



**HAL**  
open science

# Efficiency of the V-fold model selection for localized bases

Fabien Navarro, Adrien Saumard

► **To cite this version:**

Fabien Navarro, Adrien Saumard. Efficiency of the V-fold model selection for localized bases. Conference of the International Society for Non-Parametric Statistics, Jun 2016, Avignon, France. 10.1007/978-3-319-96941-1\_4. hal-01505514v2

**HAL Id: hal-01505514**

**<https://hal.science/hal-01505514v2>**

Submitted on 16 Dec 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Efficiency of the $V$ -fold model selection for localized bases

Fabien Navarro and Adrien Saumard

**Abstract** Many interesting functional bases, such as piecewise polynomials or wavelets, are examples of localized bases. We investigate the optimality of  $V$ -fold cross-validation and a variant called  $V$ -fold penalization in the context of the selection of linear models generated by localized bases in a heteroscedastic framework. It appears that while  $V$ -fold cross-validation is not asymptotically optimal when  $V$  is fixed, the  $V$ -fold penalization procedure is optimal. Simulation studies are also presented.

**Key words:**  $V$ -fold cross-validation,  $V$ -fold penalization, model selection, non-parametric regression, heteroscedastic noise, random design, wavelets

## 1 Introduction

$V$ -fold cross-validation type procedures are extremely used in Statistics and machine learning, with however a rather small set of theoretical results on it ([3]). This paper aims at investigating from the theoretical point of view and on simulations, the efficiency of two  $V$ -fold strategies for model selection in a heteroscedastic regression setting, with random design. From the one hand, we investigate the behaviour of the classical  $V$ -fold cross-validation to select, among other examples, linear models of wavelets. As pointed out in the case of histogram selection in [2], this procedure is not asymptotically optimal when  $V$  is fixed, as it is the case in practice where  $V$  is usually taken to be equal to 5 or 10. On the other hand, we study the  $V$ -fold penalization proposed by Arlot [2] and show its efficiency in our general context.

---

Fabien Navarro  
CREST-ENSAI-UBL, BRUZ, FRANCE e-mail: fabien.navarro@ensai.fr

Adrien Saumard  
CREST-ENSAI-UBL, BRUZ, FRANCE e-mail: adrien.saumard@ensai.fr

More precisely, the present paper is devoted to an extension of some results obtained in [16] related to efficiency of cross-validation type procedures. Indeed, as remarked in [16] (see Remark 5.1 therein) our results obtained for the selection of linear models endowed with a strongly localized basis (see Definition **(Aslb)**, Section 2.1 of [16]) can be extended to more general and more classical localized bases, at the price of considering only models with sufficiently small dimensions. Rigorous proofs are given here and further simulation studies are explored.

The paper is organized as follows. In Section 2, we describe our model selection setting. Then  $V$ -fold cross-validation is considered in Section 3, while the efficiency of  $V$ -fold penalization is tackled in Section 4. A simulation study is reported in Section 5. The proofs are exposed in Section 6.

## 2 Model selection setting

Assume that we observe  $n$  independent pairs of random variables  $\xi_i = (X_i, Y_i) \in \mathcal{X} \times \mathbb{R}$  with common distribution  $P$ . For convenience, we also denote by  $\xi = (X, Y)$  a random pair, independent of the sample  $(\xi_1, \dots, \xi_n)$ , following the same distribution  $P$ . The set  $\mathcal{X}$  is called the feature space and we assume  $\mathcal{X} \subset \mathbb{R}^d$ ,  $d \geq 1$ . We denote  $P^X$  the marginal distribution of the design  $X$ . We assume that the following regression relation is valid,

$$Y = s_*(X) + \sigma(X) \varepsilon ,$$

with  $s_* \in L_2(P^X)$  the regression function that we aim at estimating. Conditionally to  $X$ , the residual  $\varepsilon$  is normalized, i.e. it has mean zero and variance one. The function  $\sigma : \mathcal{X} \rightarrow \mathbb{R}_+$  is a heteroscedastic noise level, assumed to be unknown.

To produce an estimator of  $s_*$ , we are given a finite collection of models  $\mathcal{M}_n$ , with cardinality depending on the amount  $n$  of data. Each model  $m \in \mathcal{M}_n$  is taken to be a finite-dimensional vector space, of linear dimension  $D_m$ . We will further detail in a few lines the analytical structure of the models.

We set  $\|s\|_2 = (\int_{\mathcal{X}} s^2 dP^X)^{1/2}$  the quadratic norm in  $L_2(P^X)$  and  $s_m$  the orthogonal - with respect to the quadratic norm - projection of  $s_*$  onto  $m$ . For a function  $f \in L_1(P)$ , we write  $P(f) = Pf = \mathbb{E}[f(\xi)]$ . We call the least squares contrast a functional  $\gamma : L_2(P^X) \rightarrow L_1(P)$ , defined by

$$\gamma(s) : (x, y) \mapsto (y - s(x))^2 , \quad s \in L_2(P^X) .$$

Using these notations, the regression function  $s_*$  is the unique minimizer of the risk,

$$s_* = \arg \min_{s \in L_2(P^X)} P(\gamma(s)) .$$

The projections  $s_m$  are also characterized by

$$s_m = \arg \min_{s \in m} P(\gamma(s)) .$$

To each model  $m \in \mathcal{M}_n$ , we associate a least squares estimator  $\hat{s}_m$ , defined by

$$\begin{aligned} \hat{s}_m &\in \arg \min_{s \in m} \{P_n(\gamma(s))\} \\ &= \arg \min_{s \in m} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - s(X_i))^2 \right\} , \end{aligned}$$

where  $P_n = n^{-1} \sum_{i=1}^n \delta_{\xi_i}$  is the empirical measure associated to the sample.

