

# PAC-Bayesian estimation and prediction in sparse additive models

Benjamin Guedj\*

*Laboratoire de Statistique Théorique et Appliquée  
Université Pierre et Marie Curie - UPMC  
Tour 25 - 2ème étage, boîte n° 158  
4, place Jussieu  
75252 Paris Cedex 05, France  
e-mail: [benjamin.guedj@upmc.fr](mailto:benjamin.guedj@upmc.fr)*

and

Pierre Alquier†

*School of Mathematical Sciences  
University College Dublin  
Room 528 - James Joyce Library  
Belfield, Dublin 4, Ireland  
e-mail: [pierre.alquier@ucd.ie](mailto:pierre.alquier@ucd.ie)*

**Abstract:** The present paper is about estimation and prediction in high-dimensional additive models under a sparsity assumption ( $p \gg n$  paradigm). A PAC-Bayesian strategy is investigated, delivering oracle inequalities in probability. The implementation is performed through recent outcomes in high-dimensional MCMC algorithms, and the performance of our method is assessed on simulated data.

**AMS 2000 subject classifications:** Primary 62G08, 62J02, 65C40.

**Keywords and phrases:** Additive models, sparsity, regression estimation, PAC-Bayesian bounds, oracle inequality, MCMC, stochastic search.

Received August 2012.

## Contents

1	Introduction . . . . .	265
2	PAC-Bayesian prediction . . . . .	267
3	MCMC implementation . . . . .	272
4	Numerical studies . . . . .	273
5	Proofs . . . . .	276
	Acknowledgements . . . . .	288
	References . . . . .	288

---

\*Corresponding author.

†Research partially supported by the French “Agence Nationale pour la Recherche” under grant ANR-09-BLAN-0128 “PARCIMONIE”.

## 1. Introduction

Substantial progress has been achieved over the last years in estimating high-dimensional regression models. A thorough introduction to this dynamic field of contemporary statistics is provided by the recent monographs Hastie, Tibshirani and Friedman (2009); Bühlmann and van de Geer (2011). In the popular framework of linear and generalized linear models, the Lasso estimator introduced by Tibshirani (1996) immediately proved successful. Its theoretical properties have been extensively studied and its popularity has never wavered since then, see for example Bunea, Tsybakov and Wegkamp (2006); van de Geer (2008); Bickel, Ritov and Tsybakov (2009); Meinshausen and Yu (2009). However, even though numerous phenomena are well captured within this linear context, restraining high-dimensional statistics to this setting is unsatisfactory. To relax the strong assumptions required in the linear framework, one idea is to investigate a more general class of models, such as nonparametric regression models of the form  $Y = f(X) + W$ , where  $Y$  denotes the response,  $X$  the predictor and  $W$  a zero-mean noise. A good compromise between complexity and effectiveness is the additive model. It has been extensively studied and formalized for thirty years now. Amongst many other references, the reader is invited to refer to Stone (1985); Hastie and Tibshirani (1986, 1990); Härdle (1990). The core of this model is that the regression function is written as a sum of univariate functions  $f = \sum_{i=1}^p f_i$ , easing its interpretation. Indeed, each covariate's effect is assessed by a unique function. This class of nonparametric models is a popular setting in statistics, despite the fact that classical estimation procedures are known to perform poorly as soon as the number of covariates  $p$  exceeds the number of observations  $n$  in that setting.

In the present paper, our goal is to investigate a PAC-Bayesian-based prediction strategy in the high-dimensional additive framework ( $p \gg n$  paradigm). In that context, estimation is essentially possible at the price of a sparsity assumption, *i.e.*, most of the  $f_i$  functions are zero. More precisely, our setting is non-asymptotic. As empirical evidences of sparse representations accumulate, high-dimensional statistics are more and more coupled with a sparsity assumption, namely that the intrinsic dimension  $p_0$  of the data is much smaller than  $p$  and  $n$ , see e.g. Giraud, Huet and Verzelen (2012). Additive modelling under a sparsity constraint has been essentially studied under the scope of the Lasso in Meier, van de Geer and Bühlmann (2009), Suzuki and Sugiyama (2012) and Koltchinskii and Yuan (2010) or of a combination of functional grouped Lasso and backfitting algorithm in Ravikumar et al. (2009). Those papers inaugurated the study of this problem and contain essential theoretical results consisting in asymptotics (see Meier, van de Geer and Bühlmann (2009); Ravikumar et al. (2009)) and non-asymptotics (see Suzuki and Sugiyama (2012); Koltchinskii and Yuan (2010)) oracle inequalities. The present article should be seen as a constructive contribution towards a deeper understanding of prediction problems in the additive framework. It should also be stressed that our work is to be seen as an attempt to relax as much as possible assumptions made on the model, such as restrictive conditions on the regressors' matrix. We consider

them too much of a non-realistic burden when it comes to prediction problems.

