



**HAL**  
open science

## Estimator selection in the Gaussian setting

Yannick Baraud, Christophe Giraud, Sylvie Huet

► **To cite this version:**

Yannick Baraud, Christophe Giraud, Sylvie Huet. Estimator selection in the Gaussian setting. Annales de l'Institut Henri Poincaré (B) Probabilités et Statistiques, 2014. hal-00502156v2

**HAL Id: hal-00502156**

**<https://hal.science/hal-00502156v2>**

Submitted on 21 Jun 2011

**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.

# Estimator selection in the Gaussian setting

Yannick Baraud, Christophe Giraud and Sylvie Huet

**Abstract:** We consider the problem of estimating the mean  $f$  of a Gaussian vector  $Y$  with independent components of common unknown variance  $\sigma^2$ . Our estimation procedure is based on estimator selection. More precisely, we start with an arbitrary and possibly infinite collection  $\mathbb{F}$  of estimators of  $f$  based on  $Y$  and, with the same data  $Y$ , aim at selecting an estimator among  $\mathbb{F}$  with the smallest Euclidean risk. No assumptions on the estimators are made and their dependencies with respect to  $Y$  may be unknown. We establish a non-asymptotic risk bound for the selected estimator. As particular cases, our approach allows to handle the problems of aggregation and model selection as well as those of choosing a window and a kernel for estimating a regression function, or tuning the parameter involved in a penalized criterion. We also derive oracle-type inequalities when  $\mathbb{F}$  consists of linear estimators. For illustration, we carry out two simulation studies. One aims at comparing our procedure to cross-validation for choosing a tuning parameter. The other shows how to implement our approach to solve the problem of variable selection in practice.

**AMS 2000 subject classifications:** Primary 62J05; secondary 62J07, 62G05, 62G08, 62F07.

**Keywords and phrases:** Gaussian linear regression, Estimator selection, Model selection, Variable selection, Linear estimator, Kernel estimator, Ridge regression, Lasso, Elastic net, Random Forest, PLS1 regression.

## 1. Introduction

### 1.1. *The setting and the approach*

We consider the Gaussian regression framework

$$Y_i = f_i + \varepsilon_i, \quad i = 1, \dots, n$$

where  $f = (f_1, \dots, f_n)$  is an unknown vector of  $\mathbb{R}^n$  and the  $\varepsilon_i$  are independent centered Gaussian random variables with common variance  $\sigma^2$ . Throughout the paper,  $\sigma^2$  is assumed to be unknown which corresponds to the practical case. Our aim is to estimate  $f$  from the observation of  $Y$ . For specific forms of  $f$ , this setting allows to deal simultaneously with the following problems.

**Example 1** (Signal denoising). *The vector  $f$  is of the form*

$$f = (F(x_1), \dots, F(x_n))$$

where  $x_1, \dots, x_n$  are distinct points of a set  $\mathcal{X}$  and  $F$  is an unknown mapping from  $\mathcal{X}$  into  $\mathbb{R}$ .

**Example 2** (Linear regression). *The vector  $f$  is assumed to be of the form*

$$f = X\beta \tag{1}$$

where  $X$  is a  $n \times p$  matrix,  $\beta$  is an unknown  $p$ -dimensional vector and  $p$  some integer larger than 1 (and possibly larger than  $n$ ). The columns of the matrix  $X$  are usually called predictors. When  $p$  is large, one may assume that the decomposition (1) is sparse in the sense that only few  $\beta_j$  are non-zero. Estimating  $f$  or finding the predictors associated to the non-zero coordinates of  $\beta$  are classical issues. The latter is called variable selection.

Our estimation strategy is based on estimator selection. More precisely, we start with an arbitrary collection  $\mathbb{F} = \{\hat{f}_\lambda, \lambda \in \Lambda\}$  of estimators of  $f$  based on  $Y$  and aim at selecting the one with the smallest Euclidean risk by using the same observation  $Y$ . The way the estimators  $\hat{f}_\lambda$  depend on  $Y$  may be arbitrary and possibly unknown. For example, the  $\hat{f}_\lambda$  may be obtained from the minimization of a criterion, a Bayesian procedure or the guess of some experts.

## 1.2. The motivation

The problem of choosing some best estimator among a family of candidate ones is central in Statistics. Let us present some examples.

**Example 3** (Choosing a tuning parameter). *Many statistical procedures depend on a (possibly multi-dimensional) parameter  $\lambda$  that needs to be tuned in view of obtaining an estimator with the best possible performance. For*

example, in the context of linear regression as described in Example 2, the Lasso estimator (see Tibshirani (1996) and Chen et al. (1998)) defined by  $\widehat{f}_\lambda = X\widehat{\beta}_\lambda$  with

$$\widehat{\beta}_\lambda = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left[ \|Y - X\beta\|^2 + \lambda \sum_{j=1}^p |\beta_j| \right]$$

depends on the choice of the parameter  $\lambda \geq 0$ . Selecting this parameter among a grid  $\Lambda \subset \mathbb{R}_+$  amounts to selecting a (suitable) estimator among the family  $\mathbb{F} = \{\widehat{f}_\lambda, \lambda \in \Lambda\}$ .

Another dilemma for Statisticians is the choice of a procedure to solve a given problem. In the context of Example 3, there exist many competitors to the Lasso estimator and one may alternatively choose a procedure based on ridge regression (see Hoerl and Kennard (1970)), random forest or PLS (see Tenenhaus (1998), Helland (2001) and Helland (2006)). Similarly, for the problem of signal denoising as described in Example 1, popular approaches include spline smoothing, wavelet decompositions and kernel estimators. The choice of a kernel may be possibly tricky.

**Example 4** (Choosing a kernel). Consider the problem described in Example 1 with  $\mathcal{X} = \mathbb{R}$ . For a kernel  $K$  and a bandwidth  $h > 0$ , the Nadaraya-Watson estimator (see Nadaraya (1964) and Watson (1964))  $\widehat{f}_{K,h} \in \mathbb{R}^n$  is defined as

$$\widehat{f}_{K,h} = \left( \widehat{F}_{K,h}(x_1), \dots, \widehat{F}_{K,h}(x_n) \right)$$

where for  $x \in \mathbb{R}$

$$\widehat{F}_{K,h}(x) = \frac{\sum_{j=1}^n K\left(\frac{x-x_j}{h}\right) Y_j}{\sum_{j=1}^n K\left(\frac{x-x_j}{h}\right)}.$$

There exist many possible choices for the kernel  $K$ , such as the Gaussian kernel  $K(x) = e^{-x^2/2}$ , the uniform kernel  $K(x) = \mathbf{1}_{|x|<1}$ , etc. Given a (finite) family  $\mathcal{K}$  of candidate kernels  $K$  and a grid  $\mathcal{H} \subset \mathbb{R}_+^*$  of possible values of  $h$ , one may consider the problem of selecting the best kernel estimator among the family  $\mathbb{F} = \{\widehat{f}_\lambda, \lambda = (K, h) \in \mathcal{K} \times \mathcal{H}\}$ .

### 1.3. A look at the literature

A common way to address the above issues is to use some cross-validation scheme such as leave-one-out or  $V$ -fold. Even though these resampling techniques are widely used in practice, little is known on their theoretical performances. For more details, we refer to Arlot and Celisse (2010) for a survey on cross-validation technics applied to model selection. Compared to these approaches, as we shall see, the procedure we propose is less time consuming and easier to implement. Moreover, it does not require to know how the estimators depend on the data  $Y$  and we can therefore handle the following problem.

**Example 5** (Selecting among mute experts). *A Statistician is given a collection  $\mathbb{F} = \{\widehat{f}_\lambda, \lambda \in \Lambda\}$  of estimators from a family  $\Lambda$  of experts  $\lambda$ , each of which keeping secret the way his/her estimator  $\widehat{f}_\lambda$  depends on the observation  $Y$ . The problem is then to find which expert  $\lambda$  is the closest to the truth.*

Given a selection rule among  $\mathbb{F}$ , an important issue is to compare the risk of the selected estimator to those of the candidate ones. Results in this direction are available in the context of model selection, which can be seen as a particular case of estimator selection. More precisely, for the purpose of selecting a suitable model one starts with a collection  $\mathbb{S}$  of those, typically linear spaces chosen for their approximation properties with respect to  $f$ , and one associates to each model  $S \in \mathbb{S}$  a suitable estimator  $\widehat{f}_S$  with values in  $S$ . Selecting a model then amounts to selecting an estimator among the collection  $\mathbb{F} = \{\widehat{f}_S, S \in \mathbb{S}\}$ . For this problem, selection rules based on the minimization of a penalized criterion have been proposed in the regression setting by Yang (1999), Baraud (2000), Birgé and Massart (2001) and Baraud et al (2009). Another way, usually called Lepski's method, appears in a series of papers by Lepski (1990; 1991; 1992a; 1992b) and was originally designed to perform model selection among collections of nested models. Finally, we mention that other procedures based on resampling have interestingly emerged from the work of Arlot (2007; 2009) and Céliste (2008). A common feature of those approaches lies in the fact that the proposed selection rules apply to specific collections of estimators only.

An alternative to *estimator selection* is *aggregation* which aims at designing a suitable combination of given estimators in order to outperform each of these separately (and even the best combination of these) up to a remaining term.

Aggregation techniques can be found in Catoni (1997; 2004), Juditsky and Nemirovski (2000), Nemirovski (2000), Yang (2000a), (2000b), (2001), Tsybakov (2003), Wegkamp (2003), Birgé (2006), Rigollet and Tsybakov (2007), Bunea, Tsybakov and Wegkamp (2007) and Goldenshluger (2009) for  $\mathbb{L}_p$ -losses. Most of the aggregation procedures are based on a sample splitting, one part of the data being used for building the estimators, the remaining part for selecting among these. Such a device requires that the observations be i.i.d. or at least that one has at disposal two independent copies of the data. From this point of view our procedure differs from classical *aggregation* procedures since we use the whole data  $Y$  to build and select. In the Gaussian regression setting that is considered here, we mention the results of Leung and Barron (2006) for the problem of mixing least-squares estimators. Their procedure uses the same data  $Y$  to estimate and to aggregate but requires the variance to be known. Giraud (2008) extends their results to the case where it is unknown.

#### 1.4. What is new here?

Our approach for solving the problem of estimator selection is new. We introduce a collection  $\mathbb{S}$  of linear subspaces of  $\mathbb{R}^n$  for approximating the estimators in  $\mathbb{F}$  and use a penalized criterion to compare them. As already mentioned and as we shall see, this approach requires no assumption on the family of estimators at hand and is easy to implement, an R-package being available on

[http://w3.jouy.inra.fr/unites/miaj/public/perso/SylvieHuet\\_en.html](http://w3.jouy.inra.fr/unites/miaj/public/perso/SylvieHuet_en.html).

A general way of comparing estimators in various statistical settings has been described in Baraud (2010). However, the procedure proposed there is mainly abstract and inadequate in the Gaussian framework we consider.

We prove a non-asymptotic risk bound for the estimator we select and show that this bound is optimal in the sense that it essentially cannot be improved (except for numerical constants maybe) by any other selection rule. For the sakes of illustration and comparison, we apply our procedure to various problems among which aggregation, model selection, variable selection and selection among linear estimators. In each of these cases, our approach allows to recover classical results in the areas as well as to establish new

ones. In the context of aggregation we compute the aggregation rates for the unknown variance case. These rates turn out to be the same as those for the known variance case. For selecting an estimator among a family of linear ones, we propose a new procedure and establish a risk bound which requires almost no assumption on the considered family. Finally, our approach provides a way of selecting a suitable variable selection procedure among a family of candidate ones. It thus provides an alternative to cross-validation for which little is known.

The paper is organized as follows. In Section 2 we present our selection rule and the theoretical properties of the resulting estimator. For illustration, we show in Sections 3, 4 and 5 respectively, how the procedure can be used to aggregate preliminary estimators, select a linear estimator among a finite collection of candidate ones, or solve the problem of variable selection. Section 6 is devoted to two simulation studies. One aims at comparing the performance of our procedure to the classical  $V$ -fold in view of selecting a tuning parameter among a grid. In the other, we evaluate the performance of the variable selection procedure we propose to some classical ones such as the Lasso, random forest, and others based on ridge and PLS regression. Finally, the proofs are postponed to Section 7.

Throughout the paper  $C$  denotes a constant that may vary from line to line.

## 2. The procedure and the main result

### 2.1. The procedure

Given a collection  $\mathbb{F} = \{\widehat{f}_\lambda, \lambda \in \Lambda\}$  of estimators of  $f$  based on  $Y$ , the selection rule we propose is based on the choices of a family  $\mathbb{S}$  of linear subspaces of  $\mathbb{R}^n$ , a collection  $\{\mathbb{S}_\lambda, \lambda \in \Lambda\}$  of (possibly random) subsets of  $\mathbb{S}$ , a weight function  $\Delta$  and a penalty function  $\text{pen}$ , both from  $\mathbb{S}$  into  $\mathbb{R}_+$ . We introduce those objects below and refer to Sections 3, 4 and 5 for examples.

#### 2.1.1. The collection of estimators $\mathbb{F}$

The collection  $\mathbb{F} = \{\widehat{f}_\lambda, \lambda \in \Lambda\}$  can be arbitrary. In particular,  $\mathbb{F}$  need not be finite nor countable and it may consist of a mix of estimators based on the

minimization of a criterion, a Bayes procedure or the guess of some experts. The dependency of these estimators with respect to  $Y$  need not be known. Nevertheless, we shall see on examples how we can use this information, when available, to improve the performance of our estimation procedure.

### 2.1.2. The families $\mathbb{S}$ and $\mathbb{S}_\lambda$

Let  $\mathbb{S}$  be a family of linear spaces of  $\mathbb{R}^n$  satisfying the following.

**Assumption 1.** *The family  $\mathbb{S}$  is finite or countable and for all  $S \in \mathbb{S}$ ,  $\dim(S) \leq n - 2$ .*

To each estimator  $\hat{f}_\lambda \in \mathbb{F}$ , we associate a (possibly random) subset  $\mathbb{S}_\lambda \subset \mathbb{S}$ .

Typically, the family  $\mathbb{S}$  should be chosen to possess good approximation properties with respect to the elements of  $\mathbb{F}$  and  $\mathbb{S}_\lambda$  with respect to  $\hat{f}_\lambda$  specifically. One may take  $\mathbb{S}_\lambda = \mathbb{S}$  but for computational reasons it will be convenient to allow  $\mathbb{S}_\lambda$  to be smaller. The choices of  $\mathbb{S}_\lambda$  may be made on the basis of the observation  $\hat{f}_\lambda$ . We provide examples of  $\mathbb{S}$  and  $\mathbb{S}_\lambda$  in various statistical settings described in Sections 3 to 5.

### 2.1.3. The weight function $\Delta$ and the associated function $\text{pen}_\Delta$

We consider a function  $\Delta$  from  $\mathbb{S}$  into  $\mathbb{R}_+$  and assume

**Assumption 2.**

$$\Sigma = \sum_{S \in \mathbb{S}} e^{-\Delta(S)} < +\infty. \quad (2)$$

Whenever  $\mathbb{S}$  is finite, inequality (2) automatically holds true. However, in practice  $\Sigma$  should be kept to a reasonable size. When  $\Sigma = 1$ ,  $e^{-\Delta(\cdot)}$  can be interpreted as a prior distribution on  $\mathbb{S}$  and gives thus a Bayesian flavor to the procedure we propose. To the weight function  $\Delta$ , we associate the function  $\text{pen}_\Delta$  mapping  $\mathbb{S}$  into  $\mathbb{R}_+$  and defined by

$$\mathbb{E} \left[ \left( U - \frac{\text{pen}_\Delta(S)}{n - \dim(S)} V \right)_+ \right] = e^{-\Delta(S)} \quad (3)$$

where  $x_+$  denotes the positive part of  $x \in \mathbb{R}$  and  $U, V$  are two independent  $\chi^2$  random variables with respectively  $\dim(S) + 1$  and  $n - \dim(S) - 1$  degrees



of freedom. This function can be easily computed from the quantiles of the Fisher distribution as we shall see in Section 8.1. From a more theoretical point of view, it is shown in Baraud *et al* (2009) that under Assumption 3 below, there exists a positive constant  $C$  (depending on  $\kappa$  only) such that

$$\text{pen}_\Delta(S) \leq C(\dim(S) \vee \Delta(S)). \quad (4)$$

**Assumption 3.** *There exists  $\kappa \in (0, 1)$  such that for all  $S \in \mathbb{S}$ ,*

$$1 \leq \dim(S) \vee \Delta(S) \leq \kappa n.$$

#### 2.1.4. The selection criterion

The selection procedure we propose involves a penalty function  $\text{pen}$  from  $\mathbb{S}$  into  $\mathbb{R}_+$  with the following property.

