

# Concentration Inequalities for Cross-validation in Scattered Data Approximation

Felix Bartel

Ralf Hielscher

Choosing models from a hypothesis space is a frequent task in approximation theory and inverse problems. Cross-validation is a classical tool in the learner's repertoire to compare the goodness of fit for different reconstruction models. Much work has been dedicated to computing this quantity in a fast manner but tackling its theoretical properties occurs to be difficult. So far, most optimality results are stated in an asymptotic fashion. In this paper we propose a concentration inequality on the difference of cross-validation score and the risk functional with respect to the squared error. This gives a pre-asymptotic bound which holds with high probability. For the assumptions we rely on bounds on the uniform error of the model which allow for a broadly applicable framework.

We support our claims by applying this machinery to Shepard's model, where we are able to determine precise constants of the concentration inequality. Numerical experiments in combination with fast algorithms indicate the applicability of our results.

*Key words.* cross-validation, scattered data approximation, model selection, parameter choice, concentration inequalities

## 1 Introduction

The general problem in scattered data approximation is the reconstruction of a function  $f: \Omega \rightarrow Y$  based on discrete samples  $\mathbf{z} = (z_i)_{i=1}^n = (\mathbf{x}_i, f(\mathbf{x}_i))_{i=1}^n \in (\Omega \times Y)^n$ . The nodes  $\mathbf{x}_i$  are independent and identically distributed according to  $\rho$  on  $\Omega$ . Extensive work has been done to develop reconstruction algorithms  $R_h: (\Omega \times Y)^n \rightarrow Y^\Omega$  which propose candidates for the approximation. Here,  $h$  resembles one of the various methods with possible parameters. Using multiple reconstruction algorithms  $R_h$ ,  $h \in H$  we end up with a hypothesis space  $\{R_h(\mathbf{z}) : h \in H\} \subset Y^\Omega$ . Even given a precise application, it remains difficult to choose reconstruction algorithms  $R_h$ ,  $h \in H$  which yields the best reconstruction  $R_h(\mathbf{z})$  of  $f$ .

In order to find an optimal  $R_h(\mathbf{z})$ ,  $h \in H$ , we would like to rank the reconstructions with respect to their goodness of fit. This is quantified by the *risk functional*. In this paper we consider the risk functional with respect to the squared loss

$$\mathcal{E}(R_h(\mathbf{z})) = \int_{\Omega} |(R_h(\mathbf{z}))(\mathbf{x}) - f(\mathbf{x})|^2 \, d\rho(\mathbf{x}). \quad (1.1)$$

Even though this is theoretically appealing we would need to know the underlying distribution  $\rho$  and the function  $f$  to compute this quantity. Since this is not the case, we seek for an alternative which only relies on the given data. The concept which struck our attention is called *cross-validation*, was initially introduced in [14], and has been widely used since then, cf. [41, 7, 31, 35, 9]. The basic idea consists of subdividing the data into a training set and a validation set for estimating the error. Doing this multiple times we obtain a reasonable estimator for the risk functional. A special case is where the partitionings seclude single nodes, then the training sets become  $\mathbf{z}_{-i} := (z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n)$  and the validation sets  $\{z_i\}$ . This leads to the so called *leave-one-out cross-validation score*

$$CV(\mathbf{z}, h) = \frac{1}{n} \sum_{i=1}^n |(R_h(\mathbf{z}_{-i}))(\mathbf{x}_i) - f(\mathbf{x}_i)|^2. \quad (1.2)$$

An immediate drawback is given by the numerical complexity of computing the  $n$  approximations  $R_h(\mathbf{z}_{-i})$ . However, this is circumvented in many cases with ideas including Monte Carlo approximations [10], matrix decomposition methods [43, 38], Krylow space methods [29], or Fourier analysis [2].

One is interested in a theoretical foundation of the cross-validation score. By the Bakushinskii veto, cf. [1], we know that there exists a realization of the samples, such that purely data-driven regularization methods have no guarantee for a good approximation without incorporating further information. One still has propositions about the goodness of the cross-validation score in asymptotic cases, cf. [27, 17, 28, 16], on average, cf. [14, 4, 5], or by restriction of noise, cf. [21, 22].

In this paper we bound the difference of cross-validation and risk pre-asymptotically, which supports the choice of cross-validation for model selection. To circumvent the Bakushinskii veto our results will hold with high probability as it is common in learning theory. We use mild assumptions on the uniform error of the reconstruction algorithm, which allow for a broadly applicable framework. These bounds improve on the results from [17, Chapter 8] in a more general setting. Other pre-asymptotic results can be found in [19, 24], where the algorithmic stability, a variance-like concept, of the cross-validation score is examined.

As for the structure of this paper, in Section 2 we repeat on an extension of McDiarmid's concentration inequality, as it will be of importance later on. The main part is Section 3, where we present our general framework. Therefore, we prove in Theorems 3.4 and 3.6 concentration inequalities for the risk functional (1.1) and the cross-validation score (1.2) with respect to the data  $\mathbf{z}$ . These concentration inequalities are used to surround the expected values of the risk functional  $\mathcal{E}(R_h(\cdot))$  and the cross-validation score

$\text{CV}(\cdot, h)$  by narrow intervals in which nearly all realizations of these quantities lie. In Lemma 3.7 we show that the expected values of  $\mathcal{E}(R_h(\cdot))$  and  $\text{CV}(\cdot, h)$  coincide. Eventually, this leads us to our main result in Theorem 3.8, which bounds the difference of risk functional and cross-validation score with high probability and, therefore, justifies the usage of cross-validation for choosing models and parameters. To exemplify the applicability of our results and reason for the stated conditions to make sense we apply the framework to Shepard's model in Section 4. As before, we bound the difference of cross-validation score and risk with high probability, now with precise constants in Theorem 4.4. We confirm our results with numerical experiments.

## 2 McDiarmid's concentration inequality

Since it will be of fundamental importance, we dedicate this section to an extension of McDiarmid's concentration inequality. We consider random variables  $\mathbf{X} = (X_1, \dots, X_n)$  on a probability space  $(\Omega^n, \mathcal{A}, \mathbb{P})$ . As usual we denote with

$$\mathbb{P}\{A|B\} = \frac{\mathbb{P}\{A \cap B\}}{\mathbb{P}\{B\}} \quad \text{and} \quad \mathbb{E}\{X|B\} = \frac{\mathbb{E}\{\mathbf{1}_B X\}}{\mathbb{P}\{B\}}$$

the conditional probability and expected value, respectively. To state McDiarmid's theorem we need the following concept.

**Definition 2.1.** *A function  $f: \Omega^n \rightarrow \mathbb{R}$  is said to be  $\mathbf{c}$ -bounded on  $\Xi \subset \Omega^n$  for  $\mathbf{c} = (c_1, \dots, c_n) \in [0, \infty)^n$  if and only if*

$$|f(\mathbf{x}) - f(\mathbf{x}')| \leq d_{\mathbf{c}}(\mathbf{x}, \mathbf{x}')$$

for all  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{x}' = (x'_1, \dots, x'_n) \in \Xi$  where the distance  $d_{\mathbf{c}}$  is defined by

$$d_{\mathbf{c}}(\mathbf{x}, \mathbf{x}') = \sum_{i: x_i \neq x'_i} c_i.$$

Note, that a function is  $\mathbf{c}$ -bounded if changing a single variable  $x_i$ ,  $1 \leq i \leq n$  changes  $f(\mathbf{x})$  only by  $c_i$ , i.e.,

$$|f(x_1, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \leq c_i$$

for all  $(x_1, \dots, x_n), (x'_1, \dots, x'_n) \in \Xi$ .

McDiarmid's inequality, cf. [30], is a generalization of Hoeffding's inequality. We will not state the original theorem, but an extension from [8].

