

Recall Reproducing Kernel Hilbert Space.

Idea for linear regression.  $W = X^T \alpha$   $\{x_1 \dots x_n\}$  are  $n$  data.  
 $f(x) = \langle w, x \rangle$

$w = \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n$   
 $f$  is linear combination of  $f(x_i, \cdot)$

For kernel space: final function  $f = \alpha_1 k_{x_1} + \alpha_2 k_{x_2} + \dots + \alpha_n k_{x_n}$

Kernel:  $k(x, y) = \langle k_x, k_y \rangle$   $f(x) = \langle f, k_x \rangle$ ,  $k_x$  is the mapping from  $f$  to  $f(x)$

$f(y) = \langle f, k_y \rangle = \langle \alpha_1 k_{x_1} + \alpha_2 k_{x_2} + \dots + \alpha_n k_{x_n}, k_y \rangle$

$= \alpha_1 \langle k_{x_1}, k_y \rangle + \alpha_2 \langle k_{x_2}, k_y \rangle + \dots + \alpha_n \langle k_{x_n}, k_y \rangle$

$= \alpha_1 k(x_1, y) + \alpha_2 k(x_2, y) + \dots + \alpha_n k(x_n, y)$

$f(\cdot) = [\alpha_i k(x_i, \cdot)]$   
 as basis function

$\mathcal{H}$  Hilbert Space (Vector space  $x_1, \dots, x_n \in V$   $\alpha_1 x_1 + \dots + \alpha_n x_n \in V$  + inner product)

$\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle = \alpha_1 \langle f_1, g \rangle + \alpha_2 \langle f_2, g \rangle$   
 $\langle f, g \rangle = \langle g, f \rangle$   
 $\langle f, f \rangle_H = \|f\|_H^2 \geq 0$

$\langle \cdot, \cdot \rangle$  is a function  
 $\langle f, g \rangle \rightarrow \mathbb{R}$   
 $f$  and  $g$  lies in a vector space

why  $K = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_n, x_1) & k(x_n, x_2) & \dots & k(x_n, x_n) \end{pmatrix}$  is P.S.D.?

Example Functions  $(f+g)(x) = f(x) + g(x) \quad \forall x \in \mathbb{R}^d$

examples  $\rightarrow \langle f, g \rangle_L = \int f(x)g(x) dx$   $\langle f, f \rangle = \|f\|_{L^2}^2 = \int f^2 dx$   
 $\rightarrow \langle f, g \rangle = \int f(x)g(x) + f'(x)g'(x) dx$

for a Reproducing kernel Hilbert space

$$- \sup_x \|k_x\| \leq B$$

$$f(x) \leq \|k_x\| \cdot \|f\|_H$$

Reproducing  $k_x$ : is a bounded mapping  $f \rightarrow f(x)$   
 $\langle f, k_x \rangle = f(x)$

- The function space is a Hilbert space

$$\Rightarrow f = \sum \alpha_i k_{x_i} \text{ means}$$

$$f(y) = \sum \alpha_i \langle k_{x_i}, k_y \rangle = \sum \alpha_i k(x_i, y)$$

- kernel:  $k(x, y) = \langle k_x, k_y \rangle$

??

Not Required If we have an innerproduct, can we write down the kernel?  $\Rightarrow$  Not always, in many cases,  $\|k_x\|$  is  $\infty$

In some of cases, we can do this.

Examples  $\langle f, g \rangle_H = \int_0^1 f'(x) g'(x) dx$ ,  $f, g: [0, 1] \rightarrow \mathbb{R}$   
 $f(0) = 0$ .