Our *modus operandi* will be based on PAC-Bayesian results, which is original in that context to our knowledge. The PAC-Bayesian theory originates in the two seminal papers Shawe-Taylor and Williamson (1997); McAllester (1999) and has been extensively formalized in the context of classification (see Catoni (2004, 2007)) and regression (see Audibert (2004a,b); Alquier (2006, 2008); Audibert and Catoni (2010, 2011)). However, the methods presented in these references are not explicitly designed to cover the high-dimensional setting under the sparsity assumption. Thus, the PAC-Bayesian theory has been worked out in the sparsity perspective lately, by Dalalyan and Tsybakov (2008, 2012); Alquier and Lounici (2011); Rigollet and Tsybakov (2012). The main message of these studies is that aggregation with a properly chosen prior is able to deal effectively with the sparsity issue. Interesting additional references addressing the aggregation outcomes would be Rigollet (2006); Audibert (2009). The former aggregation procedures rely on an exponential weights approach, achieving good statistical properties. Our method should be seen as an extension of these techniques, and is particularly focused on additive modelling specificities. Contrary to procedures such as the Lasso, the Dantzig selector and other penalized methods which are provably consistent under restrictive assumptions on the Gram matrix associated to the predictors, PAC-Bayesian aggregation requires only minimal assumptions on the model. Our method is supported by oracle inequalities in probability, that are valid in both asymptotic and non-asymptotic settings. We also show that our estimators achieve the optimal rate of convergence over traditional smoothing classes such as Sobolev ellipsoids. It should be stressed that our work is inspired by Alquier and Biau (2011), which addresses the celebrated single-index model with similar tools and philosophy. Let us also mention that although the use of PAC-Bayesian techniques are original in this context, parallel work has been conducted in the deterministic design case by Suzuki (2012).

A major difficulty when considering high-dimensional problems is to achieve a favorable compromise between statistical and computational performances. The recent and thorough monograph Bühlmann and van de Geer (2011) shall provide the reader with valuable insights that address this drawback. As a consequence, the explicit implementation of PAC-Bayesian techniques remains unsatisfactory as existing routines are only put to test with small values of  $p$  (typically  $p < 100$ ), contradicting with the high-dimensional framework. In the meantime, as a solution of a convex problem the Lasso proves computable for large values of  $p$  in reasonable amounts of time. We therefore focused on improving the computational aspect of our PAC-Bayesian strategy. Monte Carlo Markov Chains (MCMC) techniques proved increasingly popular in the Bayesian community, for they probably are the best way of sampling from potentially complex probability distributions. The reader willing to find a thorough introduction to such techniques is invited to refer to the comprehensive monographs Marin and Robert (2007); Meyn and Tweedie (2009). While Alquier and Biau (2011); Alquier and Lounici (2011) explore versions of the reversible jump MCMC method (RJMCMC) introduced by Green (1995), Dalalyan and Tsybakov (2008, 2012) in-

investigate a Langevin-Monte Carlo-based method, however only a deterministic design is considered. We shall try to overcome those limitations by considering adaptations of a recent procedure whose comprehensive description is to be found in Petralias (2010); Petralias and Dellaportas (2012). This procedure called Subspace Carlin and Chib algorithm originates in the seminal paper by Carlin and Chib (1995), and has a close philosophy of Hans, Dobra and West (2007), as it favors local moves for the Markov chain. We provide numerical evidence that our method is computationally efficient, on simulated data.

The paper is organized as follows. Section 2 presents our PAC-Bayesian prediction strategy in additive models. In particular, it contains the main theoretical results of this paper which consist in oracle inequalities. Section 3 is devoted to the implementation of our procedure, along with numerical experiments on simulated data, presented in Section 4. Finally, and for the sake of clarity, proofs have been postponed to Section 5.