**Theorem 2.2.** *Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a vector of independent random variables taking values in a set  $\Omega$ . Furthermore, let  $f: \Omega^n \rightarrow \mathbb{R}$  be  $\mathbf{c}$ -bounded on  $\Xi \subset \Omega^n$ ,  $m = \mathbb{E}\{f(\mathbf{X})|\mathbf{X} \in \Xi\}$  be the expected value of  $f$  restricted to  $\Xi$ , and  $\gamma = 1 - \mathbb{P}\{\mathbf{X} \in \Xi\}$  the probability of  $\mathbf{X}$  not being in  $\Xi$ .*

*Then we have for  $\varepsilon > \gamma\|\mathbf{c}\|_1$  the concentration of  $f(\mathbf{X})$  around its expected value*

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

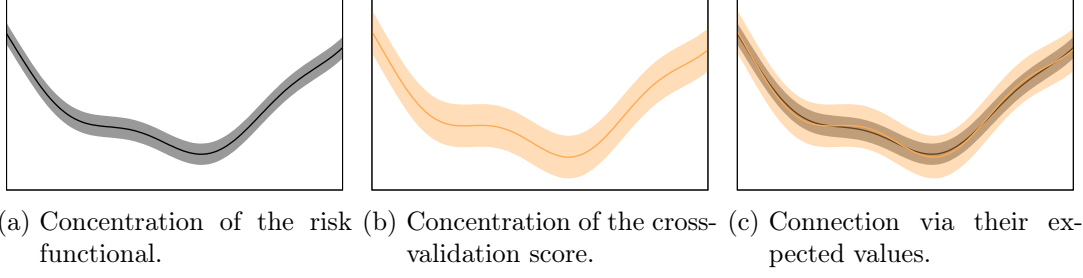


Figure 3.1: Intuition of Theorems 3.4, 3.6, and 3.8

### 3 General framework

Throughout this section we consider an arbitrary domain  $\Omega$  equipped with some probability measure  $\rho$  and a function  $f: \Omega \rightarrow Y$  which we want to approximate from a finite sampling  $\mathbf{z} = (x_i, f(x_i))_{i=1}^n$ . We consider the sampling  $\mathbf{z} \in (\Omega \times Y)^n$  as a realization of the random variable  $\mathbf{Z} = (\mathbf{X}_i, f(\mathbf{X}_i))_{i=1}^n$  with  $\mathbf{X}_i$  being independently and identically  $\rho$ -distributed random variables with values in  $\Omega$ . This includes the generality of data-driven approximation methods.

The goal of this section is to relate, for an arbitrary approximation operator  $R_h: (\Omega \times Y)^n \rightarrow Y^\Omega$ , the risk functional (1.1) and the cross-validation score (1.2). This is done in three steps: First we prove concentration inequalities for the risk functional and cross-validation score in Theorem 3.4 and 3.6, respectively. For every reconstruction algorithm  $R_h$ , this restricts their values to an interval around their expected values with high probability as depicted in Figure 3.1 (a) and (b). In Lemma 3.7 we state the connection of these two expected values. These three facts allow us to overlap the two concentrations, cf. Figure 3.1 (c), and lead to Theorem 3.8 which is a concentration inequality for the difference of risk functional and cross-validation score.

Dealing with reconstruction algorithms  $R_h: (\Omega \times Y)^n \rightarrow Y^\Omega$  in scattered data approximation settings, there may exist possible realizations  $\mathbf{z} \in (\Omega \times Y)^n$  of the samples such that we cannot bound the error of the approximation in a small manner. An example for that would be polynomial interpolation where all nodes  $\mathbf{x}_i$  coincide. To handle these outliers we define a subset of all samples excluding the outliers without uniform bound on the reconstruction error.

**Definition 3.1.** For a reconstruction method  $R_h$  we define a subset of all possible samples

$$\Xi = \Xi(h, C_1, C_2) = \{\mathbf{z} \in (\Omega \times Y)^n : (i) \text{ and } (ii) \text{ hold}\},$$

where the two stated conditions are:

- (i) The uniform error of the reconstruction  $R_h(\mathbf{z})$  is bounded, i.e., for  $1 \leq i \leq n$

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

(ii) Changing one node will not do much damage, i.e., for all  $\mathbf{x} \in \Omega$  we assume for every  $1 \leq i \leq n$  the  $C_2\mathbb{1}$ -boundedness of  $\mathbf{z}_{-i} \mapsto R_h(\mathbf{z}_{-i})(\mathbf{x})$ .

**Remark 3.2.** (i) Note that, by applying the triangle inequality, we could use  $C_2 \leq 2C_1$  and only rely on the first assumption. For that reason we will state all results in two ways: one version using only  $C_1$  for simplicity and another using both constants to allow for fine-tuning of the bounds.

(ii) For many reconstruction methods one has a bound on the uniform error in a probabilistic fashion in the form of

$$\mathbb{P}\{\|R_h(\mathbf{Z}') - f\|_\infty > C_1\} \leq \gamma$$

for some small  $\gamma$ , e.g. [40, Section 6.3 and 6.4] or one of [25, 26, 34, 23]. To extend this to the context of assumption (i), we apply this bound for  $\mathbf{Z}_{-i}$  and  $1 \leq i \leq n$ . Union bound then gives

$$\begin{aligned} \mathbb{P}\{\mathbf{Z} \notin \Xi(h, \varepsilon, 2\varepsilon)\} &= \mathbb{P}\{\exists 1 \leq i \leq n : \|R_h(\mathbf{Z}_{-i}) - f\|_\infty > C_1\} \\ &\leq \sum_{i=1}^n \mathbb{P}\{\|R_h(\mathbf{Z}_{-i}) - f\|_\infty > C_1\} \\ &\leq n\gamma. \end{aligned}$$

For instance, in reconstructing functions via least squares, it has been shown that  $\gamma$  decays faster than  $1/n$  and the overall probability gets small, cf. [34]. This supports the sanity of the stated set.

We now want to show the  $\mathbf{c}$ -boundedness of the risk functional on  $\Xi$  in order to apply Theorem 2.2 for a concentration inequality.

**Lemma 3.3.** Let  $\Xi = \Xi(h, C_1, C_2)$  be the set of samples from Definition 3.1 and  $\mathbf{c} = 2C_1C_2\mathbb{1} \in \mathbb{R}^n$ . Then the risk functionals  $\mathbf{z} \mapsto \mathcal{E}(R_h(\mathbf{z}_{-i}))$  are  $\mathbf{c}$ -bounded.

*Proof.* We have to check what happens if we change one component. For that let  $\mathbf{z}$  and  $\mathbf{z}' \in \Xi$  be such that they differ in one sample. By the definition of the risk functional and the third binomial formula we have

$$\begin{aligned} &|\mathcal{E}(R_h(\mathbf{z}_{-i})) - \mathcal{E}(R_h(\mathbf{z}'_{-i}))| \\ &= \left| \int_{\Omega} |R_h(\mathbf{z}_{-i})(\mathbf{x}) - f(\mathbf{x})|^2 d\rho(\mathbf{x}) - \int_{\Omega} |R_h(\mathbf{z}'_{-i})(\mathbf{x}) - f(\mathbf{x})|^2 d\rho(\mathbf{x}) \right| \\ &\leq \int_{\Omega} |R_h(\mathbf{z}'_{-i})(\mathbf{x}) - f(\mathbf{x}) + R_h(\mathbf{z}_{-i}) - f(\mathbf{x})| \cdot |R_h(\mathbf{z}'_{-i})(\mathbf{x}) - R_h(\mathbf{z}_{-i})(\mathbf{x})| d\rho(\mathbf{x}). \end{aligned}$$

Using property (i) and (ii) of  $\Xi$  leads to

$$|\mathcal{E}(R_h(\mathbf{z}'_{-i})) - \mathcal{E}(R_h(\mathbf{z}_{-i}))| \leq 2C_1C_2 \int_{\Omega} d\rho(\mathbf{x}).$$

Since  $\rho$  is a probability measure the above integral evaluates to one and we obtain the desired constant of  $2C_1C_2$ .

In  $\mathcal{E}(R_h(\mathbf{z}_{-i}))$  the variable  $z_i$  does not occur and, therefore, the corresponding  $c_i$  is arbitrary. To have a general  $\mathbf{c}$  for all  $1 \leq i \leq n$ , we use  $c_i = 2C_1C_2$  anyways and obtain the assertion.  $\blacksquare$

Now we state the theorem on the concentration of the risk functional.

