_{1,j} \\ \vdots \\ \mu_{N,i} \mu_{N,j} \end{pmatrix}$$

5.2. Reproducing kernel Hilbert spaces (RKHS)

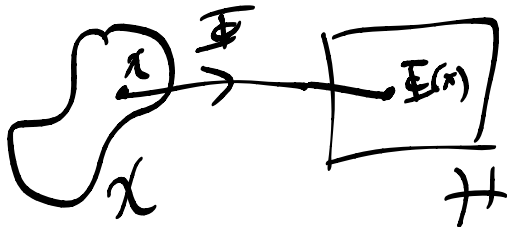


Function spaces

very nice

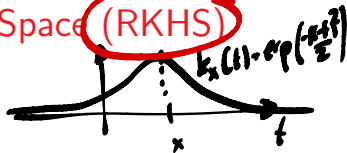
- ▶ among all spaces in the previous statement, one of them has interesting properties
- ▶ in particular, it is a **space of functions**
- ▶ i.e., we can map each point $x \in X$ to a *function* $\Phi(x) = k_x \in \mathcal{H}$
- ▶ **Example:** $X = \mathbb{R}$, we map each x to the function $t \mapsto xt$
- ▶ \rightarrow space of linear functions
- ▶ more complicated in general...

~~$F(x)(y)$~~



Reproducing Kernel Hilbert Space (RKHS)

Ex: $k(x, y) = \exp(-(y-x)^2)$



Definition: let \mathcal{X} be a set and \mathcal{H} be a function space forming a Hilbert space with inner product $\langle \cdot, \cdot \rangle$. The function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is called a *reproducing kernel* of \mathcal{H} if

- ▶ \mathcal{H} contains all functions of the form $k_x : t \mapsto k(x, t)$
- ▶ for every $x \in \mathcal{X}$ and $f \in \mathcal{H}$, the *reproducing property* holds:

$$\boxed{f} \in \mathcal{H}$$

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

- ▶ if a reproducing kernel exists, then \mathcal{H} is called a *reproducing kernel Hilbert space* (RKHS)



Equivalent definition

Theorem: the Hilbert space $\mathcal{H} \subseteq \mathbb{R}^X$ is a RKHS if, and only if, for any $x \in X$, the mapping $f \mapsto f(x)$ is continuous.

Bonus

- ▶ Proof of \Rightarrow : let k be a reproducing kernel, $x \in X$ and $f_n \rightarrow f$ in \mathcal{H}
- ▶ we write

$$|f_n(x) - f(x)| = |\langle f_n - f, k_x \rangle|$$

$$\leq \|f_n - f\|_{\mathcal{H}} \|k_x\|_{\mathcal{H}}$$

by Cauchy-Schwarz inequality.

- ▶ $\|f_n - f\| \rightarrow 0$ and we can conclude

- ▶ **Remark:** $\|k_x\|^2 = \langle k_x, k_x \rangle = k(x, x)$, thus $|f(x)| \leq \|f\|_{\mathcal{H}} \cdot k(x, x)^{1/2}$

$\exp(-(x-x')^2) = 1$

$= 1$

reproducing prop: $f(x) = \langle f, k_x \rangle$
 $|f(x)| = \langle f, k_x \rangle$
 fixed, finite
 $|f(x)| \leq \|f\|_{\mathcal{H}} \cdot \|k_x\|_{\mathcal{H}}$

Continuity ctd.

Bonus.

- ▶ Proof of \Leftarrow : let $x \in \mathcal{X}$
- ▶ by the reproducing property, $L : x \mapsto f(x)$ is a *linear functional*
- ▶ Riesz theorem there exists ℓ_x such that $L(x) = \langle f, \ell_x \rangle$
- ▶ define $k(x, y) := \ell_y(x)$
- ▶ one can check readily the RKHS properties.





→ Riesz-Fréchet

Uniqueness

Theorem: if \mathcal{H} is a RKHS, then it has a unique reproducing kernel. Conversely, a function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ can be the reproducing kernel of at most one RKHS.