## 2. PAC-Bayesian prediction

Let  $(\Omega, \mathcal{A}, \mathbb{P})$  be a probability space on which we denote by  $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$  a sample of  $n$  independent and identically distributed (i.i.d.) random vectors in  $(-1, 1)^p \times \mathbb{R}$ , with  $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})$ , satisfying

$$Y_i = \psi^*(\mathbf{X}_i) + \xi_i = \sum_{j=1}^p \psi_j^*(X_{ij}) + \xi_i, \quad i \in \{1, \dots, p\},$$

where  $\psi_1^*, \dots, \psi_p^*$  are  $p$  continuous functions  $(-1, 1) \rightarrow \mathbb{R}$  and  $\{\xi_i\}_{i=1}^n$  is a set of i.i.d. (conditionally to  $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ ) real-valued random variables. Let  $\mathcal{P}$  denote the distribution of the sample  $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ . Denote by  $\mathbb{E}$  the expectation computed with respect to  $\mathbb{P}$  and let  $\|\cdot\|_\infty$  be the supremum norm. We make the two following assumptions.

(A1) For any integer  $k$ ,  $\mathbb{E}[|\xi_1|^k] < \infty$ ,  $\mathbb{E}[\xi_1 | \mathbf{X}_1] = 0$  and there exist two positive constants  $L$  and  $\sigma^2$  such that for any integer  $k \geq 2$ ,

$$\mathbb{E}[|\xi_1|^k | \mathbf{X}_1] \leq \frac{k!}{2} \sigma^2 L^{k-2}.$$

(A2) There exists a constant  $C > \max(1, \sigma)$  such that  $\|\psi^*\|_\infty \leq C$ .

Note that A1 implies that  $\mathbb{E} \xi_1 = 0$  and that the distribution of  $\xi_1$  may depend on  $\mathbf{X}_1$ . In particular, A1 holds if  $\xi_1$  is a zero-mean gaussian with variance  $\gamma^2(\mathbf{X}_1)$  where  $x \mapsto \gamma^2(x)$  is bounded.

Further, note that the boundedness assumption A2 plays a key role in our approach, as it allows to use a version of Bernstein's inequality which is one of the two main technical tools we use to state our results. This assumption is not only a technical prerequisite since it proved crucial for critical regimes: indeed, if the intrinsic dimension  $p_0$  of the regression function  $\psi^*$  is still large, the boundedness of the function class allows much faster estimation rates. This point is profusely discussed in Raskutti, Wainwright and Yu (2012).

We are mostly interested in sparse additive models, in which only a few  $\{\psi_j^*\}_{j=1}^p$  are not identically zero. Let  $\{\varphi_k\}_{k=1}^\infty$  be a known countable set of continuous functions  $\mathbb{R} \rightarrow (-1, 1)$  called the dictionary. In the sequel,  $|\mathcal{H}|$  stands for the cardinality of a set  $\mathcal{H}$ . For any  $p$ -th tuple  $\mathbf{m} = (m_1, \dots, m_p) \in \mathbb{N}^p$ , denote by  $S(\mathbf{m}) \subset \{1, \dots, p\}$  the set of indices of nonzero elements of  $\mathbf{m}$ , *i.e.*,

$$|S(\mathbf{m})| = \sum_{j=1}^p \mathbb{1}[m_j > 0],$$

and define

$$\Theta_{\mathbf{m}} = \{\theta \in \mathbb{R}^{m_1} \times \dots \times \mathbb{R}^{m_p}\},$$

with the convention  $\mathbb{R}^0 = \emptyset$ . The set  $\Theta_{\mathbf{m}}$  is embedded with its canonical Borel field  $\mathcal{B}(\Theta_{\mathbf{m}}) = \mathcal{B}(\mathbb{R}^{m_1}) \otimes \dots \otimes \mathcal{B}(\mathbb{R}^{m_p})$ . Denote by

$$\Theta \stackrel{\text{def}}{=} \bigcup_{\mathbf{m} \in \mathcal{M}} \Theta_{\mathbf{m}},$$