**Assumption 4.** *The penalty function  $\text{pen}$  satisfies for some  $K > 1$ ,*

$$\text{pen}(S) \geq K \text{pen}_\Delta(S) \quad \text{for all } S \in \mathbb{S}. \quad (5)$$

Whenever equality holds in (5), it derives from (4) that  $\text{pen}(S)$  measures the complexity of the model  $S$  in terms of dimension and weight.

Denoting  $\Pi_S$  the projection operator onto a linear space  $S \subset \mathbb{R}^n$ , given the families  $\mathbb{S}_\lambda$ , the penalty function  $\text{pen}$  and some positive number  $\alpha$ , we define

$$\text{crit}_\alpha(\hat{f}_\lambda) = \inf_{S \in \mathbb{S}_\lambda} \left[ \left\| Y - \Pi_S \hat{f}_\lambda \right\|^2 + \alpha \left\| \hat{f}_\lambda - \Pi_S \hat{f}_\lambda \right\|^2 + \text{pen}(S) \hat{\sigma}_S^2 \right], \quad (6)$$

where

$$\hat{\sigma}_S^2 = \frac{\|Y - \Pi_S Y\|^2}{n - \dim(S)}. \quad (7)$$

## 2.2. The main result

For all  $\lambda \in \Lambda$  let us set

$$A(\hat{f}_\lambda, \mathbb{S}_\lambda) = \inf_{S \in \mathbb{S}_\lambda} \left[ \left\| \hat{f}_\lambda - \Pi_S \hat{f}_\lambda \right\|^2 + \text{pen}(S) \hat{\sigma}_S^2 \right]. \quad (8)$$

This quantity corresponds to an accuracy index for the estimator  $\hat{f}_\lambda$  with respect to the family  $\mathbb{S}_\lambda$ . The following result holds.

**Theorem 1.** *Let  $K > 1, \alpha > 0, \delta \geq 0$ . Assume that Assumptions 1, 2 and 4 hold. There exists a constant  $C$  (given by (33)) depending on  $K$  and  $\alpha$  only such that for any  $\widehat{f}_{\widehat{\lambda}}$  in  $\mathbb{F}$  satisfying*

$$\text{crit}_{\alpha}(\widehat{f}_{\widehat{\lambda}}) \leq \inf_{\lambda \in \Lambda} \text{crit}_{\alpha}(\widehat{f}_{\lambda}) + \delta, \quad (9)$$

we have the following bounds

$$\begin{aligned} C\mathbb{E} \left[ \left\| f - \widehat{f}_{\widehat{\lambda}} \right\|^2 \right] &\leq \mathbb{E} \left[ \inf_{\lambda \in \Lambda} \left[ \left\| f - \widehat{f}_{\lambda} \right\|^2 + A(\widehat{f}_{\lambda}, \mathbb{S}_{\lambda}) \right] \right] + \Sigma\sigma^2 + \delta \quad (10) \\ &\leq \inf_{\lambda \in \Lambda} \left\{ \mathbb{E} \left[ \left\| f - \widehat{f}_{\lambda} \right\|^2 \right] + \mathbb{E} \left[ A(\widehat{f}_{\lambda}, \mathbb{S}_{\lambda}) \right] \right\} + \Sigma\sigma^2 + \delta \quad (11) \end{aligned}$$

(provided that the quantity involved in the expectation in (10) is measurable). Furthermore, if equality holds in (5) and Assumption 3 is satisfied, for each  $\lambda \in \Lambda$

- if the set  $\mathbb{S}_{\lambda}$  is non-random,

$$\begin{aligned} C'\mathbb{E} \left[ A(\widehat{f}_{\lambda}, \mathbb{S}_{\lambda}) \right] \\ \leq \mathbb{E} \left[ \left\| f - \widehat{f}_{\lambda} \right\|^2 \right] + \inf_{S \in \mathbb{S}_{\lambda}} \left[ \mathbb{E} \left[ \left\| \widehat{f}_{\lambda} - \Pi_S \widehat{f}_{\lambda} \right\|^2 \right] + (\dim(S) \vee \Delta(S))\sigma^2 \right] \quad (12) \end{aligned}$$

- if there exists a (possibly random) linear space  $\widehat{S}_{\lambda} \in \mathbb{S}_{\lambda}$  such that  $\widehat{f}_{\lambda} \in \widehat{S}_{\lambda}$  with probability 1,

$$C'\mathbb{E} \left[ A(\widehat{f}_{\lambda}, \mathbb{S}_{\lambda}) \right] \leq \mathbb{E} \left[ \left\| f - \widehat{f}_{\lambda} \right\|^2 \right] + \mathbb{E} \left[ \dim(\widehat{S}_{\lambda}) \vee \Delta(\widehat{S}_{\lambda}) \right] \sigma^2, \quad (13)$$

where  $C'$  is a positive constant only depending on  $\kappa$  and  $K$ .

Let us now comment Theorem 1.

It turns out that inequality (10) leaves no place for a substantial improvement in the sense that the bound we get is essentially optimal and cannot be improved (apart from constants) by any other selection rule among  $\mathbb{F}$ . To see this, let us assume for simplicity that  $\mathbb{F}$  is finite so that a measurable minimizer of  $\text{crit}_{\alpha}$  always exists and  $\delta$  can be chosen as 0. Let  $K = 1.1$ ,  $\alpha = 1/2$  (to fix up the ideas),  $\mathbb{S}$  a family of linear spaces satisfying the assumptions

of Theorem 1 and  $\text{pen}$ , the penalty function achieving equality in (5). Besides, assume that  $\mathbb{S}$  contains a linear space  $S$  such that  $1 \leq \dim(S) \leq n/2$  and associate to  $S$  the weight  $\Delta(S) = \dim(S)$ . If  $\mathbb{S}_\lambda = \mathbb{S}$  for all  $\lambda$ , we deduce from (4) and (10) that for some universal constant  $C'$ , whatever  $\mathbb{F}$  and  $f \in \mathbb{R}^n$

$$\begin{aligned} & C' \mathbb{E} \left[ \left\| f - \widehat{f}_{\widehat{\lambda}} \right\|^2 \right] \\ & \leq \mathbb{E} \left[ \inf_{\lambda \in \Lambda} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 + \inf_{S \in \mathbb{S}} \left( \left\| \widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda \right\|^2 + \text{pen}(S) \widehat{\sigma}_S^2 \right) \right] \right] \\ & \leq \mathbb{E} \left[ \inf_{\lambda \in \Lambda} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 + \left\| \widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda \right\|^2 + \dim(S) \widehat{\sigma}_S^2 \right] \right]. \end{aligned} \quad (14)$$

In the opposite direction, the following result holds.

**Proposition 1.** *There exists a universal constant  $C$ , such that for any finite family  $\mathbb{F} = \{\widehat{f}_\lambda, \lambda \in \Lambda\}$  of estimators and any selection rule  $\widehat{\lambda}$  based on  $Y$  among  $\Lambda$ , there exists  $f \in S$  such that*

$$C \mathbb{E} \left[ \left\| f - \widehat{f}_{\widehat{\lambda}} \right\|^2 \right] \geq \mathbb{E} \left[ \inf_{\lambda \in \Lambda} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 + \left\| \widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda \right\|^2 + \dim(S) \sigma^2 \right] \right]. \quad (15)$$

We see that, up to the estimator  $\widehat{\sigma}_S^2$  in place of  $\sigma^2$  and numerical constants, the left-hand sides of (14) and (15) coincide.

In view of commenting (11) further, we continue assuming that  $\mathbb{F}$  is finite so that we can keep  $\delta = 0$  in (11). A particular feature of (11) lies in the fact that the risk bound pays no price for considering a large collection  $\mathbb{F}$  of estimators. In fact, it is actually decreasing with respect to  $\mathbb{F}$  (or equivalently  $\Lambda$ ) for the inclusion. This means that if one adds a new estimator to the collection  $\mathbb{F}$  (without changing neither  $\mathbb{S}$  nor the families  $\mathbb{S}_\lambda$  associated to the former estimators), the risk bound for  $\widehat{f}_{\widehat{\lambda}}$  can only be improved. In contrast, the computation of the estimator  $\widehat{f}_{\widehat{\lambda}}$  is all the more difficult that  $|\mathbb{F}|$  is large. More precisely, if the cardinalities of the families  $\mathbb{S}_\lambda$  are not too large, the computation of  $\widehat{f}_{\widehat{\lambda}}$  requires around  $|\mathbb{F}|$  steps.

The selection rule we use does not require to know how the estimators depend on  $Y$ . In fact, as we shall see, a more important piece of information is the ranges of the estimators  $\widehat{f}_\lambda = \widehat{f}_\lambda(Y)$  as  $Y$  varies in  $\mathbb{R}^n$ . A situation of special interest occurs when each  $\widehat{f}_\lambda$  belongs to some (possibly random)

linear space  $\widehat{S}_\lambda$  in  $\mathbb{S}$  with probability one. By taking  $\mathbb{S}_\lambda$  such that  $\widehat{S}_\lambda \in \mathbb{S}_\lambda$  for all  $\lambda$ , we deduce from Theorem 1 by using (11) and (13) the following corollary.

**Corollary 1.** *Assume that the Assumptions of Theorem 1 are satisfied, that Assumption 3 holds and that equality holds in (5). If for all  $\lambda \in \Lambda$  there exists a (possibly random) linear space  $\widehat{S}_\lambda \in \mathbb{S}_\lambda$  such that  $\widehat{f}_\lambda \in \widehat{S}_\lambda$  with probability 1, then  $\widehat{f}_{\widehat{\lambda}}$  satisfies*

$$C\mathbb{E} \left[ \left\| f - \widehat{f}_{\widehat{\lambda}} \right\|^2 \right] \leq \inf_{\lambda \in \Lambda} \left[ \mathbb{E} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 \right] + \mathbb{E} \left[ \dim(\widehat{S}_\lambda) \vee \Delta(\widehat{S}_\lambda) \right] \sigma^2 \right] + \delta, \quad (16)$$

for some  $C$  depending on  $K$  and  $\kappa$  only.

One may apply this result in the context of model selection. One starts with a collection of models  $\mathbb{S} = \{S_m, m \in \mathcal{M}\}$  and associate to each  $S_m$  an estimator  $\widehat{f}_m$  with values in  $S_m$ . By taking  $\mathbb{F} = \{\widehat{f}_m, m \in \mathcal{M}\}$  (here  $\Lambda = \mathcal{M}$ ) and  $\mathbb{S}_m = \{S_m\}$  for all  $m \in \mathcal{M}$ , our selection procedure leads to an estimator  $\widehat{f}_{\widehat{m}}$  which satisfies

$$C\mathbb{E} \left[ \left\| f - \widehat{f}_{\widehat{m}} \right\|^2 \right] \leq \inf_{m \in \mathcal{M}} \left[ \mathbb{E} \left[ \left\| f - \widehat{f}_m \right\|^2 \right] + (\dim(S_m) \vee \Delta(S_m)) \sigma^2 \right]. \quad (17)$$

When  $\widehat{f}_m = \Pi_{S_m} Y$  for all  $m \in \mathcal{M}$ , our selection rule becomes

$$\widehat{m} = \arg \min_{m \in \mathcal{M}} \left[ \left\| Y - \widehat{f}_m \right\|^2 + \text{pen}(S_m) \widehat{\sigma}_{S_m}^2 \right] \quad (18)$$

and turns out to coincide with that described in Baraud *et al* (2009). Interestingly, Corollary 1 shows that this selection rule can still be used for families  $\mathbb{F}$  of (non-linear) estimators of the form  $\Pi_{S_{\widehat{m}}} Y$  where the  $S_{\widehat{m}}$  are chosen randomly among  $\mathbb{S}$  on the basis of  $Y$ , doing thus as if the linear spaces  $S_{\widehat{m}}$  were non-random. An estimator of the form  $\Pi_{S_{\widehat{m}}} Y$  can be interpreted as resulting from a model selection procedures among the family of projection estimators  $\{\Pi_m Y, m \in \mathcal{M}\}$  and hence, (18) can be used to choose some best model selection rule among a collection of candidate ones.

### 3. Aggregation

In this section, we consider the problems of *Model Selection Aggregation* (MS), *Convex Aggregation* (Cv) and *Linear Aggregation* (L) defined below.

Given  $M \geq 2$  preliminary estimators of  $f$ , denoted  $\{\phi_k, k = 1, \dots, M\}$ , our aim is to build an estimator  $\widehat{f}$  based on  $Y$  whose risk is as close as possible to  $\inf_{g \in \mathbb{F}_\Lambda} \|f - g\|^2$  where

$$\mathbb{F}_\Lambda = \left\{ f_\lambda = \sum_{j=1}^M \lambda_j \phi_j, \lambda \in \Lambda \right\}$$

and, according to the aggregation problem at hand,  $\Lambda$  is one of the three sets

$$\Lambda_{\text{MS}} = \left\{ \lambda \in \{0, 1\}^M, \sum_{j=1}^M \lambda_j = 1 \right\}, \Lambda_{\text{Cv}} = \left\{ \lambda \in \mathbb{R}_+^M, \sum_{j=1}^M \lambda_j = 1 \right\}, \Lambda_{\text{L}} = \mathbb{R}^M.$$

When  $\Lambda = \Lambda_{\text{MS}}$ ,  $\mathbb{F}_\Lambda$  is the set  $\{\phi_1, \dots, \phi_M\}$  consisting of the initial estimators. When  $\Lambda = \Lambda_{\text{Cv}}$ ,  $\mathbb{F}_\Lambda$  is the convex hull of the  $\phi_j$ . In the literature, one may also find

$$\Lambda'_{\text{Cv}} = \left\{ \lambda \in [0, 1]^M, \sum_{j=1}^M \lambda_j \leq 1 \right\}$$

in place of  $\Lambda_{\text{Cv}}$  in which case  $\mathbb{F}_\Lambda$  is the convex hull of  $\{0, \phi_1, \dots, \phi_M\}$ . Finally, when  $\Lambda = \Lambda_{\text{L}}$ ,  $\mathbb{F}_\Lambda$  is the linear span of the  $\phi_j$ .

Each of these three aggregation problems are solved *separately* if for each  $\Lambda \in \{\Lambda_{\text{MS}}, \Lambda_{\text{Cv}}, \Lambda_{\text{L}}\}$  one can design an estimator  $\widehat{f} = \widehat{f}(\Lambda)$  satisfying

$$\mathbb{E} \left[ \left\| f - \widehat{f} \right\|^2 \right] - C \inf_{g \in \mathbb{F}_\Lambda} \|f - g\|^2 \leq C' \psi_{n,\Lambda} \sigma^2 \quad (19)$$

with  $C = 1$ ,  $C' > 0$  free of  $f, n, M$  and

$$\psi_{n,\Lambda} = \begin{cases} M & \text{if } \Lambda = \Lambda_{\text{L}} \\ \sqrt{n \log(eM/\sqrt{n})} & \text{if } \Lambda = \Lambda_{\text{Cv}} \text{ and } \sqrt{n} \leq M \\ M & \text{if } \Lambda = \Lambda_{\text{Cv}} \text{ and } \sqrt{n} \geq M \\ \log M & \text{if } \Lambda = \Lambda_{\text{MS}}. \end{cases} \quad (20)$$

These problems have only been considered when the variance is known. The quantity  $\psi_{n,\Lambda}$  then corresponds to the best possible upper bound in (19) over all possible  $f \in \mathbb{R}^n$  and preliminary estimators  $\phi_j$  and is called the *optimal rate of aggregation*. For a more precise definition, we refer the reader to Tsybakov (2003). Bunea et al (2007) considered the problem of solving

these three problems *simultaneously* by building an estimator  $\widehat{f}$  which satisfies (19) simultaneously for all  $\Lambda \in \{\Lambda_{\text{MS}}, \Lambda_{\text{CV}}, \Lambda_{\text{L}}\}$  and some constant  $C > 1$ . This is an interesting issue since it is impossible to know in practice which aggregation device should be used to achieve the smallest risk bound: as  $\Lambda$  grows (for the inclusion), the bias  $\inf_{g \in \mathbb{F}_\Lambda} \|f - g\|^2$  decreases while the rate  $\psi_{n,\Lambda}$  increases.

The aim of this section is to show that our procedure provides a way of solving (or nearly solving) the three aggregation problems both *separately* and *simultaneously* when the variance is unknown.