- ▶ we talk about the RKHS associated to k
- ▶ *Proof:* let k and k' be two reproducing kernels
- ▶ then for all $x \in \mathcal{X}$,


$$\begin{aligned}\|k_x - k'_x\|^2 &= \langle k_x - k'_x, k_x - k'_x \rangle = \langle k_x, k'_x \rangle - \langle k'_x, k_x \rangle \\ &= \cancel{k_x(x)} - \cancel{k'_x(x)} - \cancel{k_x(x)} + \cancel{k'_x(x)} = 0.\end{aligned}$$

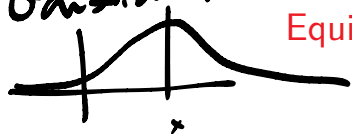
my. rep



Ex: Gaussian kernel.

Equivalence psd / RKHS

k_x



Theorem: a function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is positive definite if, and only if, it is a reproducing kernel. $(\mathcal{H}, \langle \cdot, \cdot \rangle)$

► **Idea:** build \mathcal{H} as the completion of

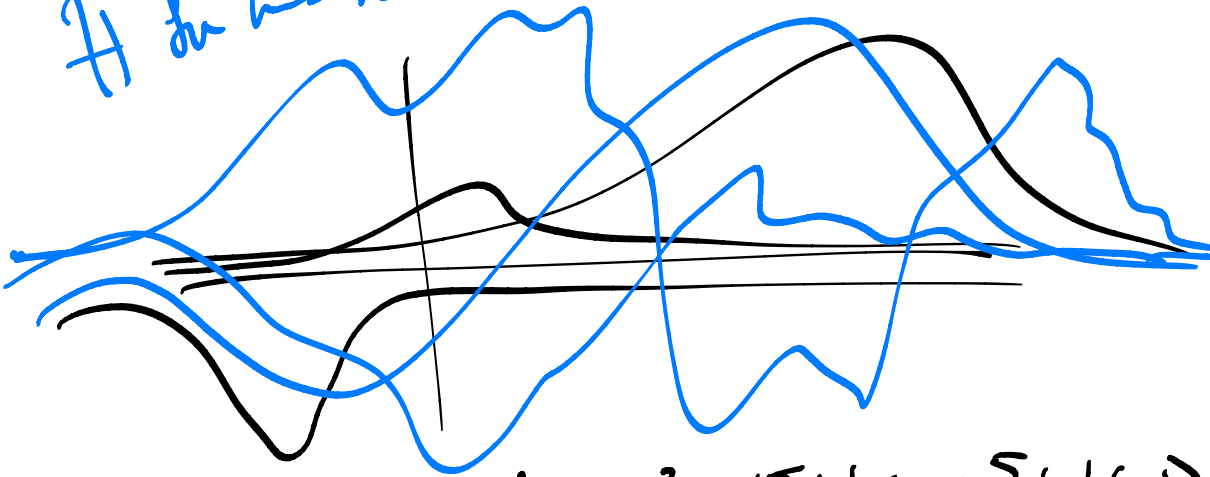
$$\mathcal{H}_0 := \left\{ \sum_{i=1}^n \alpha_i \underbrace{k(\cdot, x_i)}_{k_{x_i}}, n \in \mathbb{N}, \alpha_i \in \mathbb{R}, x_i \in \mathcal{X} \right\}$$

► **Remark:** showing that a kernel is positive definite is enough to get Φ and \mathcal{H} with the reproducing property “for free”

$$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

\mathcal{H}

\mathcal{H} der kernel.



$$\mathcal{H} = \left\{ \sum \alpha_i k(\cdot, x_i) \right\} \quad \|\mathbf{f}\|_{\mathcal{H}}^2 = \left\langle \sum \alpha_i k(\cdot, x_i), \sum \alpha_j k(\cdot, x_j) \right\rangle \\ = \sum \sum \alpha_i \alpha_j \langle k(\cdot, x_i), k(\cdot, x_j) \rangle \rightarrow$$

$$\forall f \in H_0, \quad \|f\|_H^2 = \sum_{i=1}^n \sum_{j=1}^n d_i d_j \langle k_{x_i}, k_{x_j} \rangle$$

$$= \sum_{i=1}^n \sum_{j=1}^n d_i d_j k(x_i, x_j)$$

$$= \sum_{i=1}^n \sum_{j=1}^n d_i d_j K_{ij} \leftarrow \text{Gram matrix}$$

$$\|f\|_H^2 = \mathbf{a}^T \mathbf{K} \mathbf{a} = \|\mathbf{a}\|_K^2$$

Matrizenb.z.