The accuracy of estimation is tackled through the excess loss of the estimators,

$$\ell(s_*, \hat{s}_m) := P(\gamma(\hat{s}_m) - \gamma(s_*)) = \|\hat{s}_m - s_*\|_2^2 .$$

The following "bias-variance" decomposition holds,

$$\ell(s_*, \hat{s}_m) = \ell(s_*, s_m) + \ell(s_m, \hat{s}_m) ,$$

where

$$\begin{aligned} \ell(s_*, s_m) &:= P(\gamma(s_m) - \gamma(s_*)) = \|s_m - s_*\|_2^2 \\ \ell(s_m, \hat{s}_m) &:= P(\gamma(\hat{s}_m) - \gamma(s_m)) \geq 0 . \end{aligned}$$

The deterministic quantity  $\ell(s_*, s_m)$  is called the bias of the model  $m$ , while the random variable  $\ell(s_m, \hat{s}_m)$  is called the excess loss of the least squares estimator  $\hat{s}_m$  on the model  $m$ . By the Pythagorean Theorem, we have

$$\ell(s_m, \hat{s}_m) = \|\hat{s}_m - s_m\|_2^2 .$$

From the collection of models  $\mathcal{M}_n$ , we aim at proposing an estimator that is as close as possible in terms of excess loss to an oracle model  $m_*$ , defined by

$$m_* \in \arg \min_{m \in \mathcal{M}_n} \{\ell(s_*, \hat{s}_m)\} .$$

We choose to select an estimator from the collection  $\{\hat{s}_m ; m \in \mathcal{M}_n\}$ . Hence, the selected model is denoted  $\hat{m}$ . The goal is to ensure that the selected estimator achieves an oracle inequality of the form

$$\ell(s_*, \hat{s}_{\hat{m}}) \leq C \times \inf_{m \in \mathcal{M}_n} \ell(s_*, \hat{s}_m) ,$$

for a constant  $C \geq 1$  as close as possible to one and on an event of probability close to one.

### 3 V-fold cross-validation

For convenience, let us denote in the following  $\widehat{s}_m(P_n)$  the least squares estimator built from the empirical distribution  $P_n = 1/n \sum_{i=1}^n \delta_{(X_i, Y_i)}$ . To perform the  $V$ -fold cross-validation (VFCV) procedure, we consider a partition  $(B_j)_{1 \leq j \leq V}$  of the index set  $\{1, \dots, n\}$  and set,

$$P_n^{(j)} = \frac{1}{\text{Card}(B_j)} \sum_{i \in B_j} \delta_{(X_i, Y_i)} \quad \text{and} \quad P_n^{(-j)} = \frac{1}{n - \text{Card}(B_j)} \sum_{i \notin B_j} \delta_{(X_i, Y_i)} .$$

We assume that the partition  $(B_j)_{1 \leq j \leq V}$  is regular: for all  $j \in \{1, \dots, V\}$ ,  $\text{Card}(B_j) = n/V$ . It is worth noting that we can always define our partition such  $\sup_j |\text{Card}(B_j) - n/V| < 1$  so that the assumption of regular partition is only a slight approximation of the general case. Let us write  $\widehat{s}_m^{(-j)} = \widehat{s}_m(P_n^{(-j)})$  the estimators built from the data in the block  $B_j$ . Now, the selected model  $\widehat{m}_{\text{VFCV}}$  is taken equal to any model optimizing the  $V$ -fold criterion,

$$\widehat{m}_{\text{VFCV}} \in \arg \min_{m \in \mathcal{M}_n} \{\text{crit}_{\text{VFCV}}(m)\} , \quad (1)$$

where

$$\text{crit}_{\text{VFCV}}(m) = \frac{1}{V} \sum_{j=1}^V P_n^{(j)} \gamma(\widehat{s}_m^{(-j)}) . \quad (2)$$

Let us now detail the set of assumptions under which we will investigate the accuracy of VFCV.

#### Set of assumptions: (SA)

- (P1) Polynomial complexity of  $\mathcal{M}_n$ : there exist some constants  $c_{\mathcal{M}}, \alpha_{\mathcal{M}} > 0$  such that  $\text{Card}(\mathcal{M}_n) \leq c_{\mathcal{M}} n^{\alpha_{\mathcal{M}}}$ .
- (A1b) There exists a constant  $r_{\mathcal{M}}$  such that for each  $m \in \mathcal{M}_n$  one can find an orthonormal basis  $(\varphi_k)_{k=1}^{D_m}$  satisfying, for all  $(\beta_k)_{k=1}^{D_m} \in \mathbb{R}^{D_m}$ ,

$$\left\| \sum_{k=1}^{D_m} \beta_k \varphi_k \right\|_{\infty} \leq r_{\mathcal{M}} \sqrt{D_m} |\beta|_{\infty} , \quad (3)$$

where  $|\beta|_{\infty} = \max\{|\beta_k|; k \in \{1, \dots, D_m\}\}$ .

- (P2) Upper bound on dimensions of models in  $\mathcal{M}_n$ : there exists a positive constant  $A_{\mathcal{M},+}$  such that for every  $m \in \mathcal{M}_n$ ,  $D_m \leq A_{\mathcal{M},+} n^{1/3} (\ln n)^{-2}$ .
- (A1b) A positive constant  $A$  exists, that bounds the data and the projections  $s_m$  of the target  $s_*$  over the models  $m$  of the collection  $\mathcal{M}_n$ :  $|Y_i| \leq A < \infty$ ,  $\|s_m\|_{\infty} \leq A < \infty$  for all  $m \in \mathcal{M}_n$ .
- (An) Uniform lower-bound on the noise level:  $\sigma(X_i) \geq \sigma_{\min} > 0$  a.s.
- (Ap<sub>u</sub>) The bias decreases as a power of  $D_m$ : there exist  $\beta_+ > 0$  and  $C_+ > 0$  such that

$$\ell(s_*, s_m) \leq C_+ D_m^{-\beta_+} .$$

Assumption **(A1b)** refers to the classical concept of localized basis (Birgé and Massart [6]). It is proved in [5], Section 3.2.1, that linear models of piecewise polynomials with bounded degree on a regular partition of a bounded domain of  $\mathbb{R}^d$  are endowed with a localized basis. It is also proved that compactly supported wavelet expansions are also fulfilled with a localized basis on  $\mathbb{R}^d$ . However, the Fourier basis is not a localized basis. For some sharp concentration results related to the excess loss of least squares estimators built from the Fourier basis, we refer to [19].