Throughout this section, we consider the family  $\overline{\mathbb{S}}$  consisting of the  $S_m$  defined for each  $m \subset \{1, \dots, M\}$  and  $m \neq \emptyset$  as the linear span of the  $\phi_j$  for  $j \in m$ . Along this section, we shall use the weight function  $\Delta$  defined on  $\overline{\mathbb{S}}$  by

$$\Delta(S_m) = |m| + \log \left[ \binom{M}{|m|} \right],$$

take  $\alpha = 1/2$  and  $\text{pen}(\cdot) = 1.1\text{pen}_\Delta(\cdot)$  taking thus  $K = 1.1$ . The choices of  $\alpha$  and  $K$  is only to fix up the ideas. Note that  $\Delta$  satisfies Assumption 2 with  $\Sigma < 1$ . To avoid trivialities, we assume all along  $n \geq 4$ .

### 3.1. Solving the three aggregation problems separately

#### 3.1.1. Linear Aggregation

Problem (L) is the easiest to solve. Let us take  $\mathbb{F} = \mathbb{F}_\Lambda$  with  $\Lambda = \Lambda_{\text{L}}$  and

$$\mathbb{S} = \mathbb{S}_{\text{L}} = \{S_{\{1, \dots, M\}}\} \quad (21)$$

and  $\mathbb{S}_\lambda = \mathbb{S}_{\text{L}}$  for all  $\lambda \in \Lambda_{\text{L}}$ . Minimizing  $\text{crit}_\alpha(f_\lambda)$  over  $f_\lambda \in \mathbb{F}_\Lambda$  amounts to minimizing  $\|Y - f_\lambda\|^2$  over  $f_\lambda \in S_{\{1, \dots, M\}}$  and hence, the resulting estimator is merely  $\widehat{f}_{\text{L}} = \Pi_{S_{\{1, \dots, M\}}} Y$ . The risk of  $\widehat{f}_{\text{L}}$  satisfies

$$\mathbb{E} \left[ \left\| f - \widehat{f}_{\text{L}} \right\|^2 \right] \leq \inf_{g \in \mathbb{F}_\Lambda} \|f - g\|^2 + M\sigma^2.$$

whatever  $n$  and  $M$  which solves the problem of *Linear Aggregation*.

### 3.1.2. Model Selection Aggregation

To tackle Problem (MS), we take  $\mathbb{F} = \mathbb{F}_\Lambda$  with  $\Lambda = \Lambda_{\text{MS}}$ , that is,  $\mathbb{F}_\Lambda = \{\phi_1, \dots, \phi_M\}$ ,

$$\mathbb{S} = \mathbb{S}_{\text{MS}} = \{S_{\{1\}}, \dots, S_{\{M\}}\} \quad (22)$$

and associate to each  $f_\lambda = \phi_j$  the collection  $\mathbb{S}_\lambda$  reduced to  $\{S_{\{j\}}\}$ . Note that  $\dim(S) \leq 1$  and  $\Delta(S) = \log(eM) \geq \dim(S)$  for all  $S \in \mathbb{S}_{\text{MS}}$ , so that under the assumption that  $\log(eM) \leq n/2$  we may apply Corollary 1 with  $\delta = 0$  (since  $\mathbb{F}_\Lambda$  is finite),  $\kappa = 1/2$  and get that for some constant  $C > 0$  the resulting estimator  $\widehat{f}_{\text{MS}}$  satisfies

$$C\mathbb{E} \left[ \left\| f - \widehat{f}_{\text{MS}} \right\|^2 \right] \leq \inf_{g \in \mathbb{F}_\Lambda} \|f - g\|^2 + \log(M)\sigma^2.$$

This risk bound is of the form (19) except for the constant  $C$  which is not equal to 1. We do not know whether Problem (MS) can be solved or not with  $C = 1$  when the variance  $\sigma^2$  is unknown and  $M$  is large (possibly larger than  $n$ ).

### 3.1.3. Convex aggregation

For this problem, we emphasize the aggregation rate with respect to the quantity

$$L = \sup_{j=1, \dots, M} \frac{\|\phi_j\|}{\sigma\sqrt{n}}. \quad (23)$$

If  $M < \sqrt{n}L$ , take again the estimator  $\widehat{f}_L$ . Since the convex hull of the  $\phi_j$  is a subset of the linear space  $S_{\{1, \dots, M\}}$ , for  $\Lambda = \Lambda_{\text{Cv}}$  we have

$$\mathbb{E} \left[ \left\| f - \widehat{f}_L \right\|^2 \right] \leq \inf_{g \in \mathbb{F}_\Lambda} \|f - g\|^2 + M\sigma^2.$$

Let us now turn to the case  $M \geq \sqrt{n}L$ . More precisely, assume that

$$2 \leq \sqrt{n}L \leq M \leq e^{-1} \min \left\{ \sqrt{n}L e^{nL^2}, e^{\sqrt{n}/(2L)} \right\} \quad (24)$$

and set  $d(n, M) = n/(2 \log(eM))$ . We consider the family of estimators  $\mathbb{F} = \mathbb{F}_\Lambda$  with  $\Lambda = \Lambda_{\text{Cv}}$  and

$$\mathbb{S} = \mathbb{S}_{\text{Cv}} = \mathbb{S}_\lambda = \{S_m \in \overline{\mathbb{S}}, |m| \leq d(n, M)\}, \quad \forall \lambda \in \Lambda_{\text{Cv}}. \quad (25)$$

The set  $\Lambda_{\text{Cv}}$  being compact,  $\lambda \mapsto \text{crit}_\alpha(f_\lambda)$  admits a minimum  $\widehat{\lambda}$  over  $\Lambda_{\text{Cv}}$  and we set  $\widehat{f}_{\text{Cv}} = \widehat{f}_{\widehat{\lambda}}$ .

**Proposition 2.** *There exists a universal constant  $C > 1$  such that*

$$\mathbb{E} \left[ \left\| f - \widehat{f}_{\text{Cv}} \right\|^2 \right] - C \inf_{g \in \mathbb{F}_\Lambda} \|f - g\|^2 \leq C \sqrt{nL^2 \log(eM/\sqrt{nL^2})} \sigma^2.$$

This risk bound is of the form (19) except for the constant  $C$  which is not equal to 1. Again, we do not know whether Problem (Cv) can be solved or not with  $C = 1$  when the variance  $\sigma^2$  is unknown and  $M$  possibly larger than  $n$ .

### 3.2. Solving the three problems simultaneously

Consider now three estimators  $\widehat{f}_L, \widehat{f}_{\text{MS}}, \widehat{f}_{\text{Cv}}$  with values respectively in  $S_{\{1, \dots, M\}}, \bigcup_{j=1}^M S_{\{j\}}$  and the convex hull  $\mathcal{C}$  of the  $\phi_j$  (we use a new notation for this convex hull to avoid ambiguity). One may take the estimators defined in Section 3.1 but any others would suit. The aim of this section is to select the one with the smallest risk to estimate  $f$ . To do so, we apply our selection procedure with  $\mathbb{F} = \{\widehat{f}_L, \widehat{f}_{\text{MS}}, \widehat{f}_{\text{Cv}}\}$ , taking thus  $\Lambda = \{L, \text{MS}, \text{Cv}\}$ , and associate to each of these three estimators the families  $\mathbb{S}_L, \mathbb{S}_{\text{MS}}, \mathbb{S}_{\text{Cv}}$  defined by (21), (22) and (25) respectively and choose  $\mathbb{S} = \mathbb{S}_L \cup \mathbb{S}_{\text{MS}} \cup \mathbb{S}_{\text{Cv}}$ .

**Proposition 3.** *Assume that (24) holds and that  $\log(eM) \leq n/2$ . There exists a universal constant  $C > 0$  such that whatever  $\widehat{f}_L, \widehat{f}_{\text{MS}}$  and  $\widehat{f}_{\text{Cv}}$  with values in  $S_{\{1, \dots, M\}}, \bigcup_{j=1}^M S_{\{j\}}$  and  $\mathcal{C}$  respectively, the selected estimator  $\widehat{f}_{\widehat{\lambda}}$  satisfies for all  $f \in \mathbb{R}^n$ ,*

$$C \mathbb{E} \left[ \left\| f - \widehat{f}_{\widehat{\lambda}} \right\|^2 \right] \leq \inf_{\lambda \in \{L, \text{MS}, \text{Cv}\}} \left[ \mathbb{E} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 \right] + B_\lambda \right],$$

where

$$B_L = \sigma^2 M, \quad B_{\text{MS}} = \sigma^2 \log M, \quad B_{\text{Cv}} = \sigma^2 \left[ M \wedge \sqrt{nL^2 \log(eM/\sqrt{nL^2})} \right].$$

In particular, if  $\widehat{f}_L, \widehat{f}_{\text{MS}}$  and  $\widehat{f}_{\text{Cv}}$  fulfills (19), then

$$C \mathbb{E} \left[ \left\| f - \widehat{f}_{\widehat{\lambda}} \right\|^2 \right] \leq \inf_{\lambda \in \{L, \text{MS}, \text{Cv}\}} \left[ \inf_{g \in \mathbb{F}_\lambda} \|f - g\|^2 + B_\lambda \right],$$

where  $\mathbb{F}_\lambda$  stands for  $\mathbb{F}_\Lambda$  when  $\Lambda = \Lambda_\lambda$ .



#### 4. Selecting among linear estimator

In this section, we consider the situation where the estimators  $\widehat{f}_\lambda$  are linear, that is, are of the form  $\widehat{f}_\lambda = A_\lambda Y$  for some known and deterministic  $n \times n$  matrix  $A_\lambda$ . As mentioned before, this setting covers many popular estimation procedures including kernel ridge estimators, spline smoothing, Nadaraya estimators,  $\lambda$ -nearest neighbors, projection estimators, low-pass filters, etc. In some cases  $A_\lambda$  is symmetric (e.g. kernel ridge, spline smoothing, projection estimators), in some others  $A_\lambda$  is non-symmetric and non-singular (as for Nadaraya estimators) and sometimes  $A_\lambda$  can be both singular and non-symmetric (low pass filters,  $\lambda$ -nearest neighbors). A common feature of those procedures lies in the fact that they depend on a tuning parameter (possibly multidimensional) and their practical performances can be quite poor if this parameter is not suitably calibrated. A series of papers have investigated the calibration of some of these procedures. To mention a few of them, Cao and Golubev (2006) focus on spline smoothing, Zhang (2005) on kernel ridge regression, Goldenshluger and Lepski (2009) on kernel estimators and Arlot and Bach (2009) propose a procedure to select among symmetric linear estimator with spectrum in  $[0, 1]$ . The procedure we present can handle all these cases in an unified framework. Throughout the section, we assume that  $\Lambda$  is finite.

##### 4.1. The families $S_\lambda$

To apply our selection procedure, we need to associate to each  $A_\lambda$  a suitable collection of approximation spaces  $S_\lambda$ . To do so, we introduce below a linear space  $S_\lambda$  which plays a key role in our analysis.

For the sake of simplicity, let us first consider the case where  $A_\lambda$  is non-singular. Then  $S_\lambda$  is defined as the linear span of the right-singular vectors of  $A_\lambda^{-1} - I$  associated to singular values smaller than 1. When  $A_\lambda$  is symmetric,  $S_\lambda$  is merely the linear span of the eigenvectors of  $A_\lambda$  associated to eigenvalues not smaller than  $1/2$ . If none of the singular values are smaller than 1, then  $S_\lambda = \{0\}$ .

Let us now extend the definition of  $S_\lambda$  to singular operators  $A_\lambda$ . Let us recall that  $\mathbb{R}^n = \ker(A_\lambda) \oplus \text{rg}(A_\lambda^*)$  where  $A_\lambda^*$  stands for the transpose of  $A_\lambda$  and  $\text{rg}(A_\lambda^*)$  for its range. The operator  $A_\lambda$  then induces a one to one operator

between  $\text{rg}(A_\lambda^*)$  and  $\text{rg}(A_\lambda)$ . Write  $A_\lambda^+$  for the inverse of this operator from  $\text{rg}(A_\lambda)$  to  $\text{rg}(A_\lambda^*)$ . The orthogonal projection operator from  $\mathbb{R}^n$  onto  $\text{rg}(A_\lambda^*)$  induces a linear operator from  $\text{rg}(A_\lambda)$  into  $\text{rg}(A_\lambda^*)$ , denoted  $\overline{\Pi}_\lambda$ . Then  $S_\lambda$  is defined as the linear span of the right-singular vectors of  $A_\lambda^+ - \overline{\Pi}_\lambda$  associated to singular values smaller than 1. Again if this set is empty,  $S_\lambda = \{0\}$ . When  $A_\lambda$  is non-singular or symmetric, we recover the definition of  $S_\lambda$  given above.

For each  $\lambda \in \Lambda$ , take  $\mathbb{S}_\lambda$  such that  $\mathbb{S}_\lambda \supset \{S_\lambda\}$ . From a theoretical point of view, it is enough to take  $\mathbb{S}_\lambda = \{S_\lambda\}$  but practically it may be wise to use a larger set and by doing so, to possibly improve the approximation of  $\widehat{f}_\lambda$  by elements of  $\mathbb{S}_\lambda$ . One may for example take  $\mathbb{S}_\lambda = \{S_\lambda^1, \dots, S_\lambda^{n-2}\}$  where  $S_\lambda^k$  is the linear span of the right-singular vectors associated to the  $k$  smallest singular values of  $A_\lambda^+ - \overline{\Pi}_\lambda$ .

#### 4.2. Choices of $\mathbb{S}$ , $\Delta$ and pen

Take  $\mathbb{S} = \bigcup_{\lambda \in \Lambda} \mathbb{S}_\lambda$  and  $\Delta$  of the form

$$\Delta(S) = a(1 \vee \dim(S)) \quad \text{for all } S \in \mathbb{S}$$

where  $a \geq 1$  satisfies Assumption 2 with  $\Sigma \leq 1$ . One may take  $a = (\log |\Lambda|) \vee 1$  even though this choice is not necessarily the best. Finally, for some  $K > 1$ , take  $\text{pen}(S) = K \text{pen}_\Delta(S)$  for all  $S \in \mathbb{S}$  and select  $\widehat{f}_\lambda$  by minimizing the criterion given by (6), taking thus  $\delta = 0$  in (9).

#### 4.3. An oracle-type inequality for linear estimators

The following holds.

**Corollary 2.** *Let  $K > 1$ ,  $\kappa \in (0, 1)$  and  $\alpha > 0$ . If Assumption 1 holds and  $\Delta(S) \leq \kappa n$  for all  $S \in \mathbb{S}$ , the estimator  $\widehat{f}_\lambda$  satisfies*

$$Ca^{-1} \mathbb{E} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 \right] \leq \inf_\lambda \mathbb{E} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 \right] + \sigma^2,$$

for some  $C$  depending on  $K, \alpha$  and  $\kappa$  only.

The problem of selecting some best linear estimator among a family of those have also been considered in Arlot and Bach (2009) in the Gaussian regression

framework, and in Goldenshluger and Lepski (2009) in the multidimensional Gaussian white noise model. Arlot and Bach proposed a penalized procedure based on random penalties. Unlike ours, their approach requires that the operators be symmetric with eigenvalues in  $[0, 1]$  and that the cardinality of  $\Lambda$  is at most polynomial with respect to  $n$ . Goldenshluger and Lepski proposed a selection rule among families of kernel estimators to solve the problem of structural adaptation. Their approach requires suitable assumptions on the kernels while ours requires nothing. Nevertheless, we restrict to the case of the Euclidean loss whereas Goldenshluger and Lepski considered more general  $\mathbb{L}_p$  ones.

## 5. Variable selection

Throughout this section, we consider the problem of variable selection introduced in Example 2 and assume that  $p \geq 2$  in order to avoid trivialities. When  $p$  is small enough (say smaller than 20), this problem can be solved by using a suitable variable selection procedure that explores all the subsets of  $\{1, \dots, p\}$ . For example, one may use the penalized criterion introduced in Birgé and Massart (2001) when the variance is known, and the one in Baraud *et al* (2009) when it is not. When  $p$  is larger, such an approach can no longer be applied since it becomes numerically intractable. To overcome this problem, algorithms based on the minimization of convex criteria have been proposed among which are the Lasso, the Dantzig selector of Candès and Tao (2007), the elastic net of Zou and Hastie (2005). An alternative to those criteria is the forward-backward algorithm described in Zhang (2008), among others. Since there seems to be no evidence that one of these procedures outperforms all the others, it may be reasonable to mix them all and let the data decide which is the more appropriate to solve the problem at hand. As enlarging  $\mathbb{F}$  can only improve the risk bound of our estimator, only the CPU resources should limit the number of candidate estimators.