**Theorem 3.4.** *Let  $\mathbf{Z} = (\mathbf{X}_i, f(\mathbf{X}_i))_{i=1}^n$  with  $\mathbf{X}_i$  distributed independent and identically according to  $\rho$  on  $\Omega$ . Further, let*

$$m = \mathbb{E}\{\mathcal{E}(R_h(\mathbf{Z}_{-i})) | \mathbf{Z} \in \Xi\},$$

be the expected value of the risk functionals  $\mathcal{E}(R_h(\mathbf{Z}_{-i}))$  restricted to  $\Xi = \Xi(h, C_1, C_2)$  from Definition 3.1, and  $\gamma = 1 - \mathbb{P}\{\mathbf{Z} \in \Xi\}$  the probability of  $\mathbf{Z}$  not being in  $\Xi$ .

Then for  $\varepsilon > 2\gamma n C_1 C_2$  and  $1 \leq i \leq n$  we obtain the concentration of the risk functionals

$$\begin{aligned} \mathbb{P}\{|\mathcal{E}(R_h(\mathbf{Z}_{-i})) - m| > \varepsilon\} &\leq 2\gamma + 2 \exp\left(-\left(\frac{\varepsilon}{\sqrt{2n}C_1C_2} - \sqrt{2n\gamma}\right)^2\right) \\ &\leq 2\gamma + 2 \exp\left(-\left(\frac{\varepsilon}{\sqrt{8n}C_1^2} - \sqrt{2n\gamma}\right)^2\right). \end{aligned}$$

*Proof.* Lemma 3.3 in combination with Theorem 2.2 yields for  $\varepsilon > 2\gamma n C_1 C_2$  the first inequality

$$\mathbb{P}\{|\mathbb{E}_{\mathbf{Z}'}\{\mathcal{E}(R_h(\mathbf{Z}'))\} - \mathcal{E}(R_h(\mathbf{Z}))| > \varepsilon\} \leq 2\gamma + 2 \exp\left(-\frac{2(\varepsilon - 2\gamma n C_1 C_2)^2}{4n C_1^2 C_2^2}\right).$$

The second inequality is due to Remark 3.2 (i).  $\blacksquare$

Next, we tackle the related problem with respect to the cross-validation score. First we take care of its  $\mathbf{c}$ -boundedness on  $\Xi$ .

**Lemma 3.5.** *Let  $\Xi = \Xi(h, C_1, C_2)$  be the set of samples from Definition 3.1 and  $\mathbf{c} = C_1(C_1/n + 2C_2)\mathbf{1} \in \mathbb{R}^n$ . Then the cross-validation score  $\mathbf{z} \mapsto \text{CV}(\mathbf{z}, h)$  is  $\mathbf{c}$ -bounded.*

*Proof.* We have to check what happens if we change one component. For symmetry reasons we only have a look at what happens if we change the first sample. Let  $\mathbf{z}, \mathbf{z}' \in \Xi$  be such that

$$\mathbf{z} = (z_1, \dots, z_n) \quad \text{and} \quad \mathbf{z}' = (z'_1, z_2, \dots, z_n).$$

By the triangle inequality we have

$$\begin{aligned} &|\text{CV}(\mathbf{z}, h) - \text{CV}(\mathbf{z}', h)| \\ &\leq \frac{1}{n} \left| |R_h(\mathbf{z}_{-1})(\mathbf{x}_1) - f(\mathbf{x}_1)|^2 - |R_h(\mathbf{z}'_{-1})(\mathbf{x}'_1) - f(\mathbf{x}'_1)|^2 \right| \\ &\quad + \frac{1}{n} \sum_{i=2}^n |R_h(\mathbf{z}_{-i})(\mathbf{x}_i) - f(\mathbf{x}_i) + R_h(\mathbf{z}'_{-i})(\mathbf{x}_i) - f(\mathbf{x}_i)| |R_h(\mathbf{z}_{-i})(\mathbf{x}_i) - R_h(\mathbf{z}'_{-i})(\mathbf{x}_i)|. \end{aligned}$$

Using the properties of  $\Xi$  and  $|a^2 - b^2| \leq \max\{a^2, b^2\}$ , we further estimate

$$\begin{aligned} |\text{CV}(\mathbf{z}, h) - \text{CV}(\mathbf{z}', h)| &\leq \frac{C_1^2 + 2(n-1)C_1C_2}{n} \\ &\leq C_1(C_1/n + 2C_2). \end{aligned}$$

■

The corresponding concentration result looks as follows.

**Theorem 3.6.** *Let  $\mathbf{Z} = (\mathbf{X}_i, f(\mathbf{X}_i))_{i=1}^n$  with  $\mathbf{X}_i$  distributed independent and identically according to  $\rho$  on  $\Omega$ . Further, let*

$$m = \mathbb{E}\{\text{CV}(\mathbf{Z}, h) | \mathbf{Z} \in \Xi\},$$

*be the expected value of the cross-validation score  $\text{CV}(\mathbf{Z}, h)$  restricted to  $\Xi = \Xi(h, C_1, C_2)$  from Definition 3.1, and  $\gamma = 1 - \mathbb{P}\{\mathbf{Z} \in \Xi\}$  the probability of  $\mathbf{Z}$  not being in  $\Xi$ .*

*Then for  $\varepsilon > 2\gamma n C_1 C_2 + \gamma C_1^2$  we obtain the concentration of the cross-validation score*

$$\begin{aligned} \mathbb{P}\{|\text{CV}(\mathbf{Z}, h) - m| > \varepsilon\} &\leq 2\gamma + 2 \exp\left(-\left(\frac{\sqrt{2}\varepsilon}{C_1(C_1/\sqrt{n} + 2\sqrt{n}C_2)} - \sqrt{2n}\gamma\right)^2\right) \\ &\leq 2\gamma + 2 \exp\left(-\left(\frac{\varepsilon}{3\sqrt{n}C_1^2} - \sqrt{2n}\gamma\right)^2\right) \end{aligned}$$

where the second inequality holds for  $n \geq 5$ .

*Proof.* Applying Lemma 3.5 and Theorem 2.2 gives the first inequality. The second one is obtained by using Remark 3.2 (i),  $n \geq 5$ , and basic calculus. ■

Next, we prepare the connection of the two previous theorems by connecting the expected values of the risk functional and the cross-validation score.

**Lemma 3.7.** *The expected value of the risk functional for  $n - 1$  nodes is equal to the expected value of the cross-validation score for  $n$  nodes, i.e.,*

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

for  $\mathbf{Z}' = (\mathbf{X}'_i, f(\mathbf{X}'_i))_{i=1}^{n-1}$  representing  $n-1$  samples and  $\mathbf{Z} = (\mathbf{X}_i, f(\mathbf{X}_i))_{i=1}^n$  representing  $n$  samples where  $\mathbf{X}_i, \mathbf{X}'_i$  are distributed independent and identically according to  $\rho$ .

*Proof.* Since for all  $1 \leq i \leq n$  the  $\mathbf{Z}_{-i}$  have the same distribution as  $\mathbf{Z}'$  we write

$$\mathbb{E}_{\mathbf{Z}'} \{\mathcal{E}(R_h(\mathbf{Z}'))\} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathbf{Z}_{-i}} \{\mathcal{E}(R_h(\mathbf{Z}_{-i}))\}.$$

Instead of using  $\mathbb{E}_{\mathbf{Z}_{-i}}$ , we use  $\mathbb{E}_{\mathbf{Z}}$  since  $Z_i$  does not occur in the corresponding terms

$$\mathbb{E}_{\mathbf{Z}'} \{ \mathcal{E}(R_h(\mathbf{Z}')) \} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathbf{Z}} \{ \mathcal{E}(R_h(\mathbf{Z}_{-i})) \} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathbf{Z}} \left\{ |R_h(\mathbf{Z}_{-i})(\mathbf{x}_i) - f(\mathbf{x}_i)|^2 \right\}.$$

By linearity of the expected value we obtain the assertion

$$\mathbb{E}_{\mathbf{Z}} \{ \mathcal{E}(R_h(\mathbf{Z}')) \} = \mathbb{E}_{\mathbf{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n |R_h(\mathbf{Z}_{-i})(\mathbf{x}_i) - f(\mathbf{x}_i)|^2 \right\} = \mathbb{E}_{\mathbf{Z}} \{ \text{CV}(\mathbf{Z}, h) \}.$$

■

Having all the necessary tools, we state a central theorem bringing together risk functional and cross-validation score.