The assumption **(A1b)** is more general than the assumption of strongly localized basis used in [16], but the price to pay for such generality is that, according to **(P2)** we can only consider models with dimensions  $D_m \ll n^{1/3}$ .

Assumption **(P1)** states that the collection as a polynomial cardinality with respect to the sample size, allowing in particular to consider a collection of models built from a basis expansion.

Then Assumption **(A1b)** is related to boundedness of the data and enables in particular to use Talagrand's type concentration inequalities for the empirical process. Going beyond the bounded setting would in particular bring much more technicalities that might darken our work. For an example of results in an unbounded setting, see for instance [4], dealing with optimal selection of regressograms (histograms being a very particular case of our general framework). Assumption **(A1n)** is essentially a technical assumption that allows to obtain sharp lower bounds for the excess losses of the estimators. Condition  $(Ap_u)$  is a very classical assumption in the model selection literature, specifying a rate of decay for the biases of the models. This assumption is classically satisfied for piecewise polynomials when the regression function belongs to a Sobolev space and for wavelet models whenever the target belongs to some Besov space (see for instance [5] for more details). The specific value of  $\beta_+$  parameter will only affect the value of the constants in the derived oracle inequalities.

**Theorem 1.** *Assume that **(SA)** holds. Let  $r \in (2, +\infty)$  and  $V \in \{2, \dots, n-1\}$  satisfying  $1 < V \leq r$ . Define the  $V$ -fold cross-validation procedure as the model selection procedure given by (1). Then, for all  $n \geq n_0((\mathbf{SA}), r)$ , with probability at least  $1 - L_{(\mathbf{SA}),r} n^{-2}$ ,*

$$\ell(s_*, \widehat{s}_{m_{VFCV}}) \leq \left(1 + \frac{L_{(\mathbf{SA}),r}}{\sqrt{\ln n}}\right) \inf_{m \in \mathcal{M}_n} \left\{ \ell(s_*, \widehat{s}_m^{(-1)}) \right\} + L_{(\mathbf{SA}),r} \frac{(\ln n)^3}{n} .$$

In Theorem 1, we prove an oracle inequality with principal constant tending to one when the sample size goes to infinity. This inequality bounds from above the excess loss of the selected estimator by the excess loss of the oracle learned with a fraction  $1 - V^{-1}$  of the original data. Ideally, one would, however, expect from an optimal procedure to recover the oracle built from the entire data. The next section is devoted to this task.

Parameter  $V$  (or  $r$ ) is considered in Theorem 1 as a constant, essentially for ease of presentation. Actually, the value of  $V$  may be allowed to depend on  $n$  but also on

the dimensions  $D_m$ , meaning that we may take different values of  $V$  according to the different models of the collection. More precisely, it can be seen from the arguments in the proofs (especially from Theorem 8 in [18]) that for each model  $m \in \mathcal{M}_n$ , it suffices to have  $V \leq \max \{D_m(\ln n)^{-\tau}; 2\}$  where  $\tau$  is any number in  $(1, 3)$  to ensure an oracle inequality with leading constant tending to one when the amount of data tends to infinity. In this case,  $r$  can not be considered as a parameter independent from the sample size anymore, but it can be checked that for the latter constraints on  $V$ , the constants  $n_0((\mathbf{SA}), r)$  and  $L_{(\mathbf{SA}), r}$  do not explode but are still uniformly bounded with respect to  $n$  and thus can be still considered as independent from  $n$ .

## 4 V-fold penalization

Now we investigate the behaviour of a penalization procedure proposed by Arlot [2] and called  $V$ -fold penalization,

$$\widehat{m}_{\text{penVF}} \in \arg \min_{n \in \mathcal{M}_n} \{ \text{crit}_{\text{penVF}}(m) \} ,$$

where

$$\text{crit}_{\text{penVF}}(m) = P_n(\gamma(\widehat{s}_m)) + \text{pen}_{\text{VF}}(m) ,$$

with

$$\text{pen}_{\text{VF}}(m) = \frac{V-1}{V} \sum_{j=1}^V \left[ P_n \gamma(\widehat{s}_m^{(-j)}) - P_n^{(-j)} \gamma(\widehat{s}_m^{(-j)}) \right] . \quad (4)$$

The property underlying the  $V$ -fold penalization is that the  $V$ -fold penalty  $\text{pen}_{\text{VF}}$  is an unbiased estimate of the ideal penalty  $\text{pen}_{\text{id}}$ , the latter allowing to identify the oracle  $m_*$ ,

$$\begin{aligned} m_* &\in \arg \min_{m \in \mathcal{M}_n} \{ P(\gamma(\widehat{s}_m)) \} \\ &= \arg \min_{m \in \mathcal{M}_n} \{ P_n(\gamma(\widehat{s}_m)) + \text{pen}_{\text{id}}(m) \} , \end{aligned}$$

where

$$\text{pen}_{\text{id}}(m) = P(\gamma(\widehat{s}_m)) - P_n(\gamma(\widehat{s}_m)) .$$

The following theorem states the asymptotic optimality of the  $V$ -fold penalization procedure for a fixed  $V$ .

**Theorem 2.** *Assume that  $(\mathbf{SA})$  holds. Let  $r \in (2, +\infty)$  and  $V \in \{2, \dots, n-1\}$  satisfying  $1 < V \leq r$ . Define the  $V$ -fold cross-validation procedure as the model selection procedure given by,*

$$\widehat{m}_{\text{penVF}} \in \arg \min_{n \in \mathcal{M}_n} \{ P_n(\gamma(\widehat{s}_m)) + \text{pen}_{\text{VF}}(m) \} .$$

*Then, for all  $n \geq n_0((\mathbf{SA}), r)$ , with probability at least  $1 - L_{(\mathbf{SA}), r} n^{-2}$ ,*

$$\ell\left(s_*, \widehat{s}_{\widehat{m}_{penVF}}\right) \leq \left(1 + \frac{L_{(SA),r}}{\sqrt{\ln n}}\right) \inf_{m \in \mathcal{M}_n} \{\ell(s_*, \widehat{s}_m)\} + L_{(SA),r} \frac{(\ln n)^3}{n}.$$

As for Theorem 1 above, parameter  $V$  (or  $r$ ) is considered in Theorem 2 as a constant but in fact, the value of  $V$  may be allowed to depend on  $n$  and even on the dimensions  $D_m$ , this case corresponding to possibly different choices  $V$  according to the models of the collection. As for Theorem 1, it is allowed to have  $V \leq \max\{D_m(\ln n)^{-\tau}; 2\}$  where  $\tau$  is any number in  $(1, 3)$  to ensure an oracle inequality with leading constant tending to one when the amount of data tends to infinity.