The procedure we propose could not only be used to select among those candidate procedures but also to select the tuning parameters they depend on. From this point of view, it provides an alternative to the cross-validation techniques which are quite popular but offer little theoretical guarantees.

### 5.1. Implementation roadmap

Start by choosing a family  $\mathcal{L}$  of variable selection procedures. Examples of such procedures are the Lasso, the Dantzig selector, the elastic net, among others. If necessary, associate to each  $\ell \in \mathcal{L}$  a family of tuning parameters  $H_\ell$ . For example, in order to use the Lasso procedure one needs to choose a tuning parameter  $h > 0$  among a grid  $H_{\text{Lasso}} \subset \mathbb{R}_+$ . If a selection procedure  $\ell$  requires no choice of tuning parameters, then one may take  $H_\ell = \{0\}$ . Let us denote by  $\widehat{m}(\ell, h)$  the subset of  $\{1, \dots, p\}$  corresponding to the predictors selected by the procedure  $\ell$  for the choice of the tuning parameter  $h$ . For  $m \subset \{1, \dots, p\}$ , let  $S_m$  be the linear span of the column vectors  $X_{\cdot, j}$  for  $j \in m$  (with the convention  $S_\emptyset = \{0\}$ ). For  $\ell \in \mathcal{L}$  and  $h \in H_\ell$ , associate to the subset  $\widehat{m}(\ell, h)$  an estimator  $\widehat{f}_{(\ell, h)}$  of  $f$  with values in  $S_{\widehat{m}(\ell, h)}$  (one may for example take the projection of  $Y$  onto the random linear space  $S_{\widehat{m}(\ell, h)}$  but any other choice would suit). Finally, consider the family  $\mathbb{F} = \{\widehat{f}_\lambda, \lambda \in \Lambda\}$  of these estimators by taking  $\Lambda = \bigcup_{\ell \in \mathcal{L}} (\{\ell\} \times H_\ell)$  and set  $\widehat{\mathcal{M}} = \{\widehat{m}(\lambda), \lambda \in \Lambda\}$ . All along we assume that  $\Lambda$  is finite (so that we take  $\delta = 0$  in (9)).

#### *The approximation spaces and the weight function*

Throughout, we shall restrict ourselves to subsets of predictors with cardinality not larger than some  $D_{\max} \leq n - 2$ . In view of approximating the estimators  $\widehat{f}_\lambda$ , we suggest the collection  $\mathbb{S}$  given by

$$\mathbb{S} = \bigcup \{S_m \mid m \subset \{1, \dots, p\}, \text{card}(m) \leq D_{\max}\}. \quad (26)$$

We associate to  $\mathbb{S}$  the weight function  $\Delta$  defined for  $S \in \mathbb{S}$  by

$$\Delta(S) = \log \left[ \binom{p}{D} \right] + \log(1 + D) \quad \text{with } D = \dim(S). \quad (27)$$

Since

$$\begin{aligned} \sum_{S \in \mathbb{S}} e^{-\Delta(S)} &= \sum_{D=0}^p \sum_{\substack{S \in \mathbb{S} \\ \dim(S) = D}} e^{-\Delta(S)} \\ &\leq \sum_{D=0}^p e^{-\log(1+D)} \leq 1 + \log(1 + p), \end{aligned}$$

Assumption 2 is satisfied with  $\Sigma = 1 + \log(1 + p)$ .

Let us now turn to the choices of the  $\mathbb{S}_\lambda \subset \mathbb{S}$ . The criterion given by (6) cannot be computed when  $\mathbb{S}_\lambda = \mathbb{S}$  for all  $\lambda$  as soon as  $p$  is too large. In such a case, one must consider a smaller subset of  $\mathbb{S}$  and we suggest for  $\lambda = (\ell, h) \in \Lambda$

$$\mathbb{S}_{(\ell, h)} = \{S_{\widehat{m}(\ell, h')}, h' \in H_\ell\}$$

(where the  $S_m$  are defined above), or preferably

$$\mathbb{S}_{(\ell, h)} = \{S_{\widehat{m}(\ell', h')}, \ell' \in \mathcal{L}, h' \in H_\ell\}$$

whenever this latter family is not too large. Note that these two families are random.

## 5.2. The results

Our choices of  $\Delta$  and  $\mathbb{S}_\lambda$  ensure that  $\widehat{f}_\lambda \in S_{\widehat{m}(\lambda)} \in \mathbb{S}_\lambda$  for all  $\lambda \in \Lambda$  and that

$$\Delta(S_{\widehat{m}(\lambda)}) \leq 2 \dim(S_{\widehat{m}(\lambda)}) \log p.$$

Hence, by applying Corollary 1 with  $\widehat{S}_\lambda = S_{\widehat{m}(\lambda)}$ , we get the following result.

**Corollary 3.** *Let  $K > 1$ ,  $\kappa \in (0, 1)$  and  $D_{\max}$  be some positive integer satisfying  $D_{\max} \leq \kappa n / (2 \log p)$ . Let  $\widehat{\mathcal{M}} = \{\widehat{m}(\lambda), \lambda \in \Lambda\}$  be a (finite) collection of random subsets of  $\{1, \dots, p\}$  with cardinality not larger than  $D_{\max}$  based on the observation  $Y$  and  $\{\widehat{f}_\lambda, \lambda \in \Lambda\}$  a family of estimators  $f$ , also based on  $Y$ , such that  $\widehat{f}_\lambda \in S_{\widehat{m}(\lambda)}$ . By applying our selection procedure, the resulting estimator  $\widehat{f}_{\widehat{\lambda}}$  satisfies*

$$C \mathbb{E} \left[ \left\| f - \widehat{f}_{\widehat{\lambda}} \right\|^2 \right] \leq \inf_{\lambda \in \Lambda} \left[ \mathbb{E} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 \right] + \mathbb{E} [\dim(S_{\widehat{m}(\lambda)})] \log(p) \sigma^2 \right],$$

where  $C$  is a constant depending on the choices of  $K$  and  $\kappa$  only.

Again, note that the risk bound we get is non-increasing with respect to  $\Lambda$ . This means that if one adds a new variable selection procedure or considers more tuning parameters to increase  $\Lambda$ , the risk bound we get can only be improved.

Without additional information on the estimators  $\widehat{f}_\lambda$  it is difficult to compare  $\mathbb{E} [\dim(S_{\widehat{m}(\lambda)})] \sigma^2$  and  $\mathbb{E} [\|f - \widehat{f}_\lambda\|^2]$ . If  $\widehat{f}_\lambda$  is of the form  $\Pi_S Y$  for some deterministic subset  $S \in \mathbb{S}$  it is well-known that

$$\mathbb{E} [\|f - \Pi_S Y\|^2] = \|f - \Pi_S f\|^2 + \dim(S)\sigma^2 \geq \dim(S)\sigma^2.$$

Under the assumption that  $f \in S_{m^*}$  and that  $m^*$  belongs to  $\widehat{\mathcal{M}}$  with probability close enough to 1, we can compare the risk of the estimator  $\widehat{f}_\lambda$  to the cardinality of  $m^*$ .

**Corollary 4.** *Assume that the assumptions of Corollary 3 hold and that  $\widehat{f}_\lambda = \Pi_{S_{\widehat{m}(\lambda)}} Y$  for all  $\lambda \in \Lambda$ . If  $f \in S_{m^*}$  for some non-void subset  $m^* \subset \{1, \dots, p\}$  with cardinality not larger than  $D_{\max}$ , then*

$$C\mathbb{E} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 \right] \leq \log(p) |m^*| \sigma^2 + R_n(m^*)$$

where  $C$  is a constant depending on  $K$  and  $\kappa$  only, and

$$R_n(m^*) = (\|f\|^2 + n\sigma^2) \left( \mathbb{P} \left[ m^* \notin \widehat{\mathcal{M}} \right] \right)^{1/2}.$$

Zhao and You (2006) gives sufficient conditions on the design  $X$  to ensure that  $\mathbb{P} \left[ m^* \notin \widehat{\mathcal{M}} \right]$  is exponentially small with respect to  $n$  when the family  $\widehat{\mathcal{M}}$  is obtained by using the LARS-Lasso algorithm with different values of the tuning parameter.

## 6. Simulation study

In the linear regression setting described in Example 2, we carry out a simulation study to evaluate the performances of our procedure to solve the two following problems.

We first consider the problem, described in Example 3, of tuning the smoothing parameter of the Lasso procedure for estimating  $f$ . The performances of our procedure are compared with those of the  $V$ -fold cross-validation method. Secondly, we consider the problem of variable selection. We solve it by using our criterion in view of selecting among a family  $\mathcal{L}$  of candidate variable selection procedures.

Our simulation study is based on a large number of examples which have been chosen in view of covering a large variety of situations. Most of these have been found in the literature in the context of Example 2 either for estimation or variable selection purposes when the number  $p$  of predictors is large.

The section is organized as follows. The simulation design is given in the following section. Then, we describe how our procedure is applied for tuning the Lasso and performing variable selection. Finally, we give the results of the simulation study.

### 6.1. Simulation design

One example is determined by the number of observations  $n$ , the number of variables  $p$ , the  $n \times p$  matrix  $X$ , the values of the parameters  $\beta$ , and the ratio signal/noise  $\rho$ . It is denoted by  $\text{ex}(n, p, X, \beta, \rho)$ , and the set of all considered examples is denoted  $\mathcal{E}$ . For each example, we carry out 400 simulations of  $Y$  as a Gaussian random vector with expectation  $f = X\beta$  and variance  $\sigma^2 I_n$ , where  $I_n$  is the  $n \times n$  identity matrix, and  $\sigma^2 = \|f\|^2/n\rho$ .

The collection  $\mathcal{E}$  is composed of several collections  $\mathcal{E}_e$  for  $e = 1, \dots, E$  where each collection  $\mathcal{E}_e$  is characterized by a vector of parameters  $\beta_e$ , and a set  $\mathcal{X}_e$  of matrices  $X$ :

$$\mathcal{E}_e = \{\text{ex}(n, p, X, \beta, \rho) : (n, p) \in \mathcal{I}, X \in \mathcal{X}_e, \beta = \beta_e, \rho \in \mathcal{R}\}$$

where  $\mathcal{R} = \{5, 10, 20\}$  and  $\mathcal{I}$  consists of pairs  $(n, p)$  such that  $p$  is smaller, equal or greater than  $n$ . The examples are described in further details in Section 8.2. They are inspired by examples found in Tibshirani (1996), Zou and Hastie (2005), Zou (2006), and Huang et al. (2008) for comparing the Lasso method to the ridge, adaptive Lasso and elastic net methods. They make up a large variety of situations. They include cases where

- the covariates are not, moderately or strongly correlated,
- the covariates with zero coefficients are weakly or highly correlated with covariates with non-zero coefficients,
- the covariates with non-zero coefficients are grouped and correlated within these groups,
- the lasso method is known to be inconsistent,
- few or many effects are present.

## 6.2. Tuning a smoothing parameter

We consider here the problem of tuning the smoothing parameter of the Lasso estimator as described in Example 3. Instead of considering the Lasso estimators for a fixed grid  $\Lambda$  of smoothing parameters  $\lambda$ , we rather focus on the sequence  $\{\widehat{f}_1, \dots, \widehat{f}_{D_{\max}}\}$  of estimators given by the  $D_{\max}$  first steps of the LARS-Lasso algorithm proposed by Efron *et al.* (2004). Hence, the tuning parameter is here the number  $h \in H = \{1, \dots, D_{\max}\}$  of steps. In our simulation study, we compare the performance of our criterion to that of the  $V$ -fold cross-validation for the problem of selecting the best estimator among the collection  $\mathbb{F} = \{\widehat{f}_1, \dots, \widehat{f}_{D_{\max}}\}$ .

### 6.2.1. The estimator of $f$ based on our procedure

We recall that our selection procedure relies on the choices of families  $\mathbb{S}$ ,  $\mathbb{S}_h$  for  $h \in H$ , a weight function  $\Delta$ , a penalty function  $\text{pen}$  and two universal constants  $K > 1$  and  $\alpha > 0$ . We choose the family  $\mathbb{S}$  defined by (26). We associate to  $\widehat{f}_h$  the family  $\mathbb{S}_h = \{S_{\widehat{m}(h')} \mid h' \in H\} \subset \mathbb{S}$  where the  $S_m$  are defined in Section 5.1 and  $\widehat{m}(h') \subset \{1, \dots, p\}$  is the set of indices corresponding to the predictors returned by the LARS-Lasso algorithm at step  $h' \in H$ . We take  $\text{pen}(S) = K\text{pen}_\Delta(S)$  with  $\Delta(S)$  defined by (27) and  $K = 1.1$ . This value of  $K$  is consistent with what is suggested in Baraud *et al.* (2009). The choice of  $\alpha$  is based on the following considerations. First, choosing  $\alpha$  around one seems reasonable since it weights similarly the term  $\|Y - \Pi_S \widehat{f}_\lambda\|^2$  which measures how well the estimator fits the data and the approximation term  $\|\widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda\|^2$  involved in our criterion (6). Second, simple calculation shows that the constant  $C^{-1} = C^{-1}(1.1, \alpha)$  involved in Theorem 1 is minimum for  $\alpha$  close to 0.6. We therefore carried out our simulations for  $\alpha$  varying from 0.2 to 1.5. The results being very similar for  $\alpha$  between 0.5 and 1.2, we choose  $\alpha = 0.5$ . We denote by  $\widehat{f}_{\text{pen}_\Delta}$  the resulting estimator of  $f$ .

### 6.2.2. The estimator of $f$ based on $V$ -fold cross-validation

For each  $h \in H$ , the prediction error is estimated using a  $V$ -fold cross-validation procedure, with  $V = n/10$ . The estimator  $\widehat{f}_{CV}$  is chosen by minimizing the estimated prediction error.



procedure	mean	std-err	quantiles				
			0%	50%	75%	99%	100%
CV	1.18	0.08	1.05	1.18	1.24	1.36	1.38
pen $_{\Delta}$	1.065	0.06	1.01	1.055	1.084	1.18	2.27

TABLE 1

Mean, standard-error and quantiles of the ratios  $R_{\text{ex}}/O_{\text{ex}}$  calculated over all  $\text{ex} \in \mathcal{E}$  such that  $O_{\text{ex}} < n\sigma^2/3$ . The number of such examples equals 654, see Section 8.2.

### 6.2.3. The results

The simulations were carried out with R ([www.r-project.org](http://www.r-project.org)) using the library `elasticnet`.

For each example  $\text{ex} \in \mathcal{E}$ , we estimate on the basis of 400 simulations the oracle risk

$$O_{\text{ex}} = \mathbb{E} \left( \min_{h \in H} \|f - \hat{f}_h\|^2 \right), \quad (28)$$

and the Euclidean risks  $R_{\text{ex}}(\hat{f}_{\text{pen}_{\Delta}})$  and  $R_{\text{ex}}(\hat{f}_{\text{CV}})$  of  $\hat{f}_{\text{pen}_{\Delta}}$  and  $\hat{f}_{\text{CV}}$  respectively.

The results presented in Table 1 show that our procedure tends to choose a better estimator than the CV in the sense that the ratios  $R_{\text{ex}}(\hat{f}_{\text{pen}_{\Delta}})/O_{\text{ex}}$  are closer to one than  $R_{\text{ex}}(\hat{f}_{\text{CV}})/O_{\text{ex}}$ .

Nevertheless, for a few examples these ratios are larger for our procedure than for the CV. These examples correspond to situations where the Lasso estimators are highly biased.

In practice, it is worth considering several estimation procedures in order to increase the chance to have good estimators of  $f$  among the family  $\mathbb{F}$ . Selecting among candidate procedures is the purpose of the following simulation experiment in the variable selection context.

### 6.3. Variable selection

In this section, we consider the problem of variable selection and use the procedure and notations introduced in Section 5. To solve this problem, we consider estimators of the form  $\hat{f}_{\hat{m}} = \Pi_{S_{\hat{m}}} Y$  where  $\hat{m}$  is a random subset of  $\{1, \dots, p\}$  depending on  $Y$ . Given a family  $\widehat{\mathcal{M}} = \{\hat{m}(\ell, h), \hat{m}(\ell, h) \in$

$\mathcal{L} \times H_\ell$  of such random sets, we consider the family  $\mathbb{F} = \{\widehat{f}_{\widehat{m}(\ell,h)} \mid (\ell, h) \in \mathcal{L} \times H_\ell\}$ . The descriptions of  $\mathcal{L}$  and  $H_\ell$  are postponed to Section 8.3. Let us merely mention that we choose  $\mathcal{L}$  which gathers variable selection procedures based on the Lasso, ridge regression, Elastic net, PLS1 regression, Adaptive Lasso, Random Forest, and on an exhaustive research among the subsets of  $\{1, \dots, p\}$  with small cardinality. For each procedure  $\ell$ , the parameter set  $H_\ell$  corresponds to different choices of tuning parameters. For each  $\lambda = (\ell, h) \in \mathcal{L} \times H_\ell$ , we take  $\mathbb{S}_\lambda = \{S_{\widehat{m}(\ell,h)}\}$  so that our selection rule over  $\mathbb{F}$  amounts to minimizing over  $\widehat{\mathcal{M}}$

$$\text{crit}(m) = \|Y - \Pi_{S_m} Y\|^2 + K \text{pen}_\Delta(S_m) \widehat{\sigma}_{S_m}^2, \quad (29)$$

where  $\text{pen}_\Delta$  is given by (3).