**Theorem 3.8.** *Let  $\mathbf{Z} = (\mathbf{X}_i, f(\mathbf{X}_i))_{i=1}^n$  with  $\mathbf{X}_i$  distributed independent and identically according to  $\rho$  on  $\Omega$  and  $R_h: (\Omega \times Y)^n \rightarrow Y^\Omega$  be a reconstruction method. Further, let*

$$M = \sup_{\mathbf{x}_1, \dots, \mathbf{x}_{n-1} \in \Omega} \|R_h((\mathbf{x}_i, f(\mathbf{x}_i))_{i=1}^{n-1})\|_\infty$$

be a uniform bound on the reconstruction for arbitrary nodes and  $\gamma = 1 - \mathbb{P}\{\mathbf{Z} \in \Xi\}$  the probability of  $\mathbf{Z}$  not being in  $\Xi = \Xi(h, C_1, C_2) \subseteq (\Omega \times Y)^n$  from Definition 3.1.

Then for  $\varepsilon > 2\gamma \max\{4nC_1C_2 + C_1^2, (M + \|f\|_\infty)^2\}$  we have the concentration bound of the difference of cross-validation score  $\text{CV}(\mathbf{Z}, h)$  and risk functional  $\mathcal{E}(R_h(\mathbf{Z}_{-1}))$

$$\begin{aligned} & \mathbb{P} \{ |\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))| > \varepsilon \} \\ & \leq 2\gamma + 2 \exp \left( - \left( \frac{\varepsilon}{\sqrt{2}C_1(C_1/\sqrt{n} + 4\sqrt{n}C_2)} - \sqrt{2n\gamma} \right)^2 \right) \\ & \leq 2\gamma + 2 \exp \left( - \left( \frac{\varepsilon}{12\sqrt{n}C_1^2} - \sqrt{2n\gamma} \right)^2 \right) \end{aligned}$$

where the second inequality holds for  $n \geq 3$ . In particular, for  $\delta > 0$ , we have with probability larger than  $1 - 2(\gamma + \delta)$

$$\begin{aligned} & |\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))| \\ & \leq \max \left\{ 2\gamma(M + \|f\|_\infty)^2, \left( \sqrt{2}C_1 \left( \frac{C_1}{\sqrt{n}} + 4\sqrt{n}C_2 \right) \right) (\sqrt{2n\gamma} + \sqrt{-\log \delta}) \right\} \\ & \leq \max \left\{ 2\gamma(M + \|f\|_\infty)^2, 12\sqrt{n}C_1^2 (\sqrt{2n\gamma} + \sqrt{-\log \delta}) \right\}. \end{aligned}$$

*Proof.* By the triangle inequality<sup>1</sup> we have for fixed  $\mathbf{z} \in (\Omega \times Y)^n$

$$\begin{aligned} & |\text{CV}(\mathbf{z}, h) - \mathcal{E}(R_h(\mathbf{z}_{-1}))| \\ & \leq |\text{CV}(\mathbf{z}, h) - \mathcal{E}(R_h(\mathbf{z}_{-1})) - \mathbb{E} \{ \text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1})) | \mathbf{Z} \in \Xi \}| \\ & \quad + |\mathbb{E} \{ \text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1})) | \mathbf{Z} \in \Xi \}|. \end{aligned}$$

<sup>1</sup>One might argue that using triangle inequality with the expected values one loses all information on the specific sample  $\mathbf{z}$ , which worsens the bound. However, [3] suggests that  $\text{CV}(\cdot, h)$  estimates  $\mathbb{E}\{\mathcal{E}(R_h(\mathbf{Z}))\}$  rather than  $\mathcal{E}(R_h(\mathbf{z}))$  itself, which reasons for our approach.



By Lemma 3.7 we have  $\mathbb{E}\{\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))\} = 0$  and, thus, estimate the second summand by

$$\begin{aligned} & |\mathbb{E}\{\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1})) | \mathbf{Z} \in \Xi\} - \mathbb{E}\{\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))\}| \\ & \leq \int_{(\Omega \times Y)^n \setminus \Xi} |\text{CV}(\mathbf{z}, h) - \mathcal{E}(R_h(\mathbf{z}_{-1}))| \, d\mathbf{z} \\ & \leq \int_{(\Omega \times Y)^n \setminus \Xi} (M + \|f\|_\infty)^2 \, d\mathbf{z} \\ & \leq (M + \|f\|_\infty)^2 \gamma \end{aligned}$$

where the last inequality follows from  $\mathbb{P}\{\mathbf{Z} \notin \Xi\} \leq \gamma$ . Thus, we obtain

$$\begin{aligned} & \mathbb{P}\{|\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))| > \varepsilon\} \\ & \leq \mathbb{P}\left\{|\text{CV}(\mathbf{z}, h) - \mathcal{E}(R_h(\mathbf{z}_{-1})) - \mathbb{E}\{\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1})) | \mathbf{Z} \in \Xi\}| > \frac{\varepsilon}{2}\right\} \\ & \quad + \mathbb{P}\left\{(M + \|f\|_\infty)^2 \gamma > \frac{\varepsilon}{2}\right\}. \end{aligned}$$

By the assumption on  $\varepsilon$  the latter probability evaluates to zero.

It is left to bound the first summand. Similar to the proofs of Lemmata 3.3 and 3.5 we will bound the remaining concentration by Theorem 2.2. For  $\mathbf{z}$  and  $\mathbf{z}' \in \Xi$ , which differ in one component, we have

$$\begin{aligned} & |\text{CV}(\mathbf{z}, h) - \mathcal{E}(R_h(\mathbf{z}_{-1})) - \text{CV}(\mathbf{z}', h) + \mathcal{E}(R_h(\mathbf{z}'_{-1}))| \\ & \leq |\text{CV}(\mathbf{z}, h) - \text{CV}(\mathbf{z}', h)| + |\mathcal{E}(R_h(\mathbf{z}_{-1})) - \mathcal{E}(R_h(\mathbf{z}'_{-1}))| \\ & \leq 4C_1C_2 + \frac{C_1^2}{n}, \end{aligned}$$

i.e.,  $\text{CV}(\mathbf{z}, h) - \mathcal{E}(R_h(\mathbf{z}_{-1}))$  is  $\mathbf{c}$ -bounded. Thus, with Theorem 2.2 we obtain

$$\begin{aligned} & \mathbb{P}\{|\text{CV}(\mathbf{z}, h) - \mathcal{E}(R_h(\mathbf{z}_{-1})) - \mathbb{E}\{\text{CV}(\mathbf{z}, h) + \mathcal{E}(R_h(\mathbf{z}_{-1}))\}| > \varepsilon\} \\ & \leq 2\gamma + 2 \exp\left(-\left(\frac{\varepsilon}{\sqrt{2}C_1(C_1/\sqrt{n} + 4\sqrt{n}C_2)} - \sqrt{2n}\gamma\right)^2\right) \end{aligned}$$

for  $\varepsilon > 2\gamma(4nC_1C_2 + C_1^2)$ . ■

**Remark 3.9.** (i) If, for a specific reconstruction method  $R_h$ , we have

- a uniform bound  $M$  on the reconstructions  $R_h(\mathbf{z})$ ,  $\mathbf{z} = (\mathbf{x}_i, f(\mathbf{x}_i))_{i=1}^n \in (\Omega \times Y)^n$  and
- a bound  $C_1$  on the reconstructions error of  $R_h(\mathbf{z})$  which holds with probability  $1 - \gamma$ ,

then Theorem 3.8 states, that with slightly smaller probability  $1 - 2(\gamma + \delta)$ , computing the cross-validation score  $\text{CV}(\mathbf{z}, h)$  is the same as computing the risk  $\mathcal{E}(R_h(\mathbf{z}))$  up to a small additive constant  $\varepsilon$  that can be computed explicitly from  $C_1$ ,  $M$ ,  $\gamma$ , and  $\delta$ .