## 5 Simulation study

In order to assess the numerical performances of the model selection procedures we have discussed, a short simulation study was conducted. Particularly, to illustrate the theory developed above for the selection of linear estimators using the  $V$ -fold cross-validation and  $V$ -fold penalization, linear wavelet models were considered.

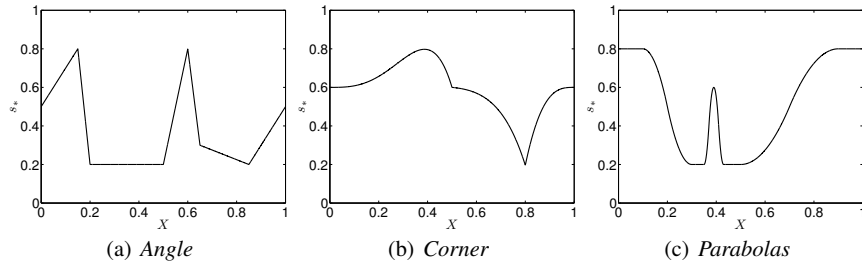
Despite the fact that a linear wavelet estimator is not as flexible, or potentially as powerful, as a nonlinear one, it still preserves the computational efficiency of wavelet methods and can provide comparative results to thresholding estimator, particularly when the unknown function is sufficiently smooth (see [1]).

The simulations were carried out using Matlab and the wavelet toolbox WaveLab850 [10]. The codes used to replicate the numerical results presented here will be available at <https://github.com/fabnavarro>. For more details on the numerical simulations and comparisons with other model selection procedures, we refer the reader to [19].

The simulated data were generated according to  $Y_i = s_*(X_i) + \sigma(X_i)\varepsilon_i$ ,  $i = 1, \dots, n$ , where  $n = 4096$ ,  $X_i$ 's are uniformly distributed on  $[0, 1]$ ,  $\varepsilon_i$ 's are independent  $\mathcal{N}(0, 1)$  variables and independent of  $X_i$ 's. The heteroscedastic noise level  $\sigma(x) = |\cos(10x)|/10$ . Daubechies' compactly-supported wavelet with 8 vanishing moments were used. Three standard regression functions with different degrees of smoothness (*Angle*, *Corner* and *Parabolas*, see [14, 7]) were considered. They are plotted in Figure 1 and a visual idea of the noise level is given in Figures 2(b).

The computation of wavelet-based estimators is straightforward and fast in the fixed design case, thanks to Mallat's pyramidal algorithm ([13]). In the case of random design, the implementation requires some changes and several strategies have been developed in the literature (see e.g. [8, 11]). In the regression with uniform design [9] have examined convergence rates when the unknown function is in a Hölder class. They showed that the standard equispaced wavelet method with universal thresholding can be directly applied to the nonequispaced data (without a loss in the rate of convergence). We have followed this approach since it preserves the computational simplicity and efficiency of the equispaced algorithm. In the context of wavelet regression in random design with heteroscedastic dependent errors [12]





**Fig. 1** (a)–(c): The three test functions used in the simulation study.

have also adopted this approach. Thus, the wavelet coefficients of the collection of models is computed by a simple application of Mallat’s algorithm using the ordered  $Y_i$ ’s as input variables. The collection is then constructed by successively adding whole resolution levels of wavelet coefficients. Thus, the considered dimensions are  $\{D_m, m \in \mathcal{M}_n\} = \{2^j, j = 1, \dots, J - 1\}$ , where  $J = \log 2(n)$  (the finest resolution level). Finally, the selected model are obtained by minimizing (2) and (4) over the set  $m \in \mathcal{M}_n$ . Note that these linear models operate in a global fashion since whole levels of coefficients are suppressed as opposed to thresholding methods.

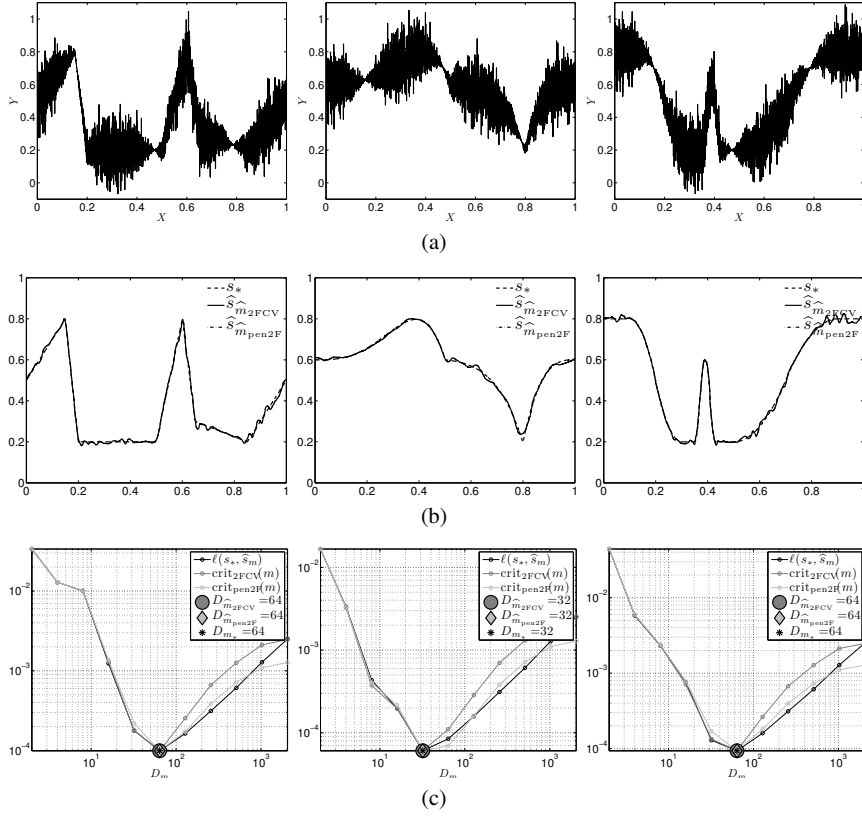
For choosing the threshold parameter in wavelet shrinkage Nason [15] adjusted the usual 2FCV method—which cannot be applied directly to wavelet estimation. In order to implement its strategy in a linear context, we test, for every model of the collection, an interpolated wavelet estimator learned from the (ordered) even-indexed data against the odd-indexed data and vice versa. More precisely, considering the data  $X_i$  are ordered, the selected model  $\hat{m}_{2FCV}$  (resp.  $\hat{m}_{pen2F}$ ) is obtained by minimizing (2) (resp. (4)) with  $V = 2$ ,  $B_1 = \{2, 4, \dots, n\}$  and  $B_2 = \{1, 3, \dots, n - 1\}$ .