### 6.3.1. Results

The simulations were carried out with R ([www.r-project.org](http://www.r-project.org)) using the libraries `elasticnet`, `randomForest`, `pls` and the program `lm.ridge` in the library `MASS`. We first select the tuning parameters associated to the procedures  $\ell$  in  $\mathcal{L}$ . More precisely, for each  $\ell$  we select an estimator among the collection  $\mathbb{F}_\ell = \{\widehat{f}_{\widehat{m}(\ell,h)} \mid h \in H_\ell\}$  by minimizing Criterion (29) over  $\widehat{\mathcal{M}}_\ell = \{\widehat{m}(\ell, h) \mid h \in H_\ell\}$ . We denote by  $\widehat{m}(\ell)$  the selected set and by  $\widehat{f}_{\widehat{m}(\ell)}$  the corresponding projection estimator. For each example  $\text{ex} \in \mathcal{E}$  and each method  $\ell \in \mathcal{L}$ , we estimate the risk

$$R_{\text{ex},\ell} = \mathbb{E} \left( \|f - \widehat{f}_{\widehat{m}(\ell)}\|^2 \right)$$

of  $\widehat{f}_{\widehat{m}(\ell)}$  on the basis of 400 simulations and we do the same to calculate that of our estimator  $\widehat{f}_{\widehat{m}}$ ,

$$R_{\text{ex},\text{all}} = \mathbb{E} \left( \|f - \widehat{f}_{\widehat{m}}\|^2 \right).$$

Let us now define the minimum of these risks over all methods:

$$R_{\text{ex},\text{min}} = \min \{R_{\text{ex},\text{all}}, R_{\text{ex},\ell}, \ell \in \mathcal{L}\}.$$

We compare the ratios  $R_{\text{ex},\ell}/R_{\text{ex},\text{min}}$  for  $\ell \in \mathcal{L} \cup \{\text{all}\}$  to judge the performances of the candidate procedures on each example  $\text{ex} \in \mathcal{E}$ . The mean, standard deviations and quantiles of the sequence  $\{R_{\text{ex},\ell}/R_{\text{ex},\text{min}}, \text{ex} \in \mathcal{E}\}$  are presented in Table 2. In particular, the results show that

method	mean	std-err	quantiles			
			50%	75%	95%	100%
Lasso	2.82	9.40	1.12	1.33	6.38	127
ridge	1.76	1.90	1.42	1.82	2.87	36.9
pls	1.50	1.20	1.22	1.50	2.58	17
en	1.46	1.90	1.12	1.33	2.57	29
ALridge	1.20	0.31	1.15	1.26	1.51	5.78
ALpls	1.29	0.87	1.14	1.29	1.75	12.7
rFmse	4.13	9.50	1.38	2.04	19.2	118
rFpurity	3.99	10.00	1.42	2.06	15.1	138
exhaustive	22.9	45	6.30	24.5	92.9	430
all	1.16	0.16	1.12	1.25	1.47	1.95

TABLE 2

For each  $\ell \in \mathcal{L} \cup \{\text{all}\}$ , mean, standard-error and quantiles of the ratios  $R_{\text{ex},\ell}/R_{\text{ex},\min}$  calculated over all  $\text{ex} \in \mathcal{E}$ . The number of examples in the collection  $\mathcal{E}$  is equal to 660.

	$\mathcal{E}_1$	$\mathcal{E}_2$	$\mathcal{E}_3$	$\mathcal{E}_4$	$\mathcal{E}_5$	$\mathcal{E}_6$	$\mathcal{E}_7$	$\mathcal{E}_8$	$\mathcal{E}_9$	$\mathcal{E}_{10}$	$\mathcal{E}_{11}$
FDR	0.045	0.026	0.004	0.026	0.018	0.041	0.012	0.026	0.042	0.15	0.014
TDR	0.74	0.63	0.18	0.63	0.17	0.99	1	1	0.98	0.29	0.20

TABLE 3

False discovery rate (FDR) and true discovery rate (TDR) using our method, for each example with  $\rho = 10$  and  $n = p = 100$ .

- none of the procedures  $\ell$  in  $\mathcal{L}$  outperforms all the others simultaneously over all examples,
- our procedure, corresponding to  $\ell = \text{all}$ , achieves the smallest mean value. Besides, this value is very close to one.
- the variability of our procedure is small compared to the others
- for all examples, our procedure selects an estimator the risk of which does not exceed twice that of the oracle.

The false discovery rate (FDR) and the true discovery rate (TDR) are also parameters of interest in the context of variable selection. These quantities are given at Table 3 for each example when  $\rho = 10$  and  $n = p = 100$ . Except for one example, the FDR is small, while the TDR is varying a lot among the examples.

## 7. Proofs

### 7.1. Proof of Theorem 1

Throughout this section, we use the following notations. For all  $\lambda \in \Lambda$  and  $S \in \mathbb{S}_\lambda$ , we write

$$\text{crit}_\alpha(\widehat{f}_\lambda, S) = \left\| Y - \Pi_S \widehat{f}_\lambda \right\|^2 + \sigma^2 \text{pen}(S) + \alpha \left\| \widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda \right\|^2,$$

where

$$\text{pen}(S) = \widehat{\text{pen}}(S) \widehat{\sigma}_S^2 / \sigma^2, \quad \text{for all } S \in \mathbb{S}. \quad (30)$$

For all  $\lambda \in \Lambda$ , let  $S(\lambda) \in \mathbb{S}_\lambda$  be such that

$$\text{crit}_\alpha(\widehat{f}_\lambda, S(\lambda)) \leq \text{crit}_\alpha(\widehat{f}_\lambda) + \delta.$$

We also write  $\varepsilon = Y - f$  and  $\overline{S}$  for the linear space generated by  $S$  and  $f$ . It follows the facts that for all  $\lambda \in \Lambda$  and  $S \in \mathbb{S}_\lambda$

$$\text{crit}_\alpha(\widehat{f}_{\widehat{\lambda}}, S(\widehat{\lambda})) \leq \text{crit}_\alpha(\widehat{f}_{\widehat{\lambda}}) + \delta \leq \text{crit}_\alpha(\widehat{f}_\lambda) + 2\delta \leq \text{crit}_\alpha(\widehat{f}_\lambda, S) + 2\delta$$

and simple algebra that

$$\begin{aligned} & \left\| f - \Pi_{S(\widehat{\lambda})} \widehat{f}_{\widehat{\lambda}} \right\|^2 + \alpha \left\| \widehat{f}_{\widehat{\lambda}} - \Pi_{S(\widehat{\lambda})} \widehat{f}_{\widehat{\lambda}} \right\|^2 \\ & \leq \left\| f - \Pi_S \widehat{f}_\lambda \right\|^2 + \alpha \left\| \widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda \right\|^2 + 2\sigma^2 \text{pen}(S) + 2\delta \\ & \quad + 2\langle \varepsilon, \Pi_{S(\widehat{\lambda})} \widehat{f}_{\widehat{\lambda}} - f \rangle - \sigma^2 \text{pen}(S(\widehat{\lambda})) + 2\langle \varepsilon, f - \Pi_S \widehat{f}_\lambda \rangle - \sigma^2 \text{pen}(S). \end{aligned}$$

For  $\lambda \in \Lambda$  and  $S \in \mathbb{S}$ , let us set  $u_{\lambda,S} = \left( \Pi_S \widehat{f}_\lambda - f \right) / \left\| \Pi_S \widehat{f}_\lambda - f \right\|$  if  $\Pi_S \widehat{f}_\lambda \neq f$  and  $u_{\lambda,S} = 0$  otherwise. For all  $\lambda$  and  $S$ , we have  $u_{\lambda,S} \in \overline{S}$  and

$$\begin{aligned}
& \left\| f - \Pi_{S(\widehat{\lambda})} \widehat{f}_{\widehat{\lambda}} \right\|^2 + \alpha \left\| \widehat{f}_{\widehat{\lambda}} - \Pi_{S(\widehat{\lambda})} \widehat{f}_{\widehat{\lambda}} \right\|^2 \\
& \leq \left\| f - \Pi_S \widehat{f}_\lambda \right\|^2 + \alpha \left\| \widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda \right\|^2 + 2\sigma^2 \mathbf{pen}(S) + 2\delta \\
& \quad + 2 \left| \langle \varepsilon, u_{\widehat{\lambda}, S(\widehat{\lambda})} \rangle \right| \left\| \Pi_{S(\widehat{\lambda})} \widehat{f}_{\widehat{\lambda}} - f \right\| - \sigma^2 \mathbf{pen}(S(\widehat{\lambda})) \\
& \quad + 2 \left| \langle \varepsilon, u_{\lambda,S} \rangle \right| \left\| \Pi_S \widehat{f}_\lambda - f \right\| - \sigma^2 \mathbf{pen}(S) \\
& \leq \left\| f - \Pi_S \widehat{f}_\lambda \right\|^2 + \alpha \left\| \widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda \right\|^2 + 2\sigma^2 \mathbf{pen}(S) + 2\delta \\
& \quad + K^{-1} \left\| f - \Pi_{S(\widehat{\lambda})} \widehat{f}_{\widehat{\lambda}} \right\|^2 + K \left\| \Pi_{\overline{S(\widehat{\lambda})}} \varepsilon \right\|^2 - \sigma^2 \mathbf{pen}(S(\widehat{\lambda})) \\
& \quad + K^{-1} \left\| f - \Pi_S \widehat{f}_\lambda \right\|^2 + K \left\| \Pi_{\overline{S}} \varepsilon \right\|^2 - \sigma^2 \mathbf{pen}(S)
\end{aligned}$$

Hence, by using (5) and (30) we get

$$\begin{aligned}
& (1 - K^{-1}) \left\| f - \Pi_{S(\widehat{\lambda})} \widehat{f}_{\widehat{\lambda}} \right\|^2 + \alpha \left\| \widehat{f}_{\widehat{\lambda}} - \Pi_{S(\widehat{\lambda})} \widehat{f}_{\widehat{\lambda}} \right\|^2 \\
& \leq (1 + K^{-1}) \left\| f - \Pi_S \widehat{f}_\lambda \right\|^2 + \alpha \left\| \widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda \right\|^2 + 2\sigma^2 \mathbf{pen}(S) + \tilde{\Sigma} + 2\delta \\
& \leq 2(1 + K^{-1}) \left\| f - \widehat{f}_\lambda \right\|^2 + 2\delta \\
& \quad + (\alpha + 2(1 + K^{-1})) \left\| \widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda \right\|^2 + 2\sigma^2 \mathbf{pen}(S) + \tilde{\Sigma} \tag{31}
\end{aligned}$$

where

$$\tilde{\Sigma} = 2K \sum_{S \in \mathbb{S}} \left( \left\| \Pi_{\overline{S}} \varepsilon \right\|^2 - \frac{\mathbf{pen}_\Delta(S)}{n - \dim(S)} \left\| Y - \Pi_{\overline{S}} Y \right\|^2 \right)_+ .$$

For each  $S \in \mathbb{S}$ ,

$$\frac{\left\| Y - \Pi_S Y \right\|^2}{n - \dim(S)} \geq \frac{\left\| Y - \Pi_{\overline{S}} Y \right\|^2}{n - \dim(S)}$$

and since the variable  $\left\| Y - \Pi_{\overline{S}} Y \right\|^2$  is independent of  $\left\| \Pi_{\overline{S}} \varepsilon \right\|^2$  and is stochastically larger than  $\left\| \varepsilon - \Pi_{\overline{S}} \varepsilon \right\|^2$ , we deduce from the definition of  $\mathbf{pen}_\Delta(S)$  and (2), that on the one hand  $\mathbb{E}(\tilde{\Sigma}) \leq 2K\sigma^2\Sigma$ .

On the other hand, since  $S$  is arbitrary among  $\mathbb{S}_\lambda$  and since

$$\left(\frac{1}{\alpha} + \frac{1}{1-K^{-1}}\right)^{-1} \|f - \widehat{f}_\lambda\|^2 \leq (1-K^{-1}) \|f - \Pi_{S(\widehat{\lambda})}\widehat{f}_\lambda\|^2 + \alpha \|\widehat{f}_\lambda - \Pi_{S(\widehat{\lambda})}\widehat{f}_\lambda\|^2$$

we deduce from (31) that for all  $\lambda \in \Lambda$ ,

$$\|f - \widehat{f}_\lambda\|^2 \leq C^{-1} \left[ \|f - \widehat{f}_\lambda\|^2 + A(\widehat{f}_\lambda, \mathbb{S}_\lambda) + \tilde{\Sigma} + \delta \right] \quad (32)$$

with

$$C^{-1} = C^{-1}(K, \alpha) = \frac{(1 + \alpha - K^{-1})(\alpha + 2(1 + K^{-1}))}{\alpha(1 - K^{-1})}, \quad (33)$$

and (11) follows by taking the expectation on both sides of (32). Note that provided that

$$\inf_{\lambda \in \Lambda} \left[ \|f - \widehat{f}_\lambda\|^2 + A(\widehat{f}_\lambda, \mathbb{S}_\lambda) \right]$$

is measurable, we have actually proved the stronger inequality

$$C\mathbb{E} \left[ \|f - \widehat{f}_\lambda\|^2 \right] \leq \mathbb{E} \left[ \inf_{\lambda \in \Lambda} \left\{ \|f - \widehat{f}_\lambda\|^2 + A(\widehat{f}_\lambda, \mathbb{S}_\lambda) \right\} \right] + \sigma^2\Sigma + \delta. \quad (34)$$

Let us now turn to the second part of the Theorem, fixing some  $\lambda \in \Lambda$ . Since equality holds in (5), under Assumption 3 by (4)

$$\text{pen}(S) = K\text{pen}_\Delta(S) \leq C(\kappa, K)(\dim(S) \vee \Delta(S)), \quad \forall S \in \mathbb{S}.$$

If  $\mathbb{S}_\lambda$  is non-random, for some  $C' = C'(\kappa, K) > 0$  and all  $S \in \mathbb{S}_\lambda$ ,

$$\begin{aligned} & C'\mathbb{E} \left[ A(\widehat{f}_\lambda, \mathbb{S}_\lambda) \right] \\ & \leq \mathbb{E} \left[ \|\widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda\|^2 \right] + (\dim(S) \vee \Delta(S))\mathbb{E} [\widehat{\sigma}_S^2], \\ & = \mathbb{E} \left[ \|\widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda\|^2 \right] + \frac{\dim(S) \vee \Delta(S)}{n - \dim(S)} [\|f - \Pi_S f\|^2 + (n - \dim(S))\sigma^2]. \end{aligned}$$

Since  $\|f - \Pi_S f\|^2 \leq \|f - \Pi_S \widehat{f}_\lambda\|^2$ , we have

$$\|f - \Pi_S f\|^2 \leq \mathbb{E} \left[ \|f - \Pi_S \widehat{f}_\lambda\|^2 \right] \leq 2\mathbb{E} \left[ \|f - \widehat{f}_\lambda\|^2 \right] + 2\mathbb{E} \left[ \|\widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda\|^2 \right],$$

and under Assumption 3,  $(\dim(S) \vee \Delta(S))/(n - \dim(S)) \leq \kappa(1 - \kappa)^{-1}$ , and hence for all  $S \in \mathbb{S}_\lambda$

$$\begin{aligned} C' \mathbb{E} \left[ A(\widehat{f}_\lambda, \mathbb{S}_\lambda) \right] &\leq \left( 1 + \frac{2\kappa}{1 - \kappa} \right) \mathbb{E} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 \right] + \frac{2\kappa}{1 - \kappa} \mathbb{E} \left[ \left\| \widehat{f}_\lambda - \Pi_S \widehat{f}_\lambda \right\|^2 \right] \\ &\quad + (\dim(S) \vee \Delta(S)) \sigma^2. \end{aligned}$$

which leads to (12).