Claim  $k(x, z) = \min\{x, z\}$

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

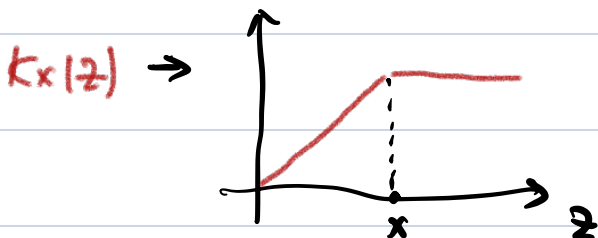
Check  $\langle f, k_x \rangle = f(x)$

$$\int_0^1 f'(z) k_x'(z) dz$$

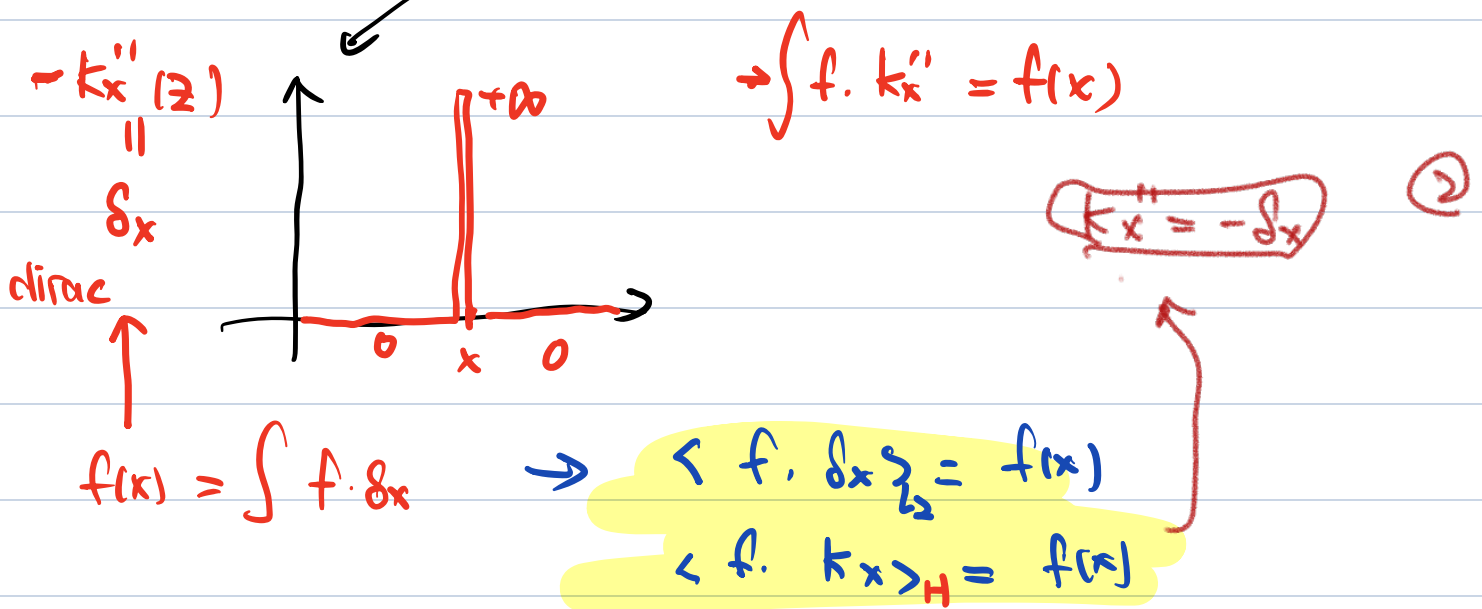
Integral by parts!!

$$-\int_0^1 f(z) k_x''(z) dz$$

$\underset{f(x)}{\parallel}$



$$\Delta \int f'(z) k_x'(z) dz = \int_0^x f'(z) dz = f(x)$$



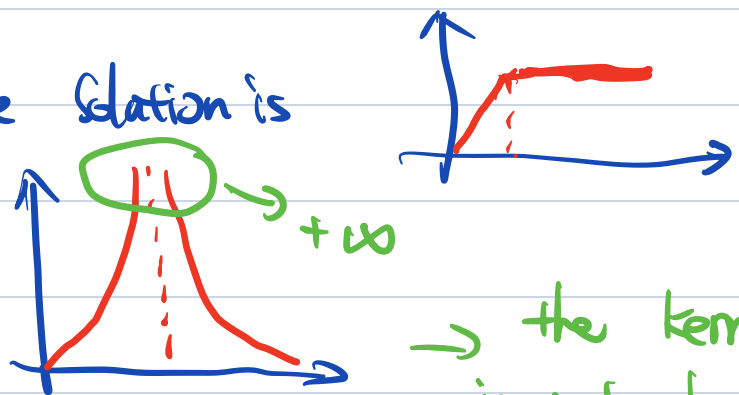
## Smoothness of Reproducing Kernel Hilbert Space (Not Required)

$\langle f, g \rangle_H = \int f' g'$   $\rightarrow k_x'' = -\delta_x$

Kernel is the  
(Green's Function)

in 1 dimension, the relation is

but in 2 dimension



$\rightarrow$  the kernel is not bounded  
 $\Rightarrow$  it's not Reprody!