For one Monte Carlo simulation with a sample size  $n = 4096$ , we display the estimation results in Figure 2(b). Plots of the excess risk  $\ell(s_*, \hat{s}_m)$  against the dimension  $D_m$  are plotted in Figure 2(c). The curve  $\text{crit}_{2FCV}(m)$  and  $\text{crit}_{pen2F}(m)$  are also displayed in Figure 2(c). It can be observed that  $\text{crit}_{2FCV}(m)$  and  $\text{crit}_{pen2F}(m)$  give very reliable estimate for the risk  $\ell(s_*, \hat{s}_m)$ , and in turn, also a high-quality estimate of the optimal model. Indeed, in this case, both methods consistently select the oracle model  $m_*$ .

## 6 Proofs

As a preliminary result, let us first prove the consistency in sup-norm of our least squares estimators. This is in fact the main change compared to the strongly localized case treated in [16].

**Theorem 3.** *Let  $\alpha > 0$ . Assume that  $m$  is a linear vector space satisfying Assumption (A1b) and use the notations given in the statement of (A1b). Assume also that Assumption (A2b) holds. If there exists  $A_+ > 0$  such that*



**Fig. 2** (a): Noisy observations. (b): Typical reconstructions from a single simulation with  $n = 4096$ . Dashed line indicates the true function  $s_*$ , solid line corresponds to the estimates  $\widehat{s}_{\widehat{m}_{2FCV}}$  and dashed-dotted line to  $\widehat{s}_{\widehat{m}_{pen2F}}$ . (c): Graph of the excess risk  $\ell(s_*, \widehat{s}_m)$  (black) against the dimension  $D_m$  and (rescaled)  $\text{crit}_{2FCV}(m)$  (gray) and  $\text{crit}_{pen2F}(m)$  (light-gray) (in a log-log scale). The gray circle represents the global minimizer  $\widehat{m}_{2FCV}$  of  $\text{crit}_{2FCV}(m)$ , the light-gray diamond corresponds to the global minimizer  $\widehat{m}_{2FCV}$  of  $\text{crit}_{pen2F}(m)$  and the black star the oracle model  $m_*$ .

$$D_m \leq A + \frac{n^{1/3}}{(\ln n)^2},$$

then there exists a positive constant  $L_{A,r,\mathcal{M},\alpha}$  such that, for all  $n \geq n_0(r,\mathcal{M},\alpha)$ ,

$$\mathbb{P}\left(\|\widehat{s}_m - s_m\|_\infty \geq L_{A,r,\mathcal{M},\alpha} \sqrt{\frac{D_m \ln n}{n}}\right) \leq n^{-\alpha}.$$

*Proof (Proof of Theorem 3).* Let  $C > 0$ . Set

$$\mathcal{F}_C^\infty := \{s \in m; \|s - s_m\|_\infty \leq C\}$$

and

$$\mathcal{F}_{>C}^\infty := \{s \in m ; \|s - s_m\|_\infty > C\} = m \setminus \mathcal{F}_C^\infty .$$

Take an orthonormal basis  $(\varphi_k)_{k=1}^{D_m}$  of  $(m, \|\cdot\|_2)$  satisfying **(Alb)**. By Lemma 19 of [16], we get that there exists  $L_{A,r_m,\alpha}^{(1)} > 0$  such that, by setting

$$\Omega_1 = \left\{ \max_{k \in \{1, \dots, D_m\}} |(P_n - P)(\psi_m \cdot \varphi_k)| \leq L_{A,r_m,\alpha}^{(1)} \sqrt{\frac{\ln n}{n}} \right\} ,$$

we have for all  $n \geq n_0(A_+)$ ,  $\mathbb{P}(\Omega_1) \geq 1 - n^{-\alpha}$ . Moreover, we set

$$\Omega_2 = \left\{ \max_{(k,l) \in \{1, \dots, D_m\}^2} |(P_n - P)(\varphi_k \cdot \varphi_l)| \leq L_{\alpha,r_m}^{(2)} \min\{\|\varphi_k\|_\infty; \|\varphi_l\|_\infty\} \sqrt{\frac{\ln n}{n}} \right\} ,$$

where  $L_{\alpha,r_m}^{(2)}$  is defined in Lemma 18 of [16]. By Lemma 18 of [16], we have that for all  $n \geq n_0(A_+)$ ,  $\mathbb{P}(\Omega_2) \geq 1 - n^{-\alpha}$  and so, for all  $n \geq n_0(A_+)$ ,

$$\mathbb{P}(\Omega_1 \cap \Omega_2) \geq 1 - 2n^{-\alpha} .$$

We thus have for all  $n \geq n_0(A_+)$ ,

$$\begin{aligned} & \mathbb{P}(\|s_n - s_m\|_\infty > C) \\ & \leq \mathbb{P}\left(\inf_{s \in \mathcal{F}_{>C}^\infty} P_n(\gamma(s) - \gamma(s_m)) \leq \inf_{s \in \mathcal{F}_C^\infty} P_n(\gamma(s) - \gamma(s_m))\right) \\ & = \mathbb{P}\left(\sup_{s \in \mathcal{F}_{>C}^\infty} P_n(\gamma(s_m) - \gamma(s)) \geq \sup_{s \in \mathcal{F}_C^\infty} P_n(\gamma(s_m) - \gamma(s))\right) \\ & \leq \mathbb{P}\left(\left\{\sup_{s \in \mathcal{F}_{>C}^\infty} P_n(\gamma(s_m) - \gamma(s)) \geq \sup_{s \in \mathcal{F}_{C/2}^\infty} P_n(\gamma(s_m) - \gamma(s))\right\} \cap \Omega_1 \cap \Omega_2\right) + 2n^{-\alpha} . \end{aligned} \tag{5}$$