(ii) For now we have a statement for one reconstruction method  $R_h$ . But we easily obtain error guarantees for the parameter  $h_{\text{CV}}$  minimizing the cross-validation score  $\text{CV}(\mathbf{z}, \cdot)$ :

Let  $h^*$  be the minimizer of  $h \mapsto \mathcal{E}(R_h(\mathbf{z}))$ . By using

$$\begin{aligned} & \mathbb{P} \{ \mathcal{E}(R_{h_{\text{CV}}}(\mathbf{Z}_{-1})) - \mathcal{E}(R_{h^*}(\mathbf{Z}_{-1})) > \varepsilon \} \\ & \leq \mathbb{P} \{ \mathcal{E}(R_{h_{\text{CV}}}(\mathbf{Z}_{-1})) - \text{CV}(\mathbf{Z}, h_{\text{CV}}) + \text{CV}(\mathbf{Z}, h^*) - \mathcal{E}(R_{h^*}(\mathbf{Z}_{-1})) > \varepsilon \} \\ & \leq \mathbb{P} \left\{ |\mathcal{E}(R_{h_{\text{CV}}}(\mathbf{Z}_{-1})) - \text{CV}(\mathbf{Z}, h_{\text{CV}})| > \frac{\varepsilon}{2} \right\} + \mathbb{P} \left\{ |\text{CV}(\mathbf{Z}, h^*) - \mathcal{E}(R_{h^*}(\mathbf{Z}_{-1}))| > \frac{\varepsilon}{2} \right\} \end{aligned}$$

we apply Theorem 3.8 twice and have that with high probability minimizing the cross-validation score is just  $\varepsilon$  worse in terms of the risk.

**Remark 3.10.** In order to derive asymptotic rates out of Theorem 3.8, we fix the probability  $\delta$  and assume that the reconstruction error of  $R_h$  decays asymptotically as  $C_1 \sim n^{-r}$  with probability at least  $1 - n^{-2r}$ . Then the difference of cross-validation score  $\text{CV}$  and the risk functional  $\mathcal{E}(R_h(\mathbf{z}))$  decays like  $n^{1/2-2r}$ .

## 4 Application using Shepard's model

Since this paper was motivated by [17, Chapter 8], where Shepard's model was used in the context of binary kernels, it seemed natural to start off with this application. Shepard's model or the Nadaraya-Watson estimator is a special case of moving least squares. It was introduced in [32, 42, 37] and is now-days widely used for solving PDEs [33, 6], manifold learning [39], or computer graphics [36]. Introductory information about this topic can be found in [12].

The crucial ingredient in Shepard's model is a, often locally supported, kernel function  $K_h$ . Given a sampling  $\mathbf{z} = (x_i, f(x_i))_{i=1}^n$  the model has the form

$$R_h(\mathbf{z}) = \frac{\sum_{i=1}^m K_h(\cdot, x_i) f(x_i)}{\sum_{i=1}^m K_h(\cdot, x_i)}. \quad (4.1)$$

A one-dimensional example for differently localized kernels is shown in Figure 4.1, which emphasizes the importance of the kernel choice. In this section we propose cross-validation as a method for choosing an optimal kernel and give an explicit error bound for the difference of risk functional (1.1) and cross-validation score (1.2). This is verified with numerical examples.

### 4.1 Theory

For simplicity, we restrict the domain to be the one-dimensional torus  $\Omega = \mathbb{T}$  and  $Y = \mathbb{R}$ . A common assumption on which we rely is to use positive, radial kernels, i.e.

$$K_h(x, x') = k_h(d(x, x'))$$

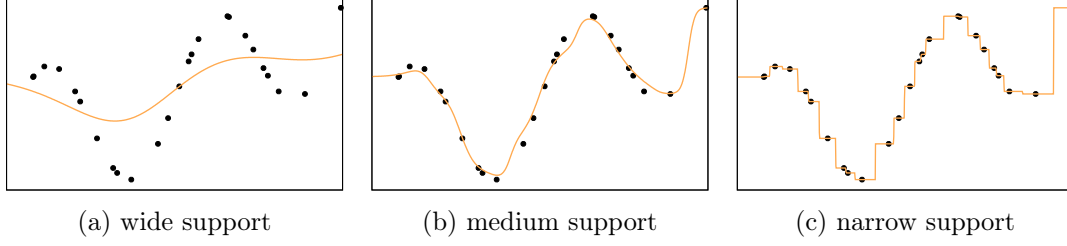


Figure 4.1: Shepard's model for different widths of the kernel support

for  $d(\cdot, \cdot)$  being the usual periodic distance on  $\mathbb{T}$  and  $k_h: [0, \infty) \rightarrow [0, \infty)$  a family of kernel functions with local support, i.e.,

$$\text{supp } k_h = \overline{\{t \in [0, \infty) : k_h(t) \neq 0\}} = [0, 1/h].$$

Note, that the range of the function  $R_h(\mathbf{z})$  is contained within the convex hull of all  $f(x_i)$ . Therefore, for samples  $\mathbf{z}$  from a bounded function  $f: \mathbb{T} \rightarrow \mathbb{R}$ , we have

$$M = \sup_{x_1, \dots, x_n \in \Omega} \|R_h((x_i, f(x_i))_{i=1}^n)\|_\infty \leq \|f\|_\infty. \quad (4.2)$$

Deterministic bounds on the approximation error are given in [12, Chapter 25]. These are based on the *mesh norm*

$$\delta_{\{x_1, \dots, x_n\}} := \max_{x \in \mathbb{T}} \min_{i=1, \dots, n} d(x, x_i).$$

For simplicity, we shall use only a simple bound which relies on stronger assumptions compared to [12, Chapter 25]. However, this still attains the same order in terms of the mesh norm.

**Lemma 4.1.** *Let  $k_h$  be supported on  $[0, 1/h]$  and  $f$  be Lipschitz continuous with constant  $L$ . Furthermore, we assume  $\delta_{\{x_1, \dots, x_n\}} < 1/h$ . Then*

$$\|R_h(\mathbf{z}) - f\|_\infty \leq \frac{L}{h}.$$

*Proof.* By the assumption on the mesh norm and the support of  $K_h$  we have

$$\sum_{i=1}^n K_h(x, x_i) > 0$$

for all  $x \in \mathbb{T}$ . Thus, we will not divide by zero in the following estimate. By the definition of Shepard's method we have

$$\begin{aligned} |R_h(\mathbf{z})(x) - f(x)| &= \left| \frac{\sum_{i=1}^m K_h(x, x_i) f(x_i)}{\sum_{i=1}^m K_h(x, x_i)} - f(x) \right| \\ &\leq \frac{\sum_{i=1}^m K_h(x, x_i) |f(x_i) - f(x)|}{\sum_{i=1}^m K_h(x, x_i)}. \end{aligned}$$

Using the Lipschitz condition and the local support we obtain

$$\begin{aligned} |R_h(\mathbf{z})(x) - f(x)| &\leq L \frac{\sum_{x_i \in [x-1/h, x+1/h]} K_h(x, x_i) |x_i - x|}{\sum_{x_i \in [x-1/h, x+1/h]} K_h(x, x_i)} \\ &\leq \frac{L}{h} \frac{\sum_{x_i \in [x-1/h, x+1/h]} K_h(x, x_i)}{\sum_{x_i \in [x-1/h, x+1/h]} K_h(x, x_i)} = \frac{L}{h}. \end{aligned}$$

■

As we draw samples randomly, we cannot guarantee an upper bound on the mesh norm  $\delta_{\{x_1, \dots, x_n\}}$ , but aim for a probabilistic result. Furthermore, in order to bound the approximation errors  $C_1$  from Definition 3.1 we actually need a bound for the mesh norms where single nodes are secluded, i.e., for  $\delta_{\{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}}$  and  $1 \leq i \leq n$ . To this end we define

$$\Xi = \{(x_i, f(x_i))_{i=1}^n : \delta_{\{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}} < 1/h \text{ for } 1 \leq i \leq n\}. \quad (4.3)$$

By the previous lemma we know, that for samples in  $\Xi$  the reconstruction error is bounded by  $L/h = C_1$ . With the following lemma we will show that the constructed set is in the paradigm of Definition 3.1 and  $\gamma = 1 - \mathbb{P}\{\mathbf{z} \in \Xi\}$  is close to zero.