Example

- ▶ **Example:** polynomial kernel of degree 2:

$$k(x, y) = (x^\top y)^2.$$

- ▶ proved during the exercise session:

$$k(x, y) = \langle xx^\top, yy^\top \rangle_F,$$

thus k is positive definite

- ▶ **Question:** what is the RKHS?
- ▶ we know that \mathcal{H} contains all the functions

$$f(x) = \sum_i a_i k(x_i, x) = \sum_i a_i \langle x_i x_i^\top, x x^\top \rangle = \langle \sum_i a_i x_i x_i^\top, x x^\top \rangle$$



Example, ctd.

- ▶ spectral theorem: any symmetric matrix can be decomposed as $\sum_i a_i x_i x_i^\top$
- ▶ candidate RKHS: set a quadratic functions

$$f_S(x) = \langle S, xx^\top \rangle = x^\top S x,$$

with S symmetric matrix of size $d \times d$

- ▶ inner product on \mathcal{H} :

$$\langle f_S, f_{S'} \rangle = \langle S, S' \rangle_F.$$

- ▶ we can check that \mathcal{H} is a Hilbert space (isomorphic to $\mathcal{S}^{d \times d}$)
- ▶ finally, we check the reproducing property

5.3. More examples

Elementary properties

Proposition: Let $k_i : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a (potentially infinite) family of s.d. kernels. Then

- ▶ for any $\lambda_1, \dots, \lambda_p \geq 0$, the sum $\sum_{i=1}^p \lambda_i k_i$ is positive definite (can add)
- ▶ for any $a_1, \dots, a_p \in \mathbb{N}$, the product $k_1^{a_1} \cdots k_p^{a_p}$ is positive definite (can multiply)
- ▶ if it exists, the limit $k = \lim_{p \rightarrow +\infty} k_p$ is positive definite

Moreover, let \mathcal{X}_i be a sequence of sets and k_i positive kernels on each \mathcal{X}_i . Then

$$\left(k((x_1, \dots, x_p), (y_1, \dots, y_p)) := \prod_{i=1}^p k_i(x_i, y_i) \right) \quad k: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$$

$x = (x_1, x_2), y = (y_1, y_2)$

$$k((x_1, \dots, x_p), (y_1, \dots, y_p)) := \sum_{i=1}^p k_i(x_i, y_i) \quad \tilde{k}(x, y) = k(x_1, y_1) + \alpha k(x_2, y_2)$$

and

are positive definite kernels.

$\left(-\|x-y\|^2 \right)$
NOT a kernel :C

Taking the exponential

$$p! = 1 \cdot 2 \cdot \dots \cdot p \\ \geq 0$$

Theorem: if k is a positive definite kernel, then e^k as well.

► *Proof:* we write

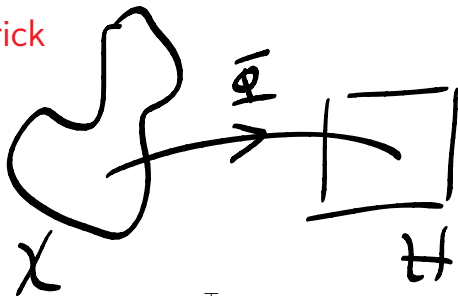
$$e^{k(x,y)} = \lim_{n \rightarrow +\infty} \sum_{p=0}^n \frac{k(x,y)^p}{p!},$$

then reason step by step.

- by the product property, $k(x,y)^p$ is a kernel for any $p \geq 0$
- as a positive linear combination of kernels, $\sum_{p=0}^n \frac{k(x,y)^p}{p!}$ is a kernel for all $n \geq 1$
- finally, e^k is a kernel as a limit of kernels. □

5.4. The kernel trick and applications

generally, $\dim \mathcal{H} = +\infty$
The kernel trick



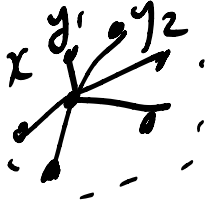
- ▶ input data $x_1, \dots, x_n \in \mathcal{X}$
- ▶ $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ kernel with associated RKHS \mathcal{H}
- ▶ we call $\Phi : \mathcal{X} \rightarrow \mathcal{H}$ the feature map
- ▶ **Idea:** imagine that our algorithm only depends on scalar products $x_i^\top x_j$
- ▶ then we can map the x_i to \mathcal{H} and replace the inner products by kernel evaluations, since

$$\langle \Phi(x_i), \Phi(x_j) \rangle = k(x_i, x_j).$$

- ▶ simple “trick” with many, many applications

$$x_i^\top x_j \longmapsto \langle \Phi(x_i), \Phi(x_j) \rangle$$

Example



- ▶ **Example:** computing distances
- ▶ suppose that our algo relies on distance computation
- ▶ that is, $\|x - y\|^2 = \langle x \cdot y, x \cdot y \rangle$
- ▶ we can write