Now, for any  $s \in m$  such that

$$s - s_m = \sum_{k=1}^{D_m} \beta_k \varphi_k, \quad \beta = (\beta_k)_{k=1}^{D_m} \in \mathbb{R}^{D_m},$$

we have

$$\begin{aligned} & P_n(\gamma(s_m) - \gamma(s)) \\ & = (P_n - P)(\psi_m \cdot (s_m - s)) - (P_n - P)\left((s - s_m)^2\right) - P(\gamma(s) - \gamma(s_m)) \\ & = \sum_{k=1}^{D_m} \beta_k (P_n - P)(\psi_m \cdot \varphi_k) - \sum_{k,l=1}^{D_m} \beta_k \beta_l (P_n - P)(\varphi_k \cdot \varphi_l) - \sum_{k=1}^{D_m} \beta_k^2 . \end{aligned}$$

We set for any  $(k, l) \in \{1, \dots, D_m\}^2$ ,

$$R_{n,k}^{(1)} = (P_n - P)(\Psi_m \cdot \varphi_k) \quad \text{and} \quad R_{n,k,l}^{(2)} = (P_n - P)(\varphi_k \cdot \varphi_l) .$$

Moreover, we set a function  $h_n$ , defined as follows,

$$h_n : \beta = (\beta_k)_{k=1}^{D_m} \mapsto \sum_{k=1}^{D_m} \beta_k R_{n,k}^{(1)} - \sum_{k,l=1}^{D_m} \beta_k \beta_l R_{n,k,l}^{(2)} - \sum_{k=1}^{D_m} \beta_k^2 .$$

We thus have for any  $s \in m$  such that  $s - s_m = \sum_{k=1}^{D_m} \beta_k \varphi_k$ ,  $\beta = (\beta_k)_{k=1}^{D_m} \in \mathbb{R}^{D_m}$ ,

$$P_n(\gamma(s_m) - \gamma(s)) = h_n(\beta) . \quad (6)$$

In addition we set for any  $\beta = (\beta_k)_{k=1}^{D_m} \in \mathbb{R}^{D_m}$ ,

$$|\beta|_{m,\infty} = r_m \sqrt{D_m} |\beta|_\infty .$$

It is straightforward to see that  $|\cdot|_{m,\infty}$  is a norm on  $\mathbb{R}^{D_m}$ , proportional to the sup-norm. We also set for a real  $D_m \times D_m$  matrix  $B$ , its operator norm  $\|A\|_m$  associated to the norm  $|\cdot|_{m,\infty}$  on the  $D_m$ -dimensional vectors. More explicitly, we set for any  $B \in \mathbb{R}^{D_m \times D_m}$ ,

$$\|B\|_m := \sup_{\beta \in \mathbb{R}^{D_m}, \beta \neq 0} \frac{|B\beta|_{m,\infty}}{|\beta|_{m,\infty}} = \sup_{\beta \in \mathbb{R}^{D_m}, \beta \neq 0} \frac{|B\beta|_\infty}{|\beta|_\infty} .$$

We have, for any  $B = (B_{k,l})_{k,l=1,\dots,D_m} \in \mathbb{R}^{D_m \times D_m}$ , the following classical formula

$$\|B\|_m = \max_{k \in \{1, \dots, D_m\}} \left\{ \left\{ \sum_{l \in \{1, \dots, D_m\}} |B_{k,l}| \right\} \right\} .$$

Notice that by inequality (3) of **(A1b)**, it holds

$$\mathcal{F}_{>C}^\infty \subset \left\{ s \in m ; s - s_m = \sum_{k=1}^{D_m} \beta_k \varphi_k \quad \& \quad |\beta|_{m,\infty} \geq C \right\} \quad (7)$$

and

$$\mathcal{F}_{C/2}^\infty \supset \left\{ s \in m ; s - s_m = \sum_{k=1}^{D_m} \beta_k \varphi_k \quad \& \quad |\beta|_{m,\infty} \leq C/2 \right\} . \quad (8)$$

Hence, from (5), (6) (8) and (7) we deduce that if we find on  $\Omega_1 \cap \Omega_2$  a value of  $C$  such that

$$\sup_{\beta \in \mathbb{R}^{D_m}, |\beta|_{m,\infty} \geq C} h_n(\beta) < \sup_{\beta \in \mathbb{R}^{D_m}, |\beta|_{m,\infty} \leq C/2} h_n(\beta) ,$$

then we will get

$$\mathbb{P}(\|\hat{s}_m - s_m\|_\infty > C) \leq 2n^{-\alpha} .$$

Taking the partial derivatives of  $h_n$  with respect to the coordinates of its arguments, it then holds for any  $(k, l) \in \{1, \dots, D_m\}^2$  and  $\beta = (\beta_i)_{i=1}^{D_m} \in \mathbb{R}^{D_m}$ ,

$$\frac{\partial h_n}{\partial \beta_k}(\beta) = R_{n,k}^{(1)} - 2 \sum_{i=1}^{D_m} \beta_i R_{n,k,i}^{(2)} - 2\beta_k \quad (9)$$

We look now at the set of solutions  $\beta$  of the following system,

$$\frac{\partial h_n}{\partial \beta_k}(\beta) = 0, \forall k \in \{1, \dots, D_m\}. \quad (10)$$