Let us turn to the proof of (13). We set  $\widehat{\sigma}_\lambda^2 = \widehat{\sigma}_{\widehat{S}_\lambda}^2$ . Since with probability one  $\widehat{f}_\lambda \in \widehat{S}_\lambda \in \mathbb{S}_\lambda$ ,

$$\mathbb{E} \left[ A(\widehat{f}_\lambda, \mathbb{S}_\lambda) \right] \leq \mathbb{E} \left[ \text{pen}(\widehat{S}_\lambda) \widehat{\sigma}_\lambda^2 \right]$$

and it suffices thus to bound the right-hand side. Since equality holds in (5) and since  $\widehat{f}_\lambda \in \widehat{S}_\lambda$

$$\begin{aligned} \text{pen}(\widehat{S}_\lambda) \widehat{\sigma}_\lambda^2 &= K \frac{\text{pen}_\Delta(\widehat{S}_\lambda)}{n - \dim(\widehat{S}_\lambda)} \left\| Y - \Pi_{\widehat{S}_\lambda} Y \right\|^2 \\ &\leq K \frac{\text{pen}_\Delta(\widehat{S}_\lambda)}{n - \dim(\widehat{S}_\lambda)} \left\| Y - \widehat{f}_\lambda \right\|^2 = K \frac{\text{pen}_\Delta(\widehat{S}_\lambda)}{n - \dim(\widehat{S}_\lambda)} \left\| f + \varepsilon - \widehat{f}_\lambda \right\|^2 \\ &\leq 2K \frac{\text{pen}_\Delta(\widehat{S}_\lambda)}{n - \dim(\widehat{S}_\lambda)} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 + \|\varepsilon\|^2 \right] \\ &\leq 2K \frac{\text{pen}_\Delta(\widehat{S}_\lambda)}{n - \dim(\widehat{S}_\lambda)} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 + (\|\varepsilon\|^2 - 2n\sigma^2)_+ + 2n\sigma^2 \right]. \end{aligned}$$

Under Assumption 3,  $1 \leq \Delta(\widehat{S}_\lambda) \vee \dim(\widehat{S}_\lambda) \leq \kappa n$  and we deduce from (4) that for some constant  $C$  depending only on  $K$  and  $\kappa$

$$C \text{pen}(\widehat{S}_\lambda) \widehat{\sigma}_\lambda^2 \leq \left\| f - \widehat{f}_\lambda \right\|^2 + \left( \dim(\widehat{S}_\lambda) \vee \Delta(\widehat{S}_\lambda) \right) \sigma^2 + (\|\varepsilon\|^2 - 2n\sigma^2)_+,$$

and the result follows from the fact that  $\mathbb{E}[(\|\varepsilon\|^2 - 2n\sigma^2)_+] \leq 3\sigma^2$  for all  $n$ .

## 7.2. Proof of Proposition 1

For all  $\lambda \in \Lambda$  and  $f \in S$ ,  $\left\| f - \widehat{f}_\lambda \right\| \geq \left\| \Pi_S \widehat{f}_\lambda - \widehat{f}_\lambda \right\|$  and hence,

$$\left\| f - \widehat{f}_\lambda \right\|^2 \geq \inf_{\lambda \in \Lambda} \left\| f - \widehat{f}_\lambda \right\|^2 \geq \frac{1}{2} \inf_{\lambda \in \Lambda} \left[ \left\| f - \widehat{f}_\lambda \right\|^2 + \left\| \Pi_S \widehat{f}_\lambda - \widehat{f}_\lambda \right\|^2 \right].$$

Besides, since the minimax rate of estimation over  $S$  is of order  $\dim(S)\sigma^2$ , for some universal constant  $C$ ,

$$C \sup_{f \in S} \mathbb{E} \left[ \left\| f - \widehat{f}_{\widehat{\lambda}} \right\|^2 \right] \geq \dim(S)\sigma^2.$$

Putting these bounds together lead to the result.

### 7.3. Proof of Proposition 2

Under (24), it is not difficult to see that  $d(n, M) = n/(2 \log(eM)) \geq 2$  so that  $\mathbb{S}$  is not empty and since for all  $S_m \in \mathbb{S}_{C_V}$

$$(\dim(S_m) \vee 1) \leq \Delta(S_m) = |m| + \log \left[ \binom{M}{|m|} \right] \leq |m|(1 + \log M) \leq \frac{n}{2},$$

Assumptions 1 to 4 are satisfied with  $\kappa = 1/2$ . Besides, the set  $\Lambda_{C_V}$  being compact,  $\lambda \mapsto \text{crit}_\alpha(f_\lambda)$  admits a minimum over  $\Lambda_{C_V}$  (we shall come back the minimization of this criterion at the end of the subsection) and hence we can take  $\delta = 0$ . By applying Theorem 1 and using (12), the resulting estimator  $\widehat{f}_{C_V} = \widehat{f}_{\widehat{\lambda}}$  satisfies for some universal constant  $C > 0$

$$C \mathbb{E} \left[ \left\| f - \widehat{f}_{C_V} \right\|^2 \right] \leq \inf_{g \in \mathbb{F}_\Lambda} \{ \|f - g\|^2 + \overline{A}(g, \mathbb{S}) \}, \quad (35)$$

where

$$\overline{A}(g, \mathbb{S}) = \inf_{S \in \mathbb{S}} [\|g - \Pi_S g\|^2 + (\dim(S) \vee \Delta(S)) \sigma^2]. \quad (36)$$

We bound  $\overline{A}(g, \mathbb{S})$  from above by using the following approximation result below the proof of which can be found in Makovoz (1996) (more precisely, we refer to the proof of his Theorem 2).

**Lemma 1.** *For all  $g$  in the convex hull  $\mathbb{F}_\Lambda$  of the  $\phi_j$  and all  $D \geq 1$ , there exists  $m \subset \{1, \dots, M\}$  such that  $|m| = (2D) \wedge M$  and*

$$\|g - \Pi_{S_m} g\|^2 \leq 4D^{-1} \sup_{j=1, \dots, M} \|\phi_j\|^2.$$



By using this lemma and the fact that  $\log\left(\frac{M}{D}\right) \leq D \log(eM/D)$  for all  $D \in \{1, \dots, M\}$ , we get

$$\bar{A}(g, \mathbb{S}) \leq \inf_{1 \leq D \leq d(n, M)/2} \left[ \frac{4nL^2}{D} + 2D(1 + \log(eM/(2D))) \right] \sigma^2.$$

Taking for  $D$  the integer part of

$$x(n, M, L) = \sqrt{\frac{nL^2}{\log(eM/\sqrt{nL^2})}}$$

which belongs to  $[1, d(n, M)/2]$  under (24), we get

$$\bar{A}(g, \mathbb{S}) \leq C' \sqrt{nL^2 \log(eM/\sqrt{nL^2})} \sigma^2 \quad (37)$$

for some universal constant  $C' > 0$  which together with (35) leads to the risk bound

$$\mathbb{E} \left[ \left\| f - \hat{f}_{Cv} \right\|^2 \right] - C \inf_{g \in \mathbb{F}_\Lambda} \|f - g\|^2 \leq C \sqrt{nL^2 \log(eM/\sqrt{nL^2})} \sigma^2.$$

Concerning the computation of  $\hat{f}_{Cv}$ , note that

$$\begin{aligned} \inf_{\lambda \in \Lambda} \text{crit}_\alpha(f_\lambda) &= \inf_{\lambda \in \Lambda} \inf_{S \in \mathbb{S}_{Cv}} \left[ \|Y - \Pi_S f_\lambda\|^2 + \alpha \|f_\lambda - \Pi_S f_\lambda\|^2 + \text{pen}(S) \hat{\sigma}_S^2 \right] \\ &= \inf_{S \in \mathbb{S}_{Cv}} \left\{ \left[ \inf_{\lambda \in \Lambda} (\|Y - \Pi_S f_\lambda\|^2 + \alpha \|f_\lambda - \Pi_S f_\lambda\|^2) \right] + \text{pen}(S) \hat{\sigma}_S^2 \right\}, \end{aligned}$$

and hence, one can solve the problem of minimizing  $\text{crit}_\alpha(f_\lambda)$  over  $\lambda \in \Lambda$  by proceeding into two steps. First, for each  $S$  in the finite set  $\mathbb{S}_{Cv}$  minimize the convex criterion

$$\text{crit}_\alpha(S, f_\lambda) = \|Y - \Pi_S f_\lambda\|^2 + \alpha \|f_\lambda - \Pi_S f_\lambda\|^2$$

over the convex (and compact set)  $\Lambda_{Cv}$ . Denote by  $\hat{f}_{Cv, S}$  the resulting minimizers. Then, minimize the quantity  $\text{crit}_\alpha(S, \hat{f}_{Cv, S}) + \text{pen}(S) \hat{\sigma}_S^2$  for  $S$  varying among  $\mathbb{S}_{Cv}$ . Denoting by  $\hat{S}$  such a minimizer, we have that  $\hat{f}_{Cv} = \hat{f}_{Cv, \hat{S}}$ .

#### 7.4. Proof of Proposition 3

By applying Theorem 1, we obtain that the selected estimator  $\widehat{f}_{\widehat{\lambda}}$  satisfies

$$C\mathbb{E} \left[ \left\| f - \widehat{f}_{\widehat{\lambda}} \right\|^2 \right] \leq \inf_{\lambda \in \{\text{L}, \text{MS}, \text{Cv}\}} \left[ \mathbb{E} \left[ \left\| f - \widehat{f}_{\lambda} \right\|^2 \right] + \mathbb{E} \left[ A(\widehat{f}_{\lambda}, \mathbb{S}_{\lambda}) \right] \right].$$

Let us now bound  $\mathbb{E} \left[ A(\widehat{f}_{\lambda}, \mathbb{S}_{\lambda}) \right]$  for each  $\lambda \in \Lambda$ .

If  $\lambda = \text{L}$ , by using (12) and the fact that  $\widehat{f}_{\text{L}} \in S_{\{1, \dots, M\}}$ , we have

$$C'\mathbb{E} \left[ A(\widehat{f}_{\text{L}}, \mathbb{S}_{\text{L}}) \right] \leq \mathbb{E} \left[ \left\| f - \widehat{f}_{\text{L}} \right\|^2 \right] + M\sigma^2.$$

If  $\lambda = \text{MS}$ , we may use (13) since with probability one  $\widehat{f}_{\text{MS}} \in \mathbb{S}_{\text{MS}}$  and since  $\dim(S) \vee \Delta(S) \leq 1 + \log(M)$  for all  $S \in \mathbb{S}_{\text{MS}}$ , we get

$$C'\mathbb{E} \left[ A(\widehat{f}_{\text{MS}}, \mathbb{S}_{\text{MS}}) \right] \leq \mathbb{E} \left[ \left\| f - \widehat{f}_{\text{MS}} \right\|^2 \right] + \log(M)\sigma^2.$$

Finally, let us turn to the case  $\lambda = \text{Cv}$  and denote by  $g$  the best approximation of  $f$  in  $\mathcal{C}$ . Since  $\widehat{f}_{\text{Cv}} \in \mathcal{C}$ , for all  $S \in \mathbb{S}_{\text{Cv}}$ ,

$$\begin{aligned} \left\| \widehat{f}_{\text{Cv}} - \Pi_S \widehat{f}_{\text{Cv}} \right\| &\leq \left\| \widehat{f}_{\text{Cv}} - \Pi_S g \right\| = \left\| \widehat{f}_{\text{Cv}} - f + f - g + g - \Pi_S g \right\| \\ &\leq 2 \left\| f - \widehat{f}_{\text{Cv}} \right\| + \left\| g - \Pi_S g \right\|, \end{aligned}$$

and hence by using (12)

$$C'\mathbb{E} \left[ A(\widehat{f}_{\text{Cv}}, \mathbb{S}_{\text{Cv}}) \right] \leq \mathbb{E} \left[ \left\| f - \widehat{f}_{\text{Cv}} \right\|^2 \right] + \overline{A}(g, \mathbb{S}_{\text{Cv}})$$

where  $\overline{A}(g, \mathbb{S}_{\text{Cv}})$  is given by (36). By arguing as in Section (3.1.3), we deduce that under (24)

$$C'\mathbb{E} \left[ A(\widehat{f}_{\text{Cv}}, \mathbb{S}_{\text{Cv}}) \right] \leq \mathbb{E} \left[ \left\| f - \widehat{f}_{\text{Cv}} \right\|^2 \right] + \sqrt{nL^2 \log(eM/\sqrt{nL^2})}\sigma^2.$$

By putting these bounds together we get the result.

### 7.5. Proof of Corollary 2

Since Assumptions 1 to 4 are fulfilled and  $\mathbb{F}$  is finite, we may apply Theorem 1 and take  $\delta = 0$ . By using (12), we have for some  $C$  depending on  $K, \alpha$  and  $\kappa$ ,

$$\begin{aligned} C\mathbb{E} \left[ \left\| f - \widehat{f}_{\widehat{\lambda}} \right\|^2 \right] \\ \leq \inf_{\lambda \in \Lambda} \left\{ \mathbb{E} \left[ \left\| f - \widehat{f}_{\lambda} \right\|^2 \right] + \mathbb{E} \left[ \left\| \widehat{f}_{\lambda} - \Pi_{S_{\lambda}} \widehat{f}_{\lambda} \right\|^2 \right] + a(1 + \dim(S_{\lambda}))\sigma^2 \right\}. \end{aligned}$$

For all  $\lambda \in \Lambda$ ,

$$\begin{aligned} \mathbb{E} \left[ \left\| f - \widehat{f}_{\lambda} \right\|^2 \right] &= \left\| f - A_{\lambda}f \right\|^2 + \mathbb{E} \left[ \left\| A_{\lambda}\varepsilon \right\|^2 \right] \\ &= \left\| f - A_{\lambda}f \right\|^2 + \text{Tr}(A_{\lambda}^*A_{\lambda})\sigma^2 \\ &\geq \max \left\{ \left\| f - A_{\lambda}f \right\|^2, \text{Tr}(A_{\lambda}^*A_{\lambda})\sigma^2 \right\} \end{aligned}$$

and

$$\begin{aligned} \mathbb{E} \left[ \left\| \widehat{f}_{\lambda} - \Pi_{S_{\lambda}} \widehat{f}_{\lambda} \right\|^2 \right] &= \left\| (I - \Pi_{S_{\lambda}})A_{\lambda}f \right\|^2 + \mathbb{E} \left[ \left\| (I - \Pi_{S_{\lambda}})A_{\lambda}\varepsilon \right\|^2 \right], \\ &\leq 2 \max \left\{ \left\| (I - \Pi_{S_{\lambda}})A_{\lambda}f \right\|^2, \mathbb{E} \left[ \left\| A_{\lambda}\varepsilon \right\|^2 \right] \right\} \\ &= 2 \max \left\{ \left\| (I - \Pi_{S_{\lambda}})A_{\lambda}f \right\|^2, \text{Tr}(A_{\lambda}^*A_{\lambda})\sigma^2 \right\} \end{aligned}$$

and hence, Corollary 2 follows from the next lemma.

**Lemma 2.** For all  $\lambda \in \Lambda$  we have

- (i)  $\left\| (I - \Pi_{S_{\lambda}})A_{\lambda}f \right\| \leq \left\| f - A_{\lambda}f \right\|$ ,
- (ii)  $\dim(S_{\lambda}) \leq 4 \text{Tr}(A_{\lambda}^*A_{\lambda})$ .

Proof of Lemma 2: Writing  $f = f_0 + f_1 \in \ker(A_{\lambda}) \oplus \text{rg}(A_{\lambda}^*)$  and using the fact that  $\text{rg}(A_{\lambda}^*) = \ker(A_{\lambda})^{\perp}$  and the definition of  $\overline{\Pi}_{\lambda}$ , we obtain

$$\begin{aligned} \left\| f - A_{\lambda}f \right\|^2 &= \left\| f_0 + f_1 - A_{\lambda}f_1 \right\|^2 \\ &= \left\| f_0 - \Pi_{\ker(A_{\lambda})}A_{\lambda}f_1 \right\|^2 + \left\| (I - \overline{\Pi}_{\lambda})A_{\lambda}f_1 \right\|^2 \\ &\geq \left\| (A_{\lambda}^+ - \overline{\Pi}_{\lambda})A_{\lambda}f_1 \right\|^2 \\ &\geq \sum_{k=1}^{m_{\lambda}} s_k^2 \langle A_{\lambda}f, v_k \rangle^2, \end{aligned}$$

where  $s_1 \geq \dots \geq s_{m_\lambda}$  are the singular values of  $A_\lambda^+ - \bar{\Pi}_\lambda$  counted with their multiplicity and  $(v_1, \dots, v_{m_\lambda})$  is an orthonormal family of right-singular vectors associated to  $(s_1, \dots, s_{m_\lambda})$ . If  $s_1 < 1$ , then  $S_\lambda = \mathbb{R}^n$  and we have  $\|f - A_\lambda f\| \geq \|(I - \Pi_{S_\lambda})A_\lambda f\| = 0$ . Otherwise,  $s_1 \geq 1$ , we may consider  $k_\lambda$  as the largest  $k$  such that  $s_k \geq 1$  and derive that

$$\begin{aligned} \|f - A_\lambda f\|^2 &\geq \sum_{k=1}^{k_\lambda} s_k^2 \langle A_\lambda f, v_k \rangle^2 \\ &\geq \sum_{k=1}^{k_\lambda} \langle A_\lambda f, v_k \rangle^2 = \|(I - \Pi_{S_\lambda})A_\lambda f\|^2, \end{aligned}$$

which proves the assertion (i).

