

# Pre-asymptotic Optimality of Cross-validation in Scattered Data Approximation

Felix Bartel   Ralf Hielscher

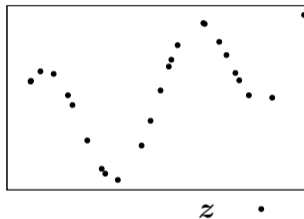
International Conference on Computational Harmonic Analysis

14th of September 2021

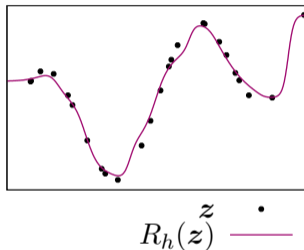


Mathematik!  
TU Chemnitz

- data  $\mathbf{z} = (\mathbf{x}_i, y_i)_{i=1}^n \in (\Omega \times Y)^n$
- $(\mathbf{x}_i, y_i)$  distributed according to  $\rho$



- data  $\mathbf{z} = (\mathbf{x}_i, y_i)_{i=1}^n \in (\Omega \times Y)^n$
- $(\mathbf{x}_i, y_i)$  distributed according to  $\rho$
- reconstruction algorithm  
 $R_h: (\Omega \times Y)^n \rightarrow Y^\Omega$



**Given:** data  $z$ , reconstruction model  $R_h(z)$

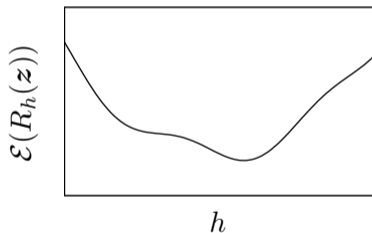
**Question:** How big is the  $L_2$ -reconstruction error of  $R_h(z)$ ?

$$\mathcal{E}(R_h(z)) = \int_{\Omega \times Y} |(R_h(z))(\mathbf{x}) - y|^2 \, d\rho(\mathbf{x}, y)$$

## Motivation 2: model selection

**Given:** data  $z$ , hypothesis space  $\{R_h(z) : h \in H\}$

**Question:** Which reconstruction model  $R_h(z)$  to choose?



Purly data-driven method to approximate the risk

- validation set



Purly data-driven method to approximate the risk

- validation set
- “Probably the simplest and most widely used method for estimation prediction error is cross-validation.” [Hastie 2001]

test	train		
train	test	train	
train		test	train
train			test train
train			test

- 1 calculate the reconstructions  $R_h(\mathbf{z}_{-i})$  from the function values  $\mathbf{z}_{-i} = (z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n)$
- 2 evaluate the residual of  $R_h(\mathbf{z}_{-i})$  in the  $i$ -th node  $\mathbf{x}_i$

$$|R_h(\mathbf{z}_{-i})(\mathbf{x}_i) - y_i|$$

- 3 calculate the mean value with respect to all nodes

$$\text{CV}(\mathbf{z}, \lambda) := \frac{1}{n} \sum_{i=1}^n |R_h(\mathbf{z}_{-i})(\mathbf{x}_i) - y_i|^2$$



For PLSE with  $1/n \text{trace } \mathbf{H}_\lambda < 1$  they showed

$$\begin{aligned} & \frac{|\mathbb{E}(\mathcal{E}(R_\lambda(\mathbf{Z})) - \mathbb{E}(\text{CV}(\mathbf{Z}, \lambda)))|}{\mathbb{E}(\text{CV}(\mathbf{Z}, \lambda))} \\ & \leq \left( \frac{2}{n} \text{trace } \mathbf{H}_\lambda + \frac{(\text{trace } \mathbf{H}_\lambda)^2}{\text{trace } \mathbf{H}_\lambda^2} \right) \frac{1}{(1 - 1/n \text{trace } \mathbf{H}_\lambda)^2}. \end{aligned}$$

Let

$$h^+ = \arg \min_h \mathcal{E}(R_h(\mathbf{z})) \quad \text{and} \quad h^* = \arg \min_h \text{CV}(\mathbf{z}, h).$$

Then under certain assumptions

$$\frac{\mathcal{E}(R_{h^*}(\mathbf{z}))}{\mathcal{E}(R_{h^+}(\mathbf{z}))} \rightarrow 1 \quad \text{for} \quad n \rightarrow \infty.$$

## Theorem (Hoeffding '63)

Let

- $Z_1, \dots, Z_n$  independent rv's with values in  $[0, 1]$  and
- $m = \mathbb{E} \left\{ \frac{1}{n} \sum_{i=1}^n Z_i \right\}$ .