We define the  $D_m \times D_m$  matrix  $R_n^{(2)}$  to be

$$R_n^{(2)} := \left( R_{n,k,l}^{(2)} \right)_{k,l=1,\dots,D_m}$$

and by (9), the system given in (10) can be written

$$2 \left( I_{D_m} + R_n^{(2)} \right) \beta = R_n^{(1)}, \quad (\text{S})$$

where  $R_n^{(1)}$  is a  $D_m$ -dimensional vector defined by

$$R_n^{(1)} = \left( R_{n,k}^{(1)} \right)_{k=1,\dots,D_m}.$$

Let us give an upper bound of the norm  $\left\| R_n^{(2)} \right\|_m$ , in order to show that the matrix  $I_{D_m} + R_n^{(2)}$  is nonsingular. On  $\Omega_2$  we have,

$$\begin{aligned} \left\| R_n^{(2)} \right\|_m &= \max_{k \in \{1, \dots, D_m\}} \left\{ \left\{ \sum_{l \in \{1, \dots, D_m\}} |(P_n - P)(\varphi_k \cdot \varphi_l)| \right\} \right\} \\ &\leq L_{\alpha, r_m}^{(2)} \max_{k \in \{1, \dots, D_m\}} \left\{ \left\{ \sum_{l \in \{1, \dots, D_m\}} \min \{ \|\varphi_k\|_\infty; \|\varphi_l\|_\infty \} \sqrt{\frac{\ln n}{n}} \right\} \right\} \\ &\leq r_m L_{\alpha, r_m}^{(2)} \sqrt{\frac{D_m^3 \ln n}{n}} \end{aligned} \quad (11)$$

Hence, from (11) and the fact that  $D_m \leq A_+ \frac{n^{1/3}}{(\ln n)^2}$ , we get that for all  $n \geq n_0(r_m, \alpha)$ , it holds on  $\Omega_2$ ,

$$\left\| R_n^{(2)} \right\|_m \leq \frac{1}{2}$$

and the matrix  $(I_d + R_n^{(2)})$  is nonsingular, of inverse  $(I_d + R_n^{(2)})^{-1} = \sum_{u=0}^{+\infty} (-R_n^{(2)})^u$ . Hence, the system (S) admits a unique solution  $\beta^{(n)}$ , given by

$$\beta^{(n)} = \frac{1}{2} \left( I_d + R_n^{(2)} \right)^{-1} R_n^{(1)}.$$

Now, on  $\Omega_1$  we have,

$$\left| R_n^{(1)} \right|_{m,\infty} \leq r_m \sqrt{D_m} \max_{k \in \{1, \dots, D_m\}} |(P_n - P)(\psi_m \cdot \varphi_k)| \leq r_m L_{A,r_m,\alpha}^{(1)} \sqrt{\frac{D_m \ln n}{n}}$$

and we deduce that for all  $n_0(r_m, \alpha)$ , it holds on  $\Omega_2 \cap \Omega_1$ ,

$$\left| \beta^{(n)} \right|_{m,\infty} \leq \frac{1}{2} \left\| \left( I_d + R_n^{(2)} \right)^{-1} \right\|_{m,m} \left| R_n^{(1)} \right|_{m,\infty} \leq r_m L_{A,r_m,\alpha}^{(1)} \sqrt{\frac{D_m \ln n}{n}}. \quad (12)$$

Moreover, by the formula (6) we have

$$h_n(\beta) = P_n(\gamma(s_m)) - P_n \left( Y - \sum_{k=1}^{D_m} \beta_k \varphi_k \right)^2$$

and we thus see that  $h_n$  is concave. Hence, for all  $n_0(r_m, \alpha)$ , we get that on  $\Omega_2$ ,  $\beta^{(n)}$  is the unique maximum of  $h_n$  and on  $\Omega_2 \cap \Omega_1$ , by (12), concavity of  $h_n$  and uniqueness of  $\beta^{(n)}$ , we get

$$h_n(\beta^{(n)}) = \sup_{\beta \in \mathbb{R}^{D_m}, |\beta|_{m,\infty} \leq C/2} h_n(\beta) > \sup_{\beta \in \mathbb{R}^{D_m}, |\beta|_{m,\infty} \geq C} h_n(\beta),$$

with  $C = 2r_m L_{A,r_m,\alpha}^{(1)} \sqrt{\frac{D_m \ln n}{n}}$ , which concludes the proof.

From Theorem 2 of [17] and Theorem 3 above, we deduce the following excess risks bounds.

**Theorem 4.** *Let  $A_+, A_-, \alpha > 0$ . Assume that  $m$  is a linear vector space of finite dimension  $D_m$  satisfying  $(\mathbf{Alb}(m))$  and use notations of  $(\mathbf{Alb}(m))$ . Assume, moreover, that the following assumption holds:*

$(\mathbf{Ab}(m))$  *There exists a constant  $A > 0$ , such that  $\|s_m\|_\infty \leq A$  and  $|Y| \leq A$  a.s.*

*If it holds*

$$A_- (\ln n)^2 \leq D_m \leq A_+ \frac{n^{1/3}}{(\ln n)^2},$$

*then a positive constant  $A_0$  exists, only depending on  $\alpha, A_-$  and on the constants  $A, \sigma_{\min}$  and  $r_m$  such that by setting*

$$\varepsilon_n = A_0 \max \left\{ \left( \frac{\ln n}{D_m} \right)^{1/4}, \left( \frac{D_m \ln n}{n} \right)^{1/4} \right\},$$

*we have for all  $n \geq n_0(A_-, A_+, A, r_m, \sigma_{\min}, \alpha)$ ,*

$$\begin{aligned} \mathbb{P} \left[ (1 - \varepsilon_n) \frac{\mathcal{C}_m}{n} \leq \ell(s_m, \widehat{s}_m) \leq (1 + \varepsilon_n) \frac{\mathcal{C}_m}{n} \right] &\geq 1 - 10n^{-\alpha}, \\ \mathbb{P} \left[ (1 - \varepsilon_n^2) \frac{\mathcal{C}_m}{n} \leq \ell_{\text{emp}}(\widehat{s}_m, s_m) \leq (1 + \varepsilon_n^2) \frac{\mathcal{C}_m}{n} \right] &\geq 1 - 5n^{-\alpha}, \end{aligned}$$

where  $\mathcal{C}_m = \sum_{k=1}^{D_m} \text{Var}((Y - s_m(X)) \cdot \varphi_k(X))$ .