which is equipped with the  $\sigma$ -algebra  $\mathcal{T} = \sigma(\bigvee_{\mathbf{m} \in \mathcal{M}} \mathcal{B}(\Theta_{\mathbf{m}}))$ , where  $\mathcal{M}$  is the collection of models  $\mathcal{M} = \{\mathbf{m} = (m_1, \dots, m_p) \in \mathbb{N}^p\}$ . Consider the span of the set  $\{\varphi_k\}_{k=1}^\infty$ , *i.e.*, the set of functions

$$\mathbb{F} = \left\{ \psi_\theta = \sum_{j \in S(\mathbf{m})} \psi_j = \sum_{j \in S(\mathbf{m})} \sum_{k=1}^{m_j} \theta_{jk} \varphi_k : \theta \in \Theta_{\mathbf{m}}, \mathbf{m} \in \mathcal{M} \right\},$$

equipped with a countable generated  $\sigma$ -algebra denoted by  $\mathcal{F}$ . The risk and empirical risk associated to any  $\psi_\theta \in \mathbb{F}$  are defined respectively as

$$R(\psi_\theta) = \mathbb{E}[Y_1 - \psi_\theta(\mathbf{X}_1)]^2 \quad \text{and} \quad R_n(\psi_\theta) = r_n(\{\mathbf{X}_i, Y_i\}_{i=1}^n, \psi_\theta),$$

where

$$r_n(\{\mathbf{x}_i, y_i\}_{i=1}^n, \psi_\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \psi_\theta(\mathbf{x}_i))^2.$$

Consider the probability  $\eta_\alpha$  on the set  $\mathcal{M}$  defined by

$$\eta_\alpha : \mathbf{m} \mapsto \frac{1 - \frac{\alpha}{1-\alpha}}{1 - \left(\frac{\alpha}{1-\alpha}\right)^{p+1}} \binom{p}{|S(\mathbf{m})|}^{-1} \alpha^{\sum_{j=1}^p m_j},$$

for some  $\alpha \in (0, 1/2)$ . Let us stress the fact that the probability  $\eta_\alpha$  acts as a penalization term over a model  $\mathbf{m}$ , on the number of its active regressors through the combinatorial term  $\binom{p}{|S(\mathbf{m})|}^{-1}$  and on their expansion through  $\alpha^{\sum_{j=1}^p m_j}$ .

Our procedure relies on the following construction of the probability  $\pi$ , referred to as the prior, in order to promote the sparsity properties of the target regression function  $\psi^*$ . For any  $\mathbf{m} \in \mathcal{M}$ ,  $\zeta > 0$  and  $\mathbf{x} \in \Theta_{\mathbf{m}}$ , denote by  $\mathcal{B}_{\mathbf{m}}^1(\mathbf{x}, \zeta)$

the  $\ell^1$ -ball centered in  $\mathbf{x}$  with radius  $\zeta$ . For any  $\mathbf{m} \in \mathcal{M}$ , denote by  $\pi_{\mathbf{m}}$  the uniform distribution on  $\mathcal{B}_{\mathbf{m}}^1(0, C)$ . Define the probability  $\pi$  on  $(\Theta, \mathcal{T})$ ,

$$\pi(A) = \sum_{\mathbf{m} \in \mathcal{M}} \eta_{\alpha}(\mathbf{m}) \pi_{\mathbf{m}}(A), \quad A \in \mathcal{T}.$$

Note that the volume  $V_{\mathbf{m}}(C)$  of  $\mathcal{B}_{\mathbf{m}}^1(0, C)$  is given by

$$V_{\mathbf{m}}(C) = \frac{(2C)^{\sum_{j \in S(\mathbf{m})} m_j}}{\Gamma\left(\sum_{j \in S(\mathbf{m})} m_j + 1\right)} = \frac{(2C)^{\sum_{j \in S(\mathbf{m})} m_j}}{\left(\sum_{j \in S(\mathbf{m})} m_j\right)!}.$$

Finally, set  $\delta > 0$  (which may be interpreted as an inverse temperature parameter) and the posterior Gibbs transition density is

$$\begin{aligned} & \rho_{\delta}(\{\mathbf{x}_i, y_i\}_{i=1}^n, \theta) \\ &= \sum_{\mathbf{m} \in \mathcal{M}} \frac{\eta_{\alpha}(\mathbf{m})}{V_{\mathbf{m}}(C)} \mathbb{1}_{\mathcal{B}_{\mathbf{m}}^1(0, C)}(\theta) \frac{\exp[-\delta r_n(\{\mathbf{x}_i, y_i\}_{i=1}^n, \psi_{\theta})]}{\int \exp[-\delta r_n(\{\mathbf{x}_i, y_i\}_{i=1}^n, \psi_{\theta})] \pi(d\theta)}. \end{aligned} \quad (2.1)$$