**Lemma 4.2.** *For  $x_1, \dots, x_n \in \mathbb{T}$  drawn uniformly at random, we have*

$$\mathbb{P}\left\{\exists 1 \leq i \leq n : \delta_{\{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}} > \frac{1}{h}\right\} \leq \sum_{k=1}^{\lfloor h \rfloor} (-1)^{k+1} \binom{n}{k} \left(1 - \frac{k}{2h}\right)^{n-1}.$$

*Proof.* The given event on the mesh norm is equivalent to saying the distance of  $x_i$  to  $x_{i+2}$  will not exceed  $1/h$ . This is certainly fulfilled for nodes where the distance of  $x_i$  to  $x_{i+1}$  will not exceed  $1/(2h)$ . Therefore,

$$\mathbb{P}\left\{\exists 1 \leq i \leq n : \delta_{\{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}} > \frac{1}{h}\right\} \leq \mathbb{P}\left\{\delta_{\{x_1, \dots, x_n\}} > \frac{1}{2h}\right\}.$$

This probability has been calculated in [18, Theorem 2.1] which gives the assertion. ■

**Remark 4.3.** (i) *Note that similar techniques, involving  $\varepsilon$ -nets, can be applied to obtain results for more general domains, cf. [15].*

(ii) *Figure 4.2 depicts the probability of all mesh norms  $\delta_{\{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}}$ ,  $1 \leq i \leq n$  being bigger than  $1/h$  for  $n = 10\,000$  nodes estimated from numerical experiments. The critical point is around 1 000, where the probability increases away from zero. The theoretical bound from Lemma 4.2 is not optimal and has its critical point around 700.*

(iii) *The binomial bound in Lemma 4.2 is difficult to evaluate. In [11] it was show that for  $n \rightarrow \infty$  it converges to the Gumbel distribution, i.e.,*

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

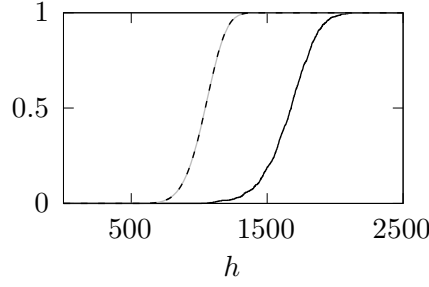


Figure 4.2: The probability of all mesh norms  $\delta_{\{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}}$ ,  $1 \leq i \leq n$  being bigger than  $1/h$  for  $n = 10\,000$  nodes. The solid line displays the numerical estimates from 1 000 experiments, the dashed line the upper bound from Lemma 4.2 and the gray line the asymptotic behaviour from Remark 4.3 (ii).

In Figure 4.2 we see that, already for 10 000 nodes, we are very close to this Gumbel distribution.

Now we have the necessary constants: the bound on the reconstruction  $M$  and the uniform bound on the reconstruction error  $C_1$  with its fail probability  $\gamma$  and are able to use the machinery of Section 3 to concentrate the difference of risk functional and cross-validation score.

**Theorem 4.4.** Let  $\mathbf{Z} = ((X_1, f(X_1)), \dots, (X_n, f(X_n)))$  represent  $n$  samples from a function  $f: \mathbb{T} \rightarrow \mathbb{R}$  with Lipschitz constant  $L$ , and  $R_h(\mathbf{Z})$  the reconstruction via Shepard's model, defined by (4.1), where the kernel  $k_h$  is supported on  $[0, 1/h]$ . Further, let

$$\gamma = \sum_{k=1}^{\lfloor h \rfloor} (-1)^{k+1} \binom{n}{k} \left(1 - \frac{k}{2h}\right)^{n-1} \quad \text{and} \quad \varepsilon > 2\gamma \max\{(4n+1)L^2/h^2, 4\|f\|_\infty^2\}.$$

Then we have the concentration bound of the difference of cross-validation score  $\text{CV}(\mathbf{Z}, h)$  and risk functional  $\mathcal{E}(R_h(\mathbf{Z}))$

$$\mathbb{P}\{|\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))| > \varepsilon\} \leq 2\gamma + 2 \exp\left(-\left(\frac{h^2\varepsilon}{12\sqrt{n}L^2} - \sqrt{2n\gamma}\right)^2\right).$$

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

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

*Proof.* By equation (4.2) we have  $M \leq \|f\|_\infty$ . With  $\Xi$  as in (4.3) we have by Lemmata 4.1 and 4.2

$$C_1 \leq \frac{L}{h} \quad \text{and} \quad \gamma \leq \sum_{k=1}^{\lfloor h \rfloor} (-1)^{k+1} \binom{n}{k} \left(1 - \frac{k}{2h}\right)^{n-1}.$$

Using these constants in Theorem 3.8 gives the assertion. ■

**Remark 4.5.** In order to interpret the error bounds in Theorem 4.4 asymptotically for  $n \rightarrow \infty$  we have to fix the desired probability  $\delta$ . Furthermore, we relate the kernel support  $1/h$  and the number of samples  $n$  via  $h = \alpha \cdot n$ . By Lemma 4.2 and Remark 4.3 we approximate the fail probability by

$$\gamma \lesssim \exp(-e^{-1/\alpha n}).$$

Inserting this bound into Theorem 4.4, we obtain with probability  $1 - 2(\exp(-e^{-1/\alpha n}) + \delta)$  that

$$|\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))| \sim \max \left\{ \exp(-e^{-1/\alpha n}), \frac{\exp(-e^{-1/\alpha n})}{n} + n^{-3/2} \right\} \lesssim n^{-3/2}.$$

**Remark 4.6.** The trade off between the constants  $C_1, C_2$ , and the fail probability  $\gamma$  is controlled by the construction of  $\Xi$ . In general, a larger set  $\Xi$  leads to a smaller fail probability  $\gamma$  but worse constants  $C_1$  and  $C_2$ .

In the extreme case we have  $\gamma = 0$  and  $\Xi$  consists of all possible data realizations, i.e.,  $\Xi = \{(x_i, f(x_i))_{i=1}^n : x_1, \dots, x_n \in \Omega\}$ . Then we have the bound  $C_1 = 2\|f\|_\infty$  as in equation (4.2). For the specific case of binary kernels, the estimate  $C_2 \sim 1/n$  can be found in [17, page 118] (with slight adaptations, as there is an individual  $C_2$  for every node  $x_i$  and one more assumption). With that, analogously to Theorem 4.4, we obtain with probability  $1 - 2\delta$

$$|\text{CV}(\mathbf{Z}, h) - \mathcal{E}(R_h(\mathbf{Z}_{-1}))| \sim \max \left\{ 0, 0 + n^{-1/2} \right\} \lesssim n^{-1/2}.$$

So, ignoring the restriction to binary kernels, the cost of improving to  $\gamma = 0$  is losing one order in  $n$ . This reasons for the construction of  $\Xi$  being a real subset of all possible data realizations.

## 4.2 Implementation

Before presenting our numerical experiments in Section 4.3, we give a brief discussion on the computational complexity of evaluating the model (4.1) as well as computing the cross-validation score  $\text{CV}(\mathbf{z}, h)$ . Evaluating the model (4.1) in nodes  $\tilde{x}_1, \dots, \tilde{x}_{\tilde{n}}$  needs two matrix-vector multiplications with

$$[K_h(x_i, \tilde{x}_j)]_{i=1, \dots, \tilde{n}, j=1, \dots, n}.$$

In [13] a method is proposed to compute (4.1) in a fast manner using the nonequispaced fast Fourier transform [20] which works for global kernels. Since we are dealing with locally supported kernels, we use sparse matrices for an efficient implementation. To compute the cross-validation score we need to compute  $R_h(\mathbf{z}_{-i})(x_i)$  for  $1 \leq i \leq n$ . To circumvent setting up  $n$  models we use the following trick. For fixed  $i$ , we obtain

$$\begin{aligned} r_i := R_h(\mathbf{z}_{-i}, h)(x_i) &= \frac{\sum_{j \in \{1, \dots, n\} \setminus \{i\}} K_h(x_j, x_i) f(x_j)}{\sum_{j \in \{1, \dots, n\} \setminus \{i\}} K_h(x_j, x_i)} \\ &= \frac{\sum_{j=1}^n K_h(x_j, x_i) f(x_j) - k_h(0) f(x_i)}{\sum_{j=1}^n K_h(x_j, x_i) - k_h(0)}. \end{aligned}$$

This favors the Algorithm 1 to compute the cross-validation score.