For the bound (ii), we set  $M_\lambda = A_\lambda^+ - \bar{\Pi}_\lambda$  and note that

$$(M_\lambda - \bar{\Pi}_\lambda)(M_\lambda - \bar{\Pi}_\lambda)^* = M_\lambda M_\lambda^* + \bar{\Pi}_\lambda \bar{\Pi}_\lambda^* - M_\lambda \bar{\Pi}_\lambda^* - \bar{\Pi}_\lambda M_\lambda^*$$

induces a semi-positive quadratic form on  $\text{rg}(A_\lambda^*)$ . As a consequence the quadratic form  $(M_\lambda + \bar{\Pi}_\lambda)(M_\lambda + \bar{\Pi}_\lambda)^*$  is dominated by the quadratic form  $2(M_\lambda M_\lambda^* + \bar{\Pi}_\lambda \bar{\Pi}_\lambda^*)$  on  $\text{rg}(A_\lambda^*)$ . Furthermore

$$(M_\lambda + \bar{\Pi}_\lambda)(M_\lambda + \bar{\Pi}_\lambda)^* = (A_\lambda^+)(A_\lambda^+)^* = (A_\lambda^* A_\lambda)^+$$

where  $(A_\lambda^* A_\lambda)^+$  is the inverse of the linear operator  $L_\lambda : \text{rg}(A_\lambda^*) \rightarrow \text{rg}(A_\lambda^*)$  induced by  $A_\lambda^* A_\lambda$  restricted on  $\text{rg}(A_\lambda^*)$ . We then have that the quadratic form induced by  $(A_\lambda^* A_\lambda)^+$  is dominated by the quadratic form

$$2(A_\lambda^+ - \bar{\Pi}_\lambda)(A_\lambda^+ - \bar{\Pi}_\lambda)^* + 2\bar{\Pi}_\lambda \bar{\Pi}_\lambda^*$$

on  $\text{rg}(A_\lambda^*)$ . In particular the sequence of the eigenvalues of  $(A_\lambda^* A_\lambda)^+$  is dominated by the sequence  $(2s_k^2 + 2)_{k=1, m_\lambda}$  so

$$\begin{aligned} \text{Tr}(A_\lambda^* A_\lambda) = \text{Tr}(L_\lambda) &\geq \sum_{k=1}^{m_\lambda} \frac{1}{2(1 + s_k^2)} \\ &\geq \sum_{k=k_\lambda+1}^{m_\lambda} \frac{1}{2(1 + s_k^2)} \geq \dim(S_\lambda)/4, \end{aligned}$$

which conclude the proof of Lemma 2.

### 7.6. Proof of Corollary 4

Along the section, we write  $S_*$  for  $S_{m^*}$  and  $\widehat{S}_\lambda$  for  $S_{\widehat{m}(\lambda)}$  for short. By using (10) with  $\delta = 0$  and since  $\Sigma \leq 1 + \log(1 + p)$ , we have

$$C\mathbb{E} \left[ \|f - \widehat{f}_\lambda\|^2 \right] \leq \mathbb{E} \left[ \inf_{\lambda \in \Lambda} \|f - \Pi_{\widehat{S}_\lambda} Y\|^2 + \text{pen}(\widehat{S}_\lambda) \widehat{\sigma}_{\widehat{S}_\lambda}^2 \right] + (1 + \log(p + 1))\sigma^2,$$

for some constant  $C > 0$  depending on  $K$  only. Writing  $B$  for the event  $B = \{m^* \notin \widehat{\mathcal{M}}\}$ , we have

$$\mathbb{E} \left[ \inf_{\lambda \in \Lambda} \left\{ \|f - \Pi_{\widehat{S}_\lambda} Y\|^2 + \text{pen}(\widehat{S}_\lambda) \widehat{\sigma}_{\widehat{S}_\lambda}^2 \right\} \right] \leq A_n + R'_n$$

where

$$\begin{aligned} A_n &= \mathbb{E} \left[ \|f - \Pi_{S_*} Y\|^2 + \text{pen}(S_*) \widehat{\sigma}_{S_*}^2 \right] \\ R'_n &= \mathbb{E} \left[ \inf_{\lambda \in \Lambda} \left\{ \|f - \Pi_{\widehat{S}_\lambda} Y\|^2 + \text{pen}(\widehat{S}_\lambda) \widehat{\sigma}_{\widehat{S}_\lambda}^2 \right\} \mathbf{1}_B \right]. \end{aligned}$$

Let us bound  $A_n$  from above. Note that  $\|f - \Pi_{S_*} Y\|^2 = \|\Pi_{S_*} \varepsilon\|^2$  and  $\widehat{\sigma}_{S_*}^2 = \|(I - \Pi_{S_*})\varepsilon\|^2 / (n - \dim(S_*))$  and since  $\dim(S_*) \leq D_{\max} \leq \kappa n / (2 \log p)$ , by using (4) we get

$$A_n \leq (\dim(S_*) + \text{pen}(S_*))\sigma^2 \leq C'(1 + \log(p)) \dim(S_*)\sigma^2,$$

for some constant  $C' > 0$  depending on  $K$  and  $\kappa$  only.

Let us now turn to  $R'_n$ . For all  $\lambda \in \Lambda$ ,  $\|f - \Pi_{\widehat{S}_\lambda} Y\|^2 \leq \|f\|^2$  and

$$\widehat{\sigma}_{\widehat{S}_\lambda}^2 = \frac{\|Y - \Pi_{\widehat{S}_\lambda} Y\|^2}{n - \dim(\widehat{S}_\lambda)} \leq 2 \frac{\|f\|^2 + \|\varepsilon\|^2}{n - \dim(\widehat{S}_\lambda)}.$$

Since for all  $S \in \mathbb{S}$ ,  $\dim(S) \leq D_{\max} \leq \kappa n / (2 \log p)$ , by using (4) again, there exists some positive constant  $c$  depending on  $K$  and  $\kappa$  only such that for all  $\lambda \in \Lambda$ ,  $\text{pen}(\widehat{S}_\lambda) / (n - \dim(\widehat{S}_\lambda)) \leq c$  and hence,

$$\inf_{\lambda \in \Lambda} \left\{ \|f - \Pi_{\widehat{S}_\lambda} Y\|^2 + \text{pen}(\widehat{S}_\lambda) \widehat{\sigma}_{\widehat{S}_\lambda}^2 \right\} \mathbf{1}_B \leq (1 + 2c) (\|f\|^2 + \|\varepsilon\|^2) \mathbf{1}_B.$$

Some calculation shows that  $\mathbb{E} \left[ (\|f\|^2 + \|\varepsilon\|^2)^2 \right] \leq (\|f\|^2 + 2n\sigma^2)^2$  and hence, by Cauchy-Schwarz inequality

$$R'_n \leq (1 + 2c) (\|f\|^2 + 2n\sigma^2) \sqrt{\mathbb{P}(B)}.$$

The result follows by putting the bounds on  $A_n$  and  $R'_n$  together.

## References

- Arlot, S. (2007). *Rééchantillonnage et Sélection de modèles*. PhD thesis, University Paris XI.
- Arlot, S. (2009). Model selection by resampling penalization. *Electron. J. Stat.*, 3:557–624.
- Arlot, S. and Bach, F. (2009). Data-driven calibration of linear estimators with minimal penalties. *Advances in Neural Information Processing Systems (NIPS)*, 22:46–54.
- Arlot, S. and Celisse, A. (2010). A survey of cross-validation procedures for model selection. *Statistics Surveys*, 4:40–79.
- Baraud, Y. (2000). Model selection for regression on a fixed design. *Probab. Theory Related Fields*, 117(4):467–493.
- Baraud, Y. (2010). Estimator selection with respect to Hellinger-type risks. *Probab. Theory Relat. Fields*.
- Baraud, Y., Giraud, C., and Huet, S. (2009). Gaussian model selection with an unknown variance. *Ann. Statist.*, 37(2):630–672.
- Birgé, L. (2006). Model selection via testing: an alternative to (penalized) maximum likelihood estimators. *Ann. Inst. H. Poincaré Probab. Statist.*, 42(3):273–325.
- Birgé, L. and Massart, P. (2001). Gaussian model selection. *J. Eur. Math. Soc. (JEMS)*, 3(3):203–268.
- Boulesteix, A. and Strimmer, K. (2006). Partial least squares: a versatile tool for the analysis of high-dimensional genomic data. *Briefings in Bioinformatics*, 8(1):32–44.
- Breiman, L. (2001). Random forests. *Machine Learning*, 45:5–32.
- Bunea, F., Tsybakov, A. B., and Wegkamp, M. H. (2007). Aggregation for Gaussian regression. *Ann. Statist.*, 35(4):1674–1697.
- Candès, E. and Tao, T. (2007). The Dantzig selector: statistical estimation when  $p$  is much larger than  $n$ . *Ann. Statist.*, 35(6):2313–2351.
- Cao, Y. and Golubev, Y. (2006). On oracle inequalities related to smoothing splines. *Math. Methods Statist.*, 15(4):398–414 (2007).
- Catoni, O. (1997). Mixture approach to universal model selection. Technical report, Ecole Normale Supérieure, France.
- Catoni, O. (2004). Statistical learning theory and stochastic optimization. In *Lecture notes from the 31st Summer School on Probability Theory held in Saint-Flour, July 8–25, 2001*. Springer-Verlag, Berlin.
- Celisse, A. (2008). *Model selection via cross-validation in density estimation*,

- regression, and change-points detection*. PhD thesis, University Paris XI.
- Chen, S. S., Donoho, D. L., and Saunders, M. A. (1998). Atomic decomposition by basis pursuit. *SIAM J. Sci. Comput.*, 20(1):33–61 (electronic).
- Díaz-Uriarte, R. and Alvares de Andrés, S. (2006). Gene selection and classification of microarray data using random forest. *BMC Bioinformatics*, 7(3).
- Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. (2004). Least angle regression. *Ann. Statist.*, 32(2):407–499. With discussion, and a rejoinder by the authors.
- Genuer, R., Poggi, J.-M., and Tuleau-Malot, C. (2010). Variable selection using random forests. *Pattern Recognition Lett.*, to appear.
- Giraud, C. (2008). Mixing least-squares estimators when the variance is unknown. *Bernoulli*, 14(4):1089–1107.
- Goldenshluger, A. (2009). A universal procedure for aggregating estimators. *Ann. Statist.*, 37(1):542–568.
- Goldenshluger, A. and Lepski, O. (2009). Structural adaptation via  $\mathbb{L}_p$ -norm oracle inequalities. *Probab. Theory Related Fields*, 143(1-2):41–71.
- Helland, I. (2006). Partial least squares regression. In Kotz, S., Balakrishnan, N., Read, C., Vidakovic, B., and Johnston, N., editors, *Encyclopedia of statistical sciences (2nd ed.)*, volume 9, pages 5957–5962, New York. Wiley.
- Helland, I. (2001). Some theoretical aspects of partial least squares regression. *Chemometrics and Intelligent Laboratory Systems*, 58:97–107.
- Hoerl, A. and Kennard, R. (1970). Ridge regression: bayes estimation for nonorthogonal problems. *Technometrics*, 12:55–67.
- Hoerl, A. and Kennard, R. (2006). Ridge regression. In Kotz, S., Balakrishnan, N., Read, C., Vidakovic, B., and Johnston, N., editors, *Encyclopedia of statistical sciences (2nd ed.)*, volume 11, pages 7273–7280, New York. Wiley.
- Huang, J., Ma, S., and Zhang, C.-H. (2008). Adaptive Lasso for sparse high-dimensional regression models. *Statistica Sinica*, 4(1603-1618).
- Juditsky, A. and Nemirovski, A. (2000). Functional aggregation for nonparametric regression. *Ann. Statist.*, 28(3):681–712.
- Lepskiĭ, O. V. (1990). A problem of adaptive estimation in Gaussian white noise. *Teor. Veroyatnost. i Primenen.*, 35(3):459–470.
- Lepskiĭ, O. V. (1991). Asymptotically minimax adaptive estimation. I. Upper bounds. Optimally adaptive estimates. *Teor. Veroyatnost. i Primenen.*, 36(4):645–659.
- Lepskiĭ, O. V. (1992a). Asymptotically minimax adaptive estimation. II.

- Schemes without optimal adaptation. Adaptive estimates. *Teor. Veroyatnost. i Primenen.*, 37(3):468–481.
- Lepskii, O. V. (1992b). On problems of adaptive estimation in white Gaussian noise. In *Topics in nonparametric estimation*, volume 12 of *Adv. Soviet Math.*, pages 87–106. Amer. Math. Soc., Providence, RI.
- Leung, G. and Barron, A. R. (2006). Information theory and mixing least-squares regressions. *IEEE Trans. Inform. Theory*, 52(8):3396–3410.
- Makovoz, Y. (1996). Random approximants and neural networks. *J. Approx. Theory*, 85(1):98–109.
- Nadaraya, E. A. (1964). On estimating regression. *Theory of Probability and its Applications*, 9(1):141–142.
- Nemirovski, A. (2000). Topics in non-parametric statistics. In *Lectures on probability theory and statistics (Saint-Flour, 1998)*, volume 1738 of *Lecture Notes in Math.*, pages 85–277. Springer, Berlin.
- Rigollet, P. and Tsybakov, A. B. (2007). Linear and convex aggregation of density estimators. *Math. Methods Statist.*, 16(3):260–280.
- Strobl, C., Boulesteix, A.-L., Kneib, T., Augustin, T., and Zeileis, A. (2008). Conditional variable importance for random forests. *BMC Bioinformatics*, 9(307).
- Strobl, C., Boulesteix, A.-L., Zeileis, A., and Hothorn, T. (2007). Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC Bioinformatics*, 8(25).
- Tenenhaus, M. (1998). *La régression PLS*. Éditions Technip, Paris. Théorie et pratique. [Theory and application].
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B*, 58(1):267–288.
- Tsybakov, A. B. (2003). Optimal rates of aggregation. In *Proceedings of the 16th Annual Conference on Learning Theory (COLT) and 7th Annual Workshop on Kernel Machines*, pages 303–313. Lecture Notes in Artificial Intelligence 2777, Springer-Verlag, Berlin.
- Watson, G. S. (1964). Smooth regression analysis. *Sankhyā Ser. A*, 26:359–372.
- Wegkamp, M. (2003). Model selection in nonparametric regression. *Ann. Statist.*, 31:252–273.
- Yang, Y. (1999). Model selection for nonparametric regression. *Statist. Sinica*, 9:475–499.
- Yang, Y. (2000a). Combining different procedures for adaptive regression. *J. Multivariate Anal.*, 74(1):135–161.



- Yang, Y. (2000b). Mixing strategies for density estimation. *Ann. Statist.*, 28(1):75–87.
- Yang, Y. (2001). Adaptive regression by mixing. *J. Amer. Statist. Assoc.*, 96(454):574–588.
- Zhang, T. (2005). Learning bounds for kernel regression using effective data dimensionality. *Neural Comput.*, 17(9):2077–2098.
- Zhang, T. (2008). Adaptive forward-backward greedy algorithm for learning sparse representations. Technical report, Rutgers University, NJ.
- Zhao, P. and Yu, B. (2006). On Model Selection Consistency of Lasso. *JMLR* 7(Nov):2541–2563.
- Zou, H. (2006). The adaptive lasso and its oracle properties. *J. Amer. Statist. Assoc.*, 101(476):1418–1429.
- Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 67(2):301–320.