Then for  $\varepsilon > 0$

$$\mathbb{P} \left\{ \left| \frac{1}{n} \sum_{i=1}^n Z_i - m \right| > \varepsilon \right\} \leq 2 \exp(-2n\varepsilon^2).$$

A function  $f: \Omega^n \rightarrow \mathbb{R}$  is  **$c$ -bounded** on  $\Xi \subset \Omega^n$  for  $c = (c_1, \dots, c_n)$ , iff

$$|f(z_1, \dots, z_n) - f(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n)| \leq c_i$$

for all  $(z_1, \dots, z_n), (z'_1, \dots, z'_n) \in \Xi$ , and  $1 \leq i \leq n$ .

## Theorem (McDiarmid '89)

Let

- $Z_1, \dots, Z_n$  be independent rv's with values in  $\Omega$ ,
- $f: \Omega^n \rightarrow \mathbb{R}$  be  $c$ -bounded on  $\Omega^n$ , and
- $m = \mathbb{E} \{f(Z_1, \dots, Z_n)\}$ .

Then for  $\varepsilon > 0$

$$\mathbb{P} \{|f(Z_1, \dots, Z_n) - m| > \varepsilon\} \leq 2 \exp\left(-\frac{2\varepsilon^2}{\|c\|_2^2}\right).$$

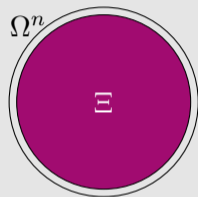
## Theorem (Extension of McDiarmid by Combes '15)

Let

- $Z_1, \dots, Z_n$  be independent rv's with values in  $\Omega$ ,
- $f: \Omega^n \rightarrow \mathbb{R}$  be  $c$ -bounded on  $\Xi \subset \Omega^n$ , and
- $m = \mathbb{E} \{f(Z_1, \dots, Z_n) | (Z_1, \dots, Z_n) \in \Xi\}$ , and
- $\gamma = 1 - \mathbb{P}\{(Z_1, \dots, Z_n) \in \Xi\}$ .

Then for  $\varepsilon > \gamma \|c\|_1$

$$\mathbb{P} \{|f(Z_1, \dots, Z_n) - m| > \varepsilon\} \leq 2\gamma + 2 \exp \left( -\frac{2(\varepsilon - \gamma \|c\|_1)^2}{\|c\|_2^2} \right).$$



Define  $\Xi = \Xi(C_1, C_2)$  as the set of  $\mathbf{z} = (z_1, \dots, z_n)$  fulfilling

- 1 a uniform bound on the reconstruction error

$$\|R_h(\mathbf{z}_{-i}) - f\|_\infty \leq C_1$$

for  $1 \leq i \leq n$ .

Define  $\Xi = \Xi(C_1, C_2)$  as the set of  $\mathbf{z} = (z_1, \dots, z_n)$  fulfilling

- 1 a uniform bound on the reconstruction error

$$\|R_h(\mathbf{z}_{-i}) - f\|_\infty \leq C_1$$

- 2  $\mathbf{z}_{-i} \mapsto R_h(\mathbf{z}_{-i})(\mathbf{x})$   $c$ -bounded for all  $\mathbf{x} \in \Omega$  on  $\Xi$  with  $c = \mathbb{1}C_2$

for  $1 \leq i \leq n$ .



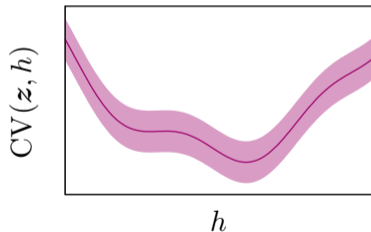
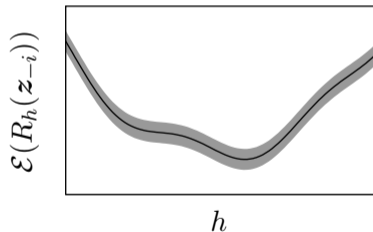
Risk functional  $z \mapsto \mathcal{E}(R_h(z_{-i}))$  is  $c$ -bounded on  $\Xi$  with  $c = 4C_1^2 \mathbf{1}$ .