We then consider two competing estimators. The first one is the randomized Gibbs estimator  $\hat{\Psi}$ , constructed with parameters  $\hat{\theta}$  sampled from the posterior Gibbs density, *i.e.*, for any  $A \in \mathcal{F}$ ,

$$\mathbb{P}(\hat{\Psi} \in A | \{\mathbf{X}_i, Y_i\}_{i=1}^n) = \int_A \rho_{\delta}(\{\mathbf{X}_i, Y_i\}_{i=1}^n, \theta) \pi(d\theta), \quad (2.2)$$

while the second one is the aggregated Gibbs estimator  $\hat{\Psi}^A$  defined as the posterior mean

$$\hat{\Psi}^A = \int \psi_{\theta} \rho_{\delta}(\{\mathbf{X}_i, Y_i\}_{i=1}^n, \theta) \pi(d\theta) = \mathbb{E}[\hat{\Psi} | \{\mathbf{X}_i, Y_i\}_{i=1}^n]. \quad (2.3)$$

These estimators have been introduced in Catoni (2004, 2007) and investigated in further work by Audibert (2004a); Alquier (2006, 2008); Dalalyan and Tsybakov (2008, 2012).

For the sake of clarity, denote by  $\mathcal{D}$  a generic numerical constant in the sequel. We are now in a position to write a PAC-Bayesian oracle inequality.

**Theorem 2.1.** *Let  $\hat{\psi}$  and  $\hat{\psi}^A$  be realizations of the Gibbs estimators defined by (2.2)–(2.3), respectively. Let A1 and A2 hold. Set  $w = 8C \max(L, C)$  and  $\delta = n\ell/[w + 4(\sigma^2 + C^2)]$ , for  $\ell \in (0, 1)$ , and let  $\varepsilon \in (0, 1)$ . Then with  $\mathbb{P}$ -probability at least  $1 - 2\varepsilon$ ,*

$$\begin{aligned} & \left. \begin{aligned} R(\hat{\psi}) - R(\psi^*) \\ R(\hat{\psi}^A) - R(\psi^*) \end{aligned} \right\} \leq \mathcal{D} \inf_{\mathbf{m} \in \mathcal{M}} \inf_{\theta \in \mathcal{B}_{\mathbf{m}}^1(0, C)} \left\{ R(\psi_{\theta}) - R(\psi^*) \right. \\ & \quad \left. + |S(\mathbf{m})| \frac{\log(p/|S(\mathbf{m})|)}{n} + \frac{\log(n)}{n} \sum_{j \in S(\mathbf{m})} m_j + \frac{\log(1/\varepsilon)}{n} \right\}, \end{aligned} \quad (2.4)$$

where  $\mathcal{D}$  depends upon  $w$ ,  $\sigma$ ,  $C$ ,  $\ell$  and  $\alpha$  defined above.

Under mild assumptions, [Theorem 2.1](#) provides inequalities which admit the following interpretation. If there exists a “small” model in the collection  $\mathcal{M}$ , *i.e.*, a model  $\mathbf{m}$  such that  $\sum_{j \in S(\mathbf{m})} m_j$  and  $|S(\mathbf{m})|$  are small, such that  $\psi_\theta$  (with  $\theta \in \Theta_{\mathbf{m}}$ ) is close to  $\psi^*$ , then  $\hat{\psi}$  and  $\hat{\psi}^A$  are also close to  $\psi^*$  up to  $\log(n)/n$  and  $\log(p)/n$  terms. However, if no such model exists, at least one of the terms  $\sum_{j \in S(\mathbf{m})} m_j/n$  and  $|S(\mathbf{m})|/n$  starts to emerge, thereby deteriorating the global quality of the bound. A satisfying estimation of  $\psi^*$  is typically possible when  $\psi^*$  admits a sparse representation.