$\langle f, g \rangle_H = \int f'' g'' \rightarrow k_x'''' = \delta_x$

In 1-3 dimension it is the bounded, in 4-dimen it's not!

Informal: In  $R^d$  dimension,  $f^{(d/2)}$  should exist

## Representer Theorem.

If one minimizing

$$f^* = \operatorname{argmin}_{f \in H} L(f(x_1), f(x_2), \dots, f(x_n)) + \Omega(\|f\|_H^2)$$

$\rightarrow$  function space (may be infinite dimensional)

Then.  $f^* = \sum_{i=1}^n \alpha_i k_{x_i}$

or  $f^*(\cdot) = \sum_{i=1}^n k(x_i, \cdot)$

"Dual Solution only depend on the number of data"

- Not just holds for RKHS. It holds for a lot of hypothesis space, even Neural Networks.

# § Rademacher Complexity of a Kernel Class

$$\hat{R}_{S_n}(\{f \in \mathcal{F} \mid \|f\|_H \leq M\}) \leq \frac{M}{n} \sqrt{\sum_{i=1}^n k(x_i, x_i)}$$

$$= \frac{M}{\sqrt{n}} \sqrt{\frac{1}{n} \sum_{i=1}^n k(x_i, x_i)}$$

Intrinsic assumption  
 $\frac{1}{n} \sum_{i=1}^n k(x_i, x_i)$  uniformly bounded

↳ This is the trace of Kernel Matrix.

Proof. LHS =  $\frac{1}{n} \mathbb{E}_\sigma \left[ \sup_{\{f \in \mathcal{F} \mid \|f\|_H \leq M\}} \sum_{i=1}^n \sigma_i f(x_i) \right]$   $x_1 \dots x_n$  is data

$$= \frac{1}{n} \mathbb{E}_\sigma \left[ \sup_{\{f \in \mathcal{F} \mid \|f\|_H \leq M\}} \sum_{i=1}^n \sigma_i \langle f, k_{x_i} \rangle \right]$$

$$= \frac{1}{n} \mathbb{E}_\sigma \left[ \sup_{\{f \in \mathcal{F} \mid \|f\|_H \leq M\}} \langle f, \sum_{i=1}^n \sigma_i k_{x_i} \rangle \right]$$

$$\leq \frac{1}{n} M \mathbb{E}_\sigma \left\| \sum_{i=1}^n \sigma_i k_{x_i} \right\|$$

(This is because  $\|f\|_H \leq M$ )

$$= \frac{1}{n} M \mathbb{E}_\sigma \sqrt{\langle \sum_{i=1}^n \sigma_i k_{x_i}, \sum_{i=1}^n \sigma_i k_{x_i} \rangle}$$

( $\|f\| = \sqrt{\langle f, f \rangle}$ )

$$\leq \frac{1}{n} M \sqrt{\mathbb{E}_\sigma \langle \sum_{i=1}^n \sigma_i k_{x_i}, \sum_{i=1}^n \sigma_i k_{x_i} \rangle}$$

( $\mathcal{F}$  is a convex Jensen's Inequality)

$$= \frac{1}{n} M \sqrt{\sum_{i=1}^n \sum_{j=1}^n \sigma_i \sigma_j \langle k_{x_i}, k_{x_j} \rangle = k(x_i, x_j)}$$

$\mathbb{E} \sigma_i \sigma_j = 0$  if  $i \neq j$

$$= \frac{1}{n} M \sqrt{\sum_{i=1}^n k(x_i, x_i)}$$

$f(\cdot)$  is concave

$$c_1 f(x) + c_2 f(y)$$

$$\leq f(c_1 x + c_2 y)$$

Matches "Least square's variance

is the trace of Covariance

matrix"

# § Spectral View of RKHS (Informal)

(HWY, Sobolev)

The Complexity is the trace(K)

$$= \sum_{i=1}^{\infty} \lambda_j(K) < \infty$$

The eigenvalues should be decay fast enough !!

Example. Shift-Invariant kernel

⇒ eigenvectors: Fourier basis !!