Cross-validation score  $z \mapsto \text{CV}(z, h)$  is  $c$ -bounded on  $\Xi$  with  $c = 5C_1^2 \mathbf{1}$ .

Risk functional  $z \mapsto \mathcal{E}(R_h(z_{-i}))$  is  $c$ -bounded on  $\Xi$  with  $c = 4C_1^2 \mathbf{1}$ .

Cross-validation score  $z \mapsto \text{CV}(z, h)$  is  $c$ -bounded on  $\Xi$  with  $c = 5C_1^2 \mathbf{1}$ .

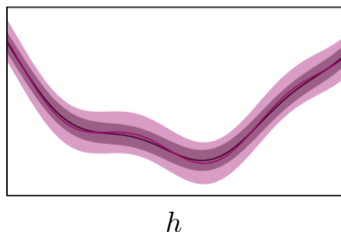
Now we apply Combes extension of McDiarmid.



## Lemma

For  $\mathbf{Z}'$  representing  $n - 1$  samples and  $\mathbf{Z}$  representing  $n$  samples element-wise distributed according to  $\rho$  we have

$$\mathbb{E}_{\mathbf{Z}'} \{ \mathcal{E}(R_h(\mathbf{Z}')) \} = \mathbb{E}_{\mathbf{Z}} \{ \text{CV}(\mathbf{Z}, h) \}.$$



$\mathcal{E}(R_h(\mathbf{z}_{-i}))$   
 $\text{CV}(\mathbf{z}, h)$

## Theorem (B., Hielscher '21)

Let  $n \geq 3$ ,

- $\mathbf{Z}$  be element-wise distributed according to  $\rho$  with values in  $(\Omega \times Y)^n$ , and
- $\gamma = 1 - \mathbb{P}\{\mathbf{Z} \in \Xi\}$ .

Then for  $\delta > 0$  we have with probability larger than  $1 - 2(\gamma + \delta)$

$$\begin{aligned} & |\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))| \\ & \leq \max \left\{ 2\gamma(M + \|f\|_\infty)^2, 12\sqrt{n}C_1^2(\sqrt{2n\gamma} + \sqrt{-\log \delta}) \right\} \end{aligned}$$

# Shepard's model

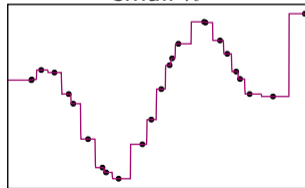
Shepard's model for  $\Omega = \mathbb{T}$ ,  $Y = \mathbb{R}$ , and  $\mathbf{z} = (x_i, f(x_i))_{i=1}^n$

$$R_h(\mathbf{z}) = \frac{\sum_{i=1}^m K_h(\cdot, x_i) f(x_i)}{\sum_{i=1}^m K_h(\cdot, x_i)}$$

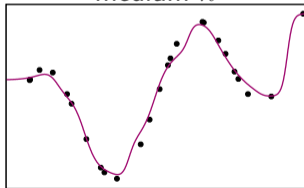
with

$$K_h(x_1, x_2) = \max\{0, 1 - h|x_1 - x_2|\}$$

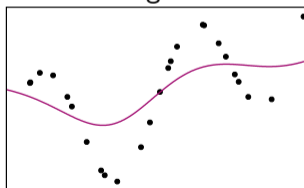
small  $h$



medium  $h$



big  $h$



We need estimates for

- $\|(R_h(\cdot))(\cdot)\|_\infty$ ,
- $\|R_h(\mathbf{z}) - f\|_\infty$  for  $\mathbf{z} \in \Xi$ , and
- $\gamma = 1 - \mathbb{P}\{\mathbf{Z} \in \Xi\}$ .

We need estimates for

- $\|(R_h(\cdot))(\cdot)\|_\infty \leq \|f\|_\infty$ ,
- $\|R_h(\mathbf{z}) - f\|_\infty$  for  $\mathbf{z} \in \Xi$ , and
- $\gamma = 1 - \mathbb{P}\{\mathbf{Z} \in \Xi\}$ .