To go further, we derive from [Theorem 2.1](#) an inequality on Sobolev ellipsoids. We show that our procedure achieves the optimal rate of convergence in this setting. For the sake of shortness, we consider Sobolev spaces, however one can easily derive the following results in other functional spaces such as Besov spaces. See [Tsybakov \(2009\)](#) and the references therein.

The notation  $\{\varphi_k\}_{k=1}^\infty$  now refers to the (non-normalized) trigonometric system, defined as

$$\varphi_1: t \mapsto 1, \quad \varphi_{2j}: t \mapsto \cos(\pi jt), \quad \varphi_{2j+1}: t \mapsto \sin(\pi jt),$$

with  $j \in \mathbb{N}^*$  and  $t \in (-1, 1)$ . Let us denote by  $S^*$  the set of indices of non-identically zero regressors. That is, the regression function  $\psi^*$  is

$$\psi^* = \sum_{j \in S^*} \psi_j^*.$$

Assume that for any  $j \in S^*$ ,  $\psi_j^*$  belongs to the Sobolev ellipsoid  $\mathcal{W}(r_j, d_j)$  defined as

$$\mathcal{W}(r_j, d_j) = \left\{ f \in L^2([-1, 1]): f = \sum_{k=1}^\infty \theta_k \varphi_k \quad \text{and} \quad \sum_{i=1}^\infty i^{2r_j} \theta_i^2 \leq d_j \right\}.$$

with  $d_j$  chosen such that  $\sum_{j \in S^*} \sqrt{d_j} \leq C\sqrt{6}/\pi$  and for unknown regularity parameters  $r_1, \dots, r_{|S^*|} \geq 1$ . Let us stress the fact that this assumption casts our results onto the adaptive setting. It also implies that  $\psi^*$  belongs to the Sobolev ellipsoid  $\mathcal{W}(r, d)$ , with  $r = \min_{j \in S^*} r_j$  and  $d = \sum_{j \in S^*} d_j$ , *i.e.*,

$$\psi^* = \sum_{j \in S^*} \sum_{k=1}^\infty \theta_{jk}^* \varphi_k. \tag{2.5}$$

It is worth pointing out that in that setting, the Sobolev ellipsoid is better approximated by the  $\ell^1$ -ball  $\mathcal{B}_{\mathbf{m}}^1(0, C)$  as the dimension of  $\mathbf{m}$  grows. Further, make the following assumption.

- (A3)** The distribution of the data  $\mathcal{P}$  has a probability density with respect to the corresponding Lebesgue measure, bounded from above by a constant  $B > 0$ .

**Theorem 2.2.** *Let  $\hat{\psi}$  and  $\hat{\psi}^\Lambda$  be realizations of the Gibbs estimators defined by (2.2)–(2.3), respectively. Let A1, A2 and A3 hold. Set  $w = 8C \max(L, C)$  and  $\delta = n\ell/[w + 4(\sigma^2 + C^2)]$ , for  $\ell \in (0, 1)$ , and let  $\varepsilon \in (0, 1)$ . Then with  $\mathbb{P}$ -probability at least  $1 - 2\varepsilon$ ,*

$$\left. \begin{array}{l} R(\hat{\psi}) - R(\psi^*) \\ R(\hat{\psi}^\Lambda) - R(\psi^*) \end{array} \right\} \leq \mathcal{D} \left\{ \sum_{j \in S^*} d_j^{\frac{1}{2r_j+1}} \left( \frac{\log(n)}{2nr_j} \right)^{\frac{2r_j}{2r_j+1}} + |S^*| \log(p/|S^*|)/n + \frac{\log(1/\varepsilon)}{n} \right\},$$

where  $\mathcal{D}$  is a constant depending only on  $w, \sigma, C, \ell, \alpha$  and  $B$ .

Theorem 2.2 illustrates that we obtain the minimax rate of convergence over Sobolev classes up to a  $\log(n)$  term. Indeed, the minimax rate to estimate a single function with regularity  $r$  is  $n^{-\frac{2r}{2r+1}}$ , see for example Tsybakov (2009, Chapter 2). Theorem 2.1 and Theorem 2.2 thus validate our method.

A salient fact about Theorem 2.2 is its links with existing work: assume that all the  $\psi_j^*$  belong to the same Sobolev ellipsoid  $\mathcal{W}