---

**Algorithm 1** Fast cross-validation for Shepard’s model

---

**Input:** data  $\mathbf{z} \in (\mathbb{T} \times \mathbb{R})^n$

**Output:** cross-validation score  $\text{CV}(\mathbf{z}, h)$

```

1: for  $i = 1, \dots, n$  do
2:    $n_i \leftarrow \sum_{j=1}^n K_h(x_j, x_i) f(x_j)$            {numerator of Shepard’s model}
3:    $d_i \leftarrow \sum_{j=1}^n K_h(x_j, x_i)$            {denominator of Shepard’s model}
4: end for
5: for  $i = 1, \dots, n$  do
6:    $r_i = (n_i - k_h(0) f(x_i)) / (d_i - k_h(0))$ 
7: end for
8:  $\text{CV}(\mathbf{z}, h) = \frac{1}{n} \sum_{i=1}^n |r_i - f(x_i)|^2$ 

```

---

In terms of complexity we obtain the same as for evaluating the model, namely, two matrix-vector multiplications.

### 4.3 Numerics

To exemplify our findings, we present some numerical experiments. We use the function  $f(x) = \sqrt{2} \sin(2\pi x)$  on  $\mathbb{T}$  with  $\|f\|_{L_2(\mathbb{T})} = 1$ ,  $\|f\|_\infty = \sqrt{2}$ , and Lipschitz constant  $L = \sqrt{2}$ . Further, we choose the simple hat kernel function

$$k_h(t) = \max\{0, 1 - ht\}.$$

We then repeat the following experiment 1 000 times for 50 different parameters  $h$ :

- (i) Choose  $n = 10\,000$  uniformly random nodes  $x_1, \dots, x_n$ .
- (ii) Compute function samples  $\mathbf{z} = (x_i, f(x_i))_{i=1}^n$ .
- (iii) Compute the reconstruction  $R_h(\mathbf{z})$  and approximate the risk  $\mathcal{E}(R_h(\mathbf{z}))$  by using evaluations in equispaced nodes.
- (iv) Compute the cross-validation score  $\text{CV}(\mathbf{z}, h)$  via Algorithm 1.

Figure 4.3 (a) shows the risk  $\mathcal{E}(R_h(\mathbf{z}))$  and (b) the cross-validation  $\text{CV}(\mathbf{z}, h)$  score for every experiment as a single dot. We observe, that both graphics resemble each other quite nicely. Both, the risk  $\mathcal{E}(R_h(\mathbf{z}))$  and the cross-validation  $\text{CV}(\mathbf{z}, h)$ , increase for small  $h$  and become increasingly unstable for  $h > 1500$  as the support of  $K_h$  gets too small.

In order to summarize the statistical behaviour we depicted in Figure 4.3 (c) the corresponding mean values and the intervals where 90% of the outcomes landed with respect the parameter  $h$ . The dashed lines depict our concentration bounds from Theorems 3.4

and 3.6. Setting the probability to 0.9, as in the experiment, we obtain the concentration bounds

$$\varepsilon \leq \alpha \frac{L^2}{h^2} \left( \sqrt{2n\gamma} + \sqrt{-n \log \left( \frac{p}{2} - \gamma \right)} \right) \quad (4.4)$$

for the risk functional with  $\alpha = \sqrt{8}$  and the cross-validation score with  $\alpha = 3$ . For the fail probability  $\gamma$  we used the numerical estimate from Remark 4.3 instead of the theoretical value from Lemma 4.2.

Finally, we depicted in Figure 4.3 (d) the 90%-quantile of the difference between the cross-validation score and risk functional. It illustrates that the risk functional and the cross-validation score coincide very well in the parameter region  $200 < h < 1500$  of interest. Our main result in Theorem 4.4 confirms this by a theoretical bound on this 90%-quantile. The theoretical bound has exactly the form (4.4) with  $\alpha = 12$  and is plotted as a dashed line.

In Figure 4.3 (c) and (d) our theoretical bounds rise rapidly at  $h \approx 1500$  which coincides with the beginning of instability in the computation of Shepard's model.

## 5 Conclusion

In this paper we presented a framework for obtaining bounds for the difference of cross-validation score and risk functional with high probability. This speaks for the use of cross-validation in parameter choice questions. In contrast to most previous results, we obtain a pre-asymptotic statement.

Along the way we proved concentration inequalities for the cross-validation score and risk functional, respectively. Connecting their expected values, we were able to combine both concentration inequalities and build a machinery to bound their difference with high probability. All those results are based on uniform bounds of the reconstruction method, which must hold in a subset of all possible samples. Estimates of this type are broadly available in learning theory.

For demonstration purposes we used Shepard's model on the one-dimensional torus with a rather simple bound of the uniform error. Numerical examples with a fast implementation support our results.

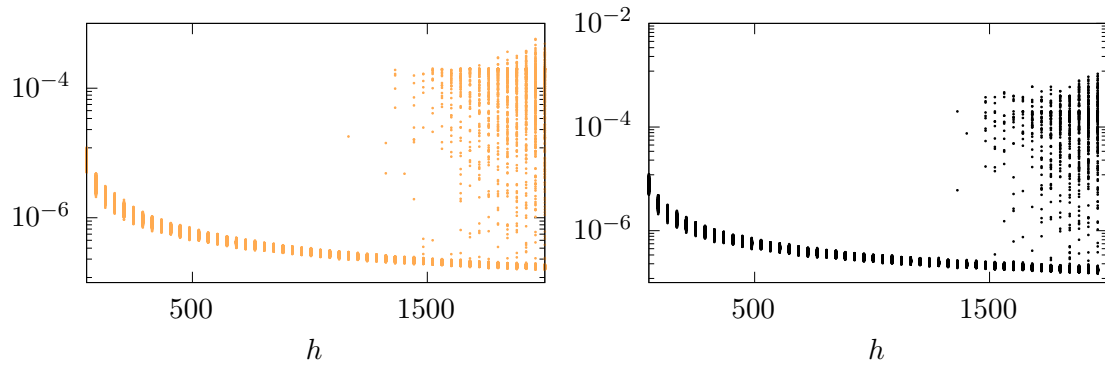
## Acknowledgments

Felix Bartel acknowledges funding by the European Social Fund (ESF), Project ID 100367298.

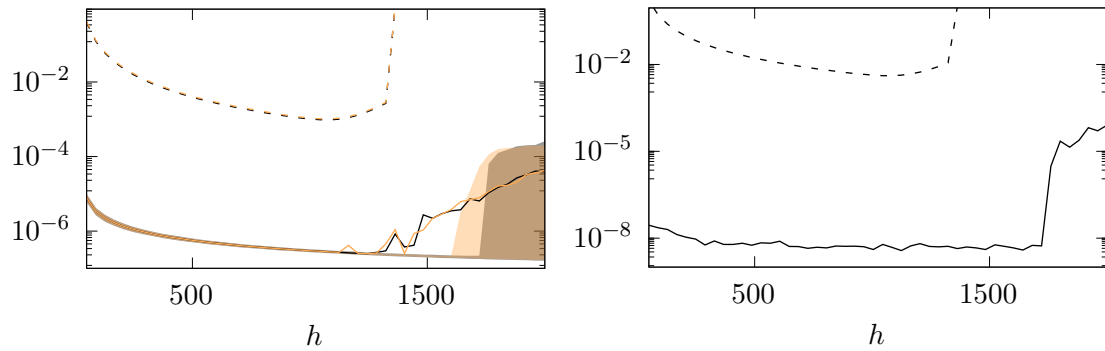
## References

- [1] A. B. Bakushinskiĭ. Remarks on the choice of regularization parameter from quasioptimality and relation tests. *Zh. Vychisl. Mat. i Mat. Fiz.*, 24(8):1258–1259, 1984.