We need estimates for

- $\|(R_h(\cdot))(\cdot)\|_\infty \leq \|f\|_\infty$ ,
- $\|R_h(z) - f\|_\infty$  for  $z \in \Xi$ , and
- $\gamma = 1 - \mathbb{P}\{\mathbf{Z} \in \Xi\}$ .

## Theorem (B., Hielscher '21)

Let  $x_1, \dots, x_n$  be uniformly random,  $K_h(x, \cdot)$  be supported on  $[x - 1/h, x + 1/h]$ , and  $f$  be Lipschitz with constant  $L$ . Then  $\|R_h(z) - f\|_\infty \leq L/h$  with probability

$$1 - \gamma \geq \sum_{k=0}^{\lfloor h \rfloor} (-1)^{k+1} \binom{n}{k} \left(1 - \frac{k}{2h}\right)^{n-1} \approx 1 - \exp\left(-n \exp\left(-\frac{n}{h}\right)\right).$$



## Theorem (B., Hielscher '21)

Let

- $f$  be Lipschitz with constant  $L$ ,
- $\mathbf{z} = (x_i, f(x_i))_{i=1}^n$  be samples with  $x_1, \dots, x_n \in \mathbb{T}$  uniformly distributed,
- $K_h(x, \cdot)$  be supported on  $[x - 1/h, x + 1/h]$ , and
- $\gamma \approx \exp(-n \exp(-n/h))$ .

Then we have for  $\delta > 0$  with probability larger than  $1 - 2(\gamma + \delta)$

$$\begin{aligned} & |\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))| \\ & \leq \max \left\{ 4\gamma \|f\|_{\infty}^2, \frac{12\sqrt{n}L^2}{h^2} (\sqrt{2n\gamma} + \sqrt{-\log \delta}) \right\} \end{aligned}$$

Using

- binary kernels  $K_h : \Omega \times \Omega \rightarrow \{0, 1\}$
- $\Xi$  all possible samples

they showed

$$\gamma = 0, \quad C_1 = 2\|f\|_\infty, \quad C_2 \sim \frac{1}{n}$$

and obtained with probability  $1 - 2\delta$

$$|\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))| \lesssim n^{-1/2}.$$

Using

- binary kernels  $K_h : \Omega \times \Omega \rightarrow \{0, 1\}$
- $\Xi$  all possible samples

they showed

$$\gamma = 0, \quad C_1 = 2\|f\|_\infty, \quad C_2 \sim \frac{1}{n}$$

and obtained with probability  $1 - 2\delta$

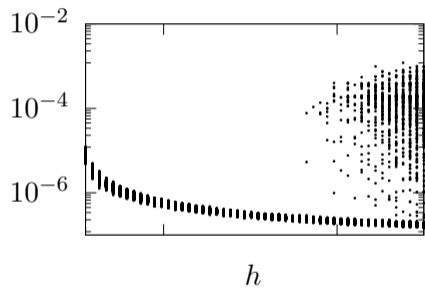
$$|\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))| \lesssim n^{-1/2}.$$

Assuming  $h \sim n$ , we obtained with probability  $1 - 2(\delta + \exp(-n \exp(-n/h)))$

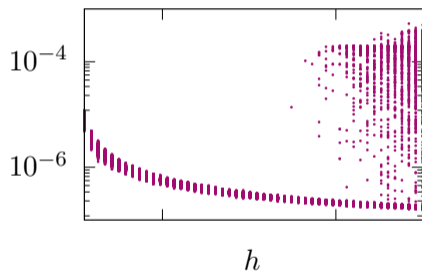
$$|\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))| \lesssim n^{-3/2}.$$

# Numerical example

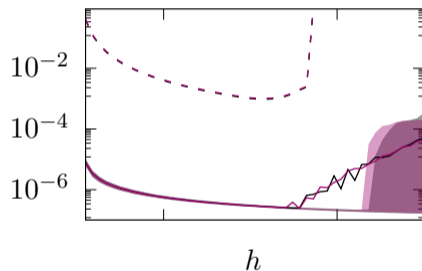
$$\mathcal{E}(R_h(\mathbf{z}_{-1}))$$



$$CV(\mathbf{z}, h)$$



$CV(z, h)$  and  $\mathcal{E}(R_h(z_{-1}))$



$|CV(z, h) - \mathcal{E}(R_h(z_{-1}))|$