## 8. Appendix

### 8.1. Computation of $\text{pen}_\Delta(S)$

The penalty  $\text{pen}_\Delta(S)$ , defined at equation (3), is linked to the EDkhi function introduced in Baraud *al* (2009) (see Definition 3), via the following formula:

$$\text{pen}_\Delta(S) = \frac{n - \dim(S)}{n - \dim(S) - 1} \text{EDkhi} \left( \dim(S) + 1, n - \dim(S) - 1, \frac{e^{-\Delta(S)}}{\dim(S) + 1} \right).$$

Therefore, according to the result given in Section 6.1 in Baraud *et al* (2009),  $\text{pen}_\Delta(S)$  is the solution in  $x$  of the equation

$$\begin{aligned} \frac{e^{-\Delta(S)}}{D+1} &= \mathbb{P} \left( F_{D+3, N-1} \geq x \frac{N-1}{N(D+3)} \right) \\ &\quad - x \frac{N-1}{N(D+1)} \mathbb{P} \left( F_{D+1, N+1} \geq x \frac{N+1}{N(D+1)} \right). \end{aligned}$$

### 8.2. Simulated examples

The collection  $\mathcal{E}$  is composed of several collections  $\mathcal{E}_1, \dots, \mathcal{E}_{11}$  that are detailed below. The collections  $\mathcal{E}_1$  to  $\mathcal{E}_{10}$  are composed of examples where  $X$  is

generated as  $n$  independent centered Gaussian vectors with covariance matrix  $C$ . For each  $e \in \{1, \dots, 10\}$ , we define a  $p \times p$  matrix  $C_e$  and a  $p$ -vector of parameters  $\beta_e$ . We denote by  $\mathcal{X}_e$  the set of 5 matrices  $X$  simulated as  $n$ -i.i.d  $\mathcal{N}_p(0, C_e)$ . The collection  $\mathcal{E}_e$  is then defined as follows:

$$\mathcal{E}_e = \{\text{ex}(n, p, X, \beta, \rho), (n, p) \in \mathcal{I}, X \in \mathcal{X}_e, \beta = \beta_e, \rho \in \mathcal{R}\}$$

where  $\mathcal{R} = \{5, 10, 20\}$  and

$$\mathcal{I} = \{(100, 50), (100, 100), (100, 1000), (200, 100), (200, 200)\} \quad (38)$$

in Section 6.2, and

$$\mathcal{I} = \{(100, 50), (100, 100), (200, 100), (200, 200)\} \quad (39)$$

in Section 6.3.

Let us now describe the collections  $\mathcal{E}_1$  to  $\mathcal{E}_{10}$ .

*Collection  $\mathcal{E}_1$*  The matrix  $C$  equals the  $p \times p$  identity matrix denoted  $I_p$ . The parameters  $\beta$  satisfy  $\beta_j = 0$  for  $j \geq 16$ ,  $\beta_j = 2.5$  for  $1 \leq j \leq 5$ ,  $\beta_j = 1.5$  for  $6 \leq j \leq 10$ ,  $\beta_j = 0.5$  for  $11 \leq j \leq 15$ .

*Collection  $\mathcal{E}_2$*  the matrix  $C$  is such that  $C_{jk} = r^{|j-k|}$ , for  $1 \leq j, k \leq 15$  and  $16 \leq j, k \leq p$  with  $r = 0.5$ . Otherwise  $C_{j,k} = 0$ . The parameters  $\beta$  are as in Collection  $\mathcal{E}_1$ .

*Collection  $\mathcal{E}_3$*  The matrix  $C$  is as in Collection  $\mathcal{E}_2$  with  $r = 0.95$ , the parameters  $\beta$  are as in Collection  $\mathcal{E}_1$ .

*Collection  $\mathcal{E}_4$*  The matrix  $C$  is such that  $C_{jk} = r^{|j-k|}$ , for  $1 \leq j, k \leq p$ , with  $r = 0.5$ , the parameters  $\beta$  are as in Collection  $\mathcal{E}_1$ .

*Collection  $\mathcal{E}_5$*  the matrix  $C$  is as in Collection  $\mathcal{E}_4$  with  $r = 0.95$ , the parameters  $\beta$  are as in Collection  $\mathcal{E}_1$ .

*Collection  $\mathcal{E}_6$*  The matrix  $C$  equals  $I_p$ . The parameters  $\beta$  satisfy  $\beta_j = 0$  for  $j \geq 16$ ,  $\beta_j = 1.5$  for  $j \leq 15$ .

*Collection  $\mathcal{E}_7$*  The matrix  $C$  satisfies  $C_{j,k} = (1 - \rho_1)\mathbb{1}_{j=k} + \rho_1$  for  $1 \leq j, k \leq 3$ ,  $C_{j,k} = C_{k,j} = \rho_2$  for  $j = 4, k = 1, 2, 3$ ,  $C_{j,k} = \mathbb{1}_{j=k}$  for  $j, k \geq 5$ , with  $\rho_1 = .39$  and  $\rho_2 = .23$ . The parameters  $\beta$  satisfy  $\beta_j = 0$  for  $j \geq 4$ ,  $\beta_j = 5.6$  for  $j \leq 3$ .

*Collection  $\mathcal{E}_8$*  The matrix  $C$  satisfies  $C_{j,k} = 0.5^{|j-k|}$  for  $j, k \leq 8$ ,  $C_{j,k} = \mathbb{1}_{j=k}$  for  $j, k \geq 9$ . The parameters  $\beta$  satisfy  $\beta_j = 0$  for  $j \notin \{1, 2, 5\}$ ,  $\beta_1 = 3$ ,  $\beta_2 = 1.5$ ,  $\beta_5 = 2$ .

*Collection  $\mathcal{E}_9$*  The matrix  $C$  is defined as in Example  $\mathcal{E}_8$ . The parameters  $\beta$  satisfy  $\beta_j = 0$  for  $j \geq 9$ ,  $\beta_j = 0.85$  for  $j \leq 8$ .

*Collection  $\mathcal{E}_{10}$*  The matrix  $C$  satisfies  $C_{j,k} = 0.5\mathbb{1}_{j \neq k} + \mathbb{1}_{j=k}$  for  $j, k \leq 40$ ,  $C_{j,k} = \mathbb{1}_{j=k}$  for  $j, k \geq 41$ . The parameters  $\beta$  satisfy  $\beta_j = 2$  for  $11 \leq j \leq 20$  and  $31 \leq j \leq 40$ ,  $\beta_j = 0$  otherwise.

*Collection  $\mathcal{E}_{11}$*  In this last example, we denote by  $\mathcal{X}_{11}$  the set of 5 matrices  $X$  simulated as follows. For  $1 \leq j \leq p$ , we denote by  $X_j$  the column  $j$  of  $X$ . Let  $E$  be generated as  $n$  i.i.d.  $\mathcal{N}_p(0, 0.01I_p)$  and let  $Z_1, Z_2, Z_3$  be generated as  $n$  i.i.d.  $\mathcal{N}_3(0, I_3)$ . Then for  $j = 1, \dots, 5$ ,  $X_j = Z_1 + E_j$ , for  $j = 6, \dots, 10$ ,  $X_j = Z_2 + E_j$ , for  $j = 11, \dots, 15$ ,  $X_j = Z_3 + E_j$ , for  $j \geq 16$ ,  $X_j = E_j$ . The parameters  $\beta$  are as in Collection  $\mathcal{E}_6$ . The collection  $\mathcal{E}_{11}$  is defined as the set of examples  $\text{ex}(n, p, X, \beta, \rho)$  for  $(n, p) \in \mathcal{I}$ ,  $X \in \mathcal{X}_{11}$ , and  $\rho \in \mathcal{R}$ .

The collection  $\mathcal{E}$  is thus composed of 660 examples for  $\mathcal{I}$  chosen as in (39), and 825 for  $\mathcal{I}$  chosen as in (38). For some of the examples, the Lasso estimators were highly biased leading to high values of the ratio  $O_{\text{ex}}/n\sigma^2$ , see Equation (28). We only keep the examples for which the Lasso estimator improves the risk of the naive estimator  $Y$  by a factor at least  $1/3$ . This convention leads us to remove 171 examples over 825. These pathological examples are coming from the collections  $\mathcal{E}_1$ ,  $\mathcal{E}_6$  and  $\mathcal{E}_7$  for  $n = 100$  and  $p \geq 100$ , and from collections  $\mathcal{E}_2$  and  $\mathcal{E}_4$  when  $p = 1000$ . The examples of collection  $\mathcal{E}_7$  were chosen by Zou to illustrate that the Lasso estimators may be highly biased. All the other examples, correspond to matrices  $X$  that are nearly orthogonal.

### 8.3. Procedures for calculating sets of predictors

Let  $\widehat{\mathcal{M}} = \bigcup_{\ell \in \mathcal{L}} \widehat{\mathcal{M}}_\ell$  where we recall that for  $\ell \in \mathcal{L}$ ,  $\widehat{\mathcal{M}}_\ell = \{\widehat{m}(\ell, h) \mid h \in H_\ell\}$ .

The Lasso procedure is described in Section 6.2. The collection  $\widehat{\mathcal{M}}_{\text{Lasso}} = \{\widehat{m}(1), \dots, \widehat{m}(D_{\max})\}$  where  $\widehat{m}(h)$  is the set of indices corresponding to the predictors returned by the LARS-Lasso algorithm at step  $h \in \{1, \dots, D_{\max}\}$  (see Section 6.2).

The ridge procedure is based on the minimization of  $\|Y - X\beta\|^2 + h\|\beta\|^2$  with respect to  $\beta$ , for some positive  $h$ , see for example Hoerl and Kennard (2006). Tibshirani (1996) noted that in the case of a large number of small effects, ridge regression gives better results than the lasso for variable selection. For each  $h \in H_{\text{ridge}}$ , the regression coefficients  $\widehat{\beta}(h)$  are calculated and a collection of predictors sets is built as follows. Let  $j_1, \dots, j_p$  be such that  $|\widehat{\beta}_{j_1}(h)| > \dots > |\widehat{\beta}_{j_p}(h)|$  and set

$$M_h = \{\{j_1, \dots, j_k\}, k = 1, \dots, D_{\max}\}.$$

Then, the collection  $\widehat{\mathcal{M}}_{\text{ridge}}$  is defined as  $\widehat{\mathcal{M}}_{\text{ridge}} = \{M_h, h \in H_{\text{ridge}}\}$ .

The elastic net procedure proposed by Zou and Hastie (2005) mixes the  $\ell_1$  and  $\ell_2$  penalties of the Lasso and the ridge procedures. Let  $H_{\text{ridge}}$  be a grid of values for the tuning parameter  $h$  of the  $\ell_2$  penalty. We choose  $\widehat{\mathcal{M}}_{\text{en}} = \{M_{(\text{en}, h)} : h \in H_{\text{ridge}}\}$  where  $M_{(\text{en}, h)}$  denotes the collection of the active sets of cardinality less than  $D_{\max}$ , selected by the elastic net procedure when the  $\ell_2$ -smoothing parameter equals  $h$ . For each  $h \in H_{\text{ridge}}$  the collection  $M_{(\text{en}, h)}$  can be conveniently computed by first calculating the ridge regression coefficients and then applying the LARS-lasso algorithm, see Zou and Hastie (2005).

The partial least squares regression (PLSR1) aims to reduce the dimensionality of the regression problem by calculating a small number of components that are useful for predicting  $Y$ . Several applications of this procedure for analysing high-dimensional genomic data have been reviewed by Boulesteix and Strimmer (2006). In particular, it can be used for calculating subsets of covariates as we did for the ridge procedure. The PLSR1 procedure constructs, for a given  $h$ , uncorrelated latent components  $t_1, \dots, t_h$  that are highly correlated with the response  $Y$ , see Helland (2006). Let  $H_{\text{pls}}$  be a grid of values for the tuning parameter  $h$ . For each  $h \in H_{\text{pls}}$ , we write  $\widehat{\beta}(h)$  for the PLS regression coefficients calculated with the first  $h$  components. We then

set  $\widehat{\mathcal{M}}_{\text{PLS}} = \{M_h : h \in H_{\text{pls}}\}$ , where  $M_h$  is build from  $\widehat{\beta}(h)$  as for the ridge procedure.

The adaptive lasso procedure proposed by Zou (2006) starts with a preliminary estimator  $\widetilde{\beta}$ . Then one applies the lasso procedure replacing the parameters  $|\beta_j|, j = 1, \dots, p$  in the  $\ell_1$  penalty by the weighted parameters  $|\beta_j|/|\widetilde{\beta}_j|^\gamma, j = 1, \dots, p$  for some positive  $\gamma$ . The idea is to increase the penalty for coefficients that are close to zero, reducing thus the bias in the estimation of  $f$  and improving the variable selection accuracy. Zou showed that, if  $\widetilde{\beta}$  is a  $\sqrt{n}$ -consistent estimator of  $\beta$ , then the adaptive lasso procedure is consistent in situations where the lasso is not. A lot of work has been done around this subject, see Huang et al. (2008) for example.

We apply the procedure with  $\gamma = 1$ , and considering two different preliminary estimators:

- using the ridge estimator,  $\widetilde{\beta}(h)$  as preliminary estimator. For each  $h \in H_{\text{ridge}}$ , the adaptive lasso procedure is applied for calculating the active sets,  $M_{\text{ALridge},h}$ , of cardinality less than  $D_{\text{max}}$ . The collection  $\widehat{\mathcal{M}}_{\text{ALridge}}$  is thus defined as  $\widehat{\mathcal{M}}_{\text{ALridge}} = \{M_{\text{ALridge},h}, h \in H_{\text{ridge}}\}$ .

- using the PLSR1 estimator,  $\widetilde{\beta}(h)$ , as preliminary estimator. The procedure is the same as described just above. The collection  $M_{\text{ALpls}}$  is defined as  $M_{\text{ALpls}} = \{M_{\text{ALpls},h}, h \in H_{\text{pls}}\}$ .

The random forest algorithm was proposed by Breiman (2001) for classification and regression problems. The procedure averages several regression trees calculated on bootstrap samples. The algorithm returns measures of variable importance that may be used for variable selection, see for example Díaz-Uriarte and Alvares de Andrés (2006), Genuer et al. (2010), Strobl et al. (2007; 2008).

Let us denote by  $h$  the number of variables randomly chosen at each split when constructing the trees and

$$H_{rF} = \{p/j \mid j \in \{3, 2, 1.5, 1\}\}.$$

For each  $h \in H_{rF}$ , we consider the set of indices

$$M_h = \{\{j_1, \dots, j_k\}, k = 1, \dots, D_{\text{max}}\},$$

where  $\{j_1, \dots, j_k\}$  are the ranks of the variable importance measures. Two importance measures are proposed. The first one is based on the decrease in

the mean square error of prediction after permutation of each of the variables. It leads to the collection  $\widehat{\mathcal{M}}_{\text{rFmse}} = \{M_h, h \in H_{rF}\}$ . The second one is based on the decrease in node impurities, and leads similarly to the collection  $\widehat{\mathcal{M}}_{\text{purity}}$ .

The exhaustive procedure considers the collection of all subsets of  $\{1, \dots, p\}$  with dimension smaller than  $D_{\max}$ . We denote this collection  $\mathcal{M}_{\text{exhaustive}}$ .

*Choice of tuning parameters* We have to choose  $D_{\max}$ , the largest number of predictors considered in the collection  $\widehat{\mathcal{M}}$ . For all methods, except the exhaustive method,  $D_{\max}$  may be large, say  $D_{\max} \leq \min(n - 2, p)$ . Nevertheless, for saving computing time, we chose  $D_{\max}$  large enough such that the dimension of the estimated subset is always smaller than  $D_{\max}$ . For the exhaustive method,  $D_{\max}$  must be chosen in order to make the calculation feasible:  $D_{\max} = 4$  for  $p = 50$ ,  $D_{\max} = 3$  for  $p = 100$  and  $D_{\max} = 2$  for  $p = 200$ .

For the ridge method we choose  $H_{\text{ridge}} = \{10^{-3}, 10^{-2}, 10^{-1}, 1, 5\}$ , and for the PLSR1 method,  $H_{\text{pls}} = 1, \dots, 5$ .

*Université de Nice Sophia-Antipolis,  
Laboratoire J-A Dieudonné, UMR CNRS 6621  
Parc Valrose  
06108, Nice cedex 02  
France  
e-mail: [baraud@unice.fr](mailto:baraud@unice.fr)*

*Ecole Polytechnique,  
CMAP, UMR CNRS 7641  
route de Saclay  
91128 Palaiseau Cedex  
France  
e-mail: [christophe.giraud@polytechnique.edu](mailto:christophe.giraud@polytechnique.edu)*

*INRA MIAJ  
78352, Jouy en Josas cedex  
France  
e-mail: [sylvie.huet@jouy.inra.fr](mailto:sylvie.huet@jouy.inra.fr)*