"eigenvalue" will mean Fourier coef

⇒ "rough understanding"

If a function lies in a RKHS, fast enough then the Fourier coef should decay

---

$$\alpha = (X X^T + \lambda I)^{-1} X^T y \quad (\text{Ridge Regression})$$

If  $X X^T = \sum \lambda_i u_i u_i^T$  is the eigendecomposition.

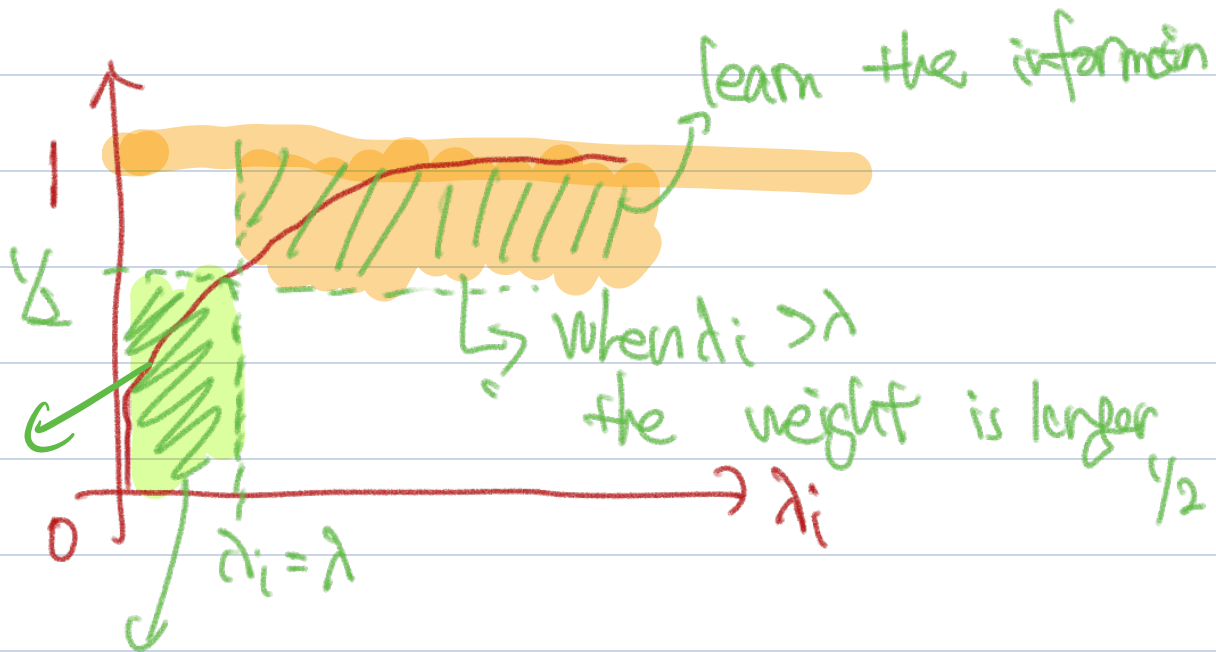
$$\langle x, \alpha \rangle = \sum_{i=1}^n \frac{\lambda_i}{\lambda_i + \lambda} \langle y, u_i \rangle$$

weight

↳ how much you have on  $u_i$

$$\frac{\lambda_i}{d_i + \lambda}$$

negative



if  $d_i < \lambda$ , weight is smaller than  $\frac{1}{2}$

---

We only learn the part  $d_i > \lambda$

$$\Rightarrow \sum_{d_i > \lambda} d_i < \infty$$

"mis-specification"  $f \notin H$

$A: x \rightarrow Ay$  Consider matrix as a linear mapping.

$$\underline{g} = Kf$$



$$K = \begin{pmatrix} k(x_1, x_1) & & \\ & \ddots & \\ & & k(x_n, x_n) \end{pmatrix}$$
$$f = \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix} \quad g = \begin{bmatrix} g(x_1) \\ \vdots \\ g(x_n) \end{bmatrix}$$

$$g(x_j) = \frac{1}{n} \sum_{i=1}^n k(x_j, x_i) f(x_i)$$



$$g(x) = \int k(x, y) f(y) dy$$

in population. "kernel matrix" as the integral operator

$$f \rightarrow \underline{g} = \int k(x, y) f(y) dy .$$