1-NN in the feature space

$$\begin{aligned}\|\Phi(x) - \Phi(y)\|^2 &= \langle \Phi(x) - \Phi(y), \Phi(x) - \Phi(y) \rangle \\ &= \langle \Phi(x), \Phi(x) \rangle - 2\langle \Phi(x), \Phi(y) \rangle + \langle \Phi(y), \Phi(y) \rangle \\ \|\Phi(x) - \Phi(y)\|^2 &= k(x, x) - 2k(x, y) + k(y, y).\end{aligned}$$

$$\begin{aligned}d(x, y_i)^2 &= k(x, x) \\ &\quad - 2k(x, y_i) \\ &\quad + k(y_i, y_i)\end{aligned}$$

- ▶ in other words,

dist. in feature space

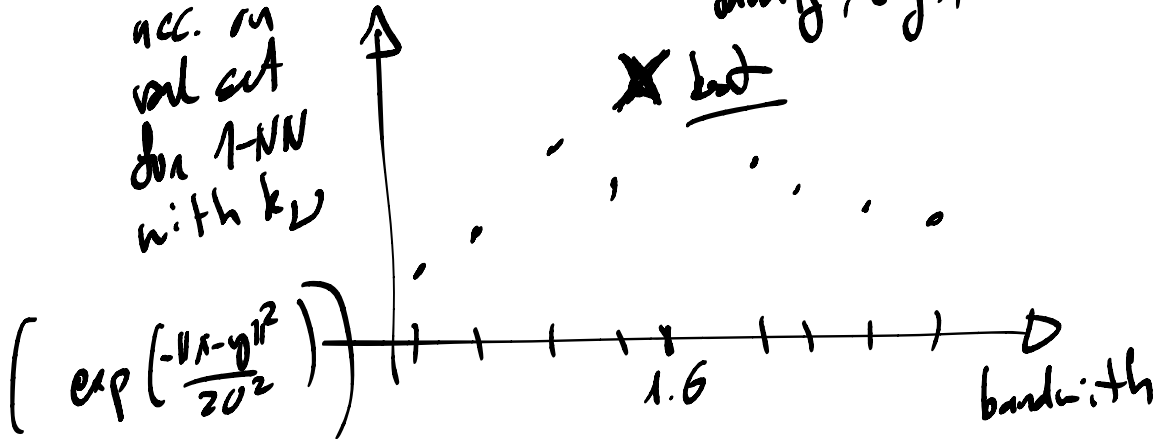
- ▶ as promised, we do not need to know Φ !

$$d_{\mathcal{H}}(x, y) = \sqrt{k(x, x) - 2k(x, y) + k(y, y)}$$

Q: How to choose the kernel?

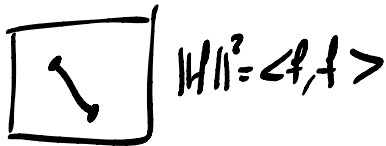
"hyperparameters of kernel family, e.g., Gaussian

acc. on
val set
for 1-NN
with k_L



5.5. The representer theorem

Motivation



- ▶ let us imagine that we take \mathcal{H} as hypothesis class
- ▶ starting from regularized ERM, our optimization problem will look like

$$\arg \min_{f \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) + \lambda \|f\|^2 \right\}. \quad (*)$$

- ▶ we penalize by the norm because it is an indicator of the *smoothness* of f
- ▶ **Why?** Cauchy-Schwarz + exercise:

$$|f(x) - f(y)| = |\langle f, k_x - k_y \rangle| \leq \|f\| \cdot \|k_x - k_y\| = \|f\| \cdot d_{\mathcal{H}}(x, y).$$

- ▶ Eq. (*) is a complicate problem, potentially *infinite-dimensional*
- ▶ **Question:** how to solve it in practice?