Having at hand Theorem 4, the proofs of Theorems 1 and 4 follow from the exact same lines as the proofs of Theorems 6 and 7 of [16]. To give a more precise view of the ideas involved, let us detail the essential arguments of the proof of Theorem 1.

We set

$$\text{crit}_{\text{VFCV}}^0(m) = \text{crit}_{\text{VFCV}}(m) - \frac{1}{V} \sum_{j=1}^V P_n^{(j)}(\gamma(s_*)).$$

The difference between  $\text{crit}_{\text{VFCV}}^0(m)$  and  $\text{crit}_{\text{VFCV}}(m)$  being a quantity independent of  $m \in \mathcal{M}_n$ , the procedure defined by  $\text{crit}_{\text{VFCV}}^0$  gives the same result as the VFCV procedure defined by  $\text{crit}_{\text{VFCV}}$ .

We get for all  $m \in \mathcal{M}_n$ ,

$$\begin{aligned} \text{crit}_{\text{VFCV}}^0(m) &= \frac{1}{V} \sum_{j=1}^V P_n^{(j)} \left( \gamma(\widehat{s}_m^{(-j)}) - \gamma(s_*) \right) \\ &= \frac{1}{V} \sum_{j=1}^V \left[ P_n^{(j)} \left( \gamma(\widehat{s}_m^{(-j)}) - \gamma(s_m) \right) \right. \\ &\quad \left. + \left( P_n^{(j)} - P \right) (\gamma(s_m) - \gamma(s_*)) + P (\gamma(s_m) - \gamma(s_*)) \right] \\ &= \ell(s_*, \widehat{s}_m^{(-1)}) + \Delta_V(m) + \bar{\delta}(m) \end{aligned} \tag{13}$$

where

$$\Delta_V(m) = \frac{1}{V} \sum_{j=1}^V P_n^{(j)} \left( \gamma(\widehat{s}_m^{(-j)}) - \gamma(s_m) \right) - P \left( \gamma(\widehat{s}_m^{(-1)}) - \gamma(s_m) \right),$$

and

$$\bar{\delta}(m) = \frac{1}{V} \sum_{j=1}^V \left( P_n^{(j)} - P \right) (\gamma(s_m) - \gamma(s_*))$$

Now, we have to show that  $\Delta_V(m)$  and  $\bar{\delta}(m)$  are negligible in front of  $\ell(s_*, \widehat{s}_m^{(-1)})$ .

For  $\bar{\delta}(m)$ , this is done by using Bernstein's concentration inequality (see Lemma 7.5 of [16]). To control  $\Delta_V(m)$ , we also make use of Bernstein's concentration inequality, but by conditioning successively on the data used to learn the estimators  $\widehat{s}_m^{(-j)}$ ,  $j = 1, \dots, V$  (see Lemma 7.3 and Corollary 7.4 of [16]).

## References

1. Antoniadis, A., Gregoire, G., McKeague, I.: Wavelet methods for curve estimation. *J. Amer. Statist. Assoc.* **89**(428), 1340–1353 (1994)
2. Arlot, S.:  $V$ -fold cross-validation improved:  $V$ -fold penalization (2008). URL <http://hal.archives-ouvertes.fr/hal-00239182/en/>. ArXiv:0802.0566v2
3. Arlot, S., Céliste, A.: A survey of cross-validation procedures for model selection. *Stat. Surv.* **4**, 40–79 (2010)
4. Arlot, S., Massart, P.: Data-driven calibration of penalties for least-squares regression. *J. Mach. Learn. Res.* **10**, 245–279 (electronic) (2009)
5. Barron, A., Birgé, L., Massart, P.: Risk bounds for model selection via penalization. *Probab. Theory Related Fields* **113**(3), 301–413 (1999)
6. Birgé, L., Massart, P.: Minimum contrast estimators on sieves: exponential bounds and rates of convergence. *Bernoulli* **4**(3), 329–375 (1998)
7. Cai, T.: Adaptive wavelet estimation: a block thresholding and oracle inequality approach. *Ann. Statist.* **27**(3), 898–924 (1999)
8. Cai, T., Brown, L.: Wavelet shrinkage for nonequispaced samples. *Ann. Statist.* **26**, 1783–1799 (1998)
9. Cai, T., Brown, L.: Wavelet estimation for samples with random uniform design. *Statist. Probab. Lett.* **42**(3), 313–321 (1999)
10. Donoho, D., Maleki, A., Shahrampour, M.: Wavelab 850 (2006). URL <http://statweb.stanford.edu/wavelab/>
11. Hall, P., Turlach, B.: Interpolation methods for nonlinear wavelet regression with irregularly spaced design. *Ann. Statist.* **25**(5), 1912–1925 (1997)
12. Kulik, R., Raimondo, M.: Wavelet regression in random design with heteroscedastic dependent errors. *Ann. Statist.* **37**(6A), 3396–3430 (2009)
13. Mallat, S.: *A wavelet tour of signal processing: the sparse way*. Academic press (2008)
14. Marron, J., Adak, S., Johnstone, I., Neumann, M., Patil, P.: Exact risk analysis of wavelet regression. *J. Comput. Graph. Statist.* **7**(3), 278–309 (1998)
15. Nason, G.: Wavelet shrinkage using cross-validation. *J. R. Stat. Soc. Ser. B* pp. 463–479 (1996)
16. Navarro, F., Saumard, A.: Slope heuristics and  $v$ -fold model selection in heteroscedastic regression using strongly localized bases. *ESAIM: Probability and Statistics* (in press)
17. Saumard, A.: Optimal upper and lower bounds for the true and empirical excess risks in heteroscedastic least-squares regression. *Electron. J. Statist.* **6**(1-2), 579–655 (2012)
18. Saumard, A.: Optimal model selection in heteroscedastic regression using piecewise polynomial functions. *Electron. J. Statist.* **7**, 1184–1223 (2013)
19. Saumard, A.: On optimality of empirical risk minimization in linear aggregation. *Bernoulli* (2017). To appear, arXiv:1605.03433