(a) Cross-validation score  $CV(\mathbf{z}, h)$  for every experiment. (b) Risk functional  $\mathcal{E}(R_h(\mathbf{z}))$  for every experiment.



(c) The solid lines are the mean values of the risk functional (black) and the cross-validation score (orange). The transparent tubes represent 90% of all outcomes. The dashed lines are our theoretical bounds for these regions. (d) The solid line is the 90%-quantile of the differences between the risk functional and cross-validation score. The dashed line is our theoretical bound for this quantity.

Figure 4.3: Numerical example on  $\mathbb{T}$

- [2] F. Bartel, R. Hielscher, and D. Potts. Fast cross-validation in harmonic approximation. *Appl. Comput. Harmon. Anal.*, 49(2):415–437, 2020.
- [3] S. Bates, T. Hastie, and R. Tibshirani. Cross-validation: what does it estimate and how well does it do it? *ArXiv e-prints*, 2021.
- [4] F. Bauer and M. Reiß. Regularization independent of the noise level: an analysis of quasi-optimality. *Inverse Problems*, 24(5):055009, 16, 2008.
- [5] S. M. A. Becker. Regularization of statistical inverse problems and the Bakushinskiĭ veto. *Inverse Problems*, 27(11):115010, 22, 2011.
- [6] T. Belytschko, Y. Y. Lu, and L. Gu. Element-free Galerkin methods. *Internat. J. Numer. Methods Engrg.*, 37(2):229–256, 1994.
- [7] H. Blockeel and J. Struyf. Efficient algorithms for decision tree cross-validation. *J. Mach. Learn. Res.*, 3:621–650, 01 2002.
- [8] R. Combes. An extension of mcdiarmid’s inequality. *ArXiv e-prints*, abs/1511.05240, 2015.
- [9] E. De Vito, S. Pereverzyev, and L. Rosasco. Adaptive kernel methods using the balancing principle. *Found. Comput. Math.*, 10(4):455–479, 2010.
- [10] L. N. Deshpande and D. Girard. Fast computation of cross-validated robust splines and other non-linear smoothing splines. *Curves and Surfaces*, pages 143–148, 1991.
- [11] L. Devroye. Laws of the iterated logarithm for order statistics of uniform spacings. *The Annals of Probability*, 9(5), Oct. 1981.
- [12] G. E. Fasshauer. *Meshfree approximation methods with MATLAB*. World Scientific Publishers, 2007.
- [13] G. E. Fasshauer and J. Zhang. Recent results for moving least squares approximation. In L. Lucian and M. Neamtu, editors, *Geometric Modeling and Computing*, pages 163–176, Brentwood, 2003. Nashboro Press.
- [14] G. H. Golub, M. Heath, and G. Wahba. Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics*, 21(2):215–223, 1979.
- [15] Y. Gordon, A. E. Litvak, A. Pajor, and N. Tomczak-Jaegermann. Random  $\epsilon$ -nets and embeddings in  $l_\infty^N$ . *Studia Math.*, 178(1):91–98, 2007.
- [16] C. Gu. *Smoothing spline ANOVA models*, volume 297 of *Springer Series in Statistics*. Springer, New York, second edition, 2013.
- [17] L. Györfi, M. Kohler, A. Krzyżak, and H. Walk. *A Distribution-Free Theory of Nonparametric Regression*. Springer Series in Statistics. Springer-Verlag, New York, 2002.

- [18] L. Holst. On the lengths of the pieces of a stick broken at random. *J. Appl. Probab.*, 17(3):623–634, 1980.
- [19] S. Kale, R. Kumar, and S. Vassilvitskii. Cross-validation and mean-square stability. In *Second Symposium on Innovations in Computer Science (ICS2011)*, pages 487–495, 2011.
- [20] J. Keiner, S. Kunis, and D. Potts. NFFT 3.5, C subroutine library. <http://www.tu-chemnitz.de/~potts/nfft>. Contributors: F. Bartel, M. Fenn, T. Görner, M. Kircheis, T. Knopp, M. Quellmalz, M. Schmiscke, T. Volkmer, A. Vollrath.
- [21] S. Kindermann and A. Neubauer. On the convergence of the quasioptimality criterion for (iterated) Tikhonov regularization. *Inverse Probl. Imaging*, 2(2):291–299, 2008.
- [22] S. Kindermann, S. Pereverzyev, Jr., and A. Pilipenko. The quasi-optimality criterion in the linear functional strategy. *Inverse Problems*, 34(7):075001, 24, 2018.
- [23] D. Krieg, E. Novak, and M. Sonnleitner. Recovery of sobolev functions restricted to iid sampling. *ArXiv e-prints*, 2021.
- [24] R. Kumar, D. Lokshtanov, S. Vassilvitskii, and A. Vattani. Near-optimal bounds for cross-validation via loss stability. *30th International Conference on Machine Learning, ICML 2013*, pages 27–35, 01 2013.
- [25] R. J. Kunsch. Breaking the curse for uniform approximation in hilbert spaces via monte carlo methods. *Journal of Complexity*, 48:15–35, Oct. 2018.
- [26] A. Lederer, J. Umlauf, and S. Hirche. Uniform error bounds for gaussian process regression with application to safe control. *ArXiv e-prints*, 2019.
- [27] K.-C. Li. Asymptotic optimality of  $C_L$  and generalized cross-validation in ridge regression with application to spline smoothing. *Ann. Statist.*, 14(3):1101–1112, 1986.
- [28] M. A. Lukas. Robust generalized cross-validation for choosing the regularization parameter. *Inverse Problems*, 22(5):1883–1902, 2006.
- [29] M. A. Lukas, F. R. de Hoog, and R. S. Anderssen. Efficient algorithms for robust generalized cross-validation spline smoothing. *J. Comput. Appl. Math.*, 235:102–107, 2010.
- [30] C. McDiarmid. On the method of bounded differences. In *Surveys in combinatorics, 1989 (Norwich, 1989)*, volume 141 of *London Math. Soc. Lecture Note Ser.*, pages 148–188. Cambridge Univ. Press, Cambridge, 1989.
- [31] M. Mullin and R. Sukthankar. Complete cross-validation for nearest neighbor classifiers. In *17th International Conference on Machine Learning (ICML)*, 2000.

- [32] E. A. Nadaraya. On estimating regression. *Theory of Probab. Appl.*, 9:141–142, 1964.
- [33] B. Nayroles, G. Touzot, and P. Villon. Generalizing the finite element method: diffuse approximation and diffuse elements. *Comput. Mech.*, 10(5):307–318, 1992.
- [34] K. Pozharska and T. Ullrich. A note on sampling recovery of multivariate functions in the uniform norm. *ArXiv e-prints*, 2021.
- [35] S. Rosset. Bi-level path following for cross validated solution of kernel quantile regression. *J. Mach. Learn. Res.*, 10:2473–2505, 2009.
- [36] S. Schaefer, T. McPhail, and J. Warren. Image deformation using moving least squares. *ACM Trans. Graph.*, 25(3):533–540, July 2006.
- [37] D. Shepard. A two-dimensional interpolation function for irregularly-spaced data. In *Proceedings of the 1968 23rd ACM National Conference*, ACM '68, page 517–524, New York, NY, USA, 1968. Association for Computing Machinery.
- [38] R. B. Sidje, A. B. Williams, and K. Burrage. Fast generalized cross validation using Krylov subspace methods. *Numer. Algor.*, 47:109–131, 2008.
- [39] B. Sober and D. Levin. Manifold approximation by moving least-squares projection (MMLS). *Constr. Approx.*, 52(3):433–478, 2020.
- [40] I. Steinwart and A. Christmann. *Support Vector Machines*. Springer Publishing Company, Incorporated, 1st edition, 2008.
- [41] M. Tasche and N. Weyrich. Smoothing inversion of Fourier series using generalized cross-validation. *Results Math.*, 29(1-2):183–195, 1996.
- [42] G. S. Watson. Smooth regression analysis. *Sankhy = a Ser. A*, 26:359–372, 1964.
- [43] H. L. Weinert. Efficient computation for Whittaker-Henderson smoothing. *Comp. Stat. & Data Analysis*, 52:959–974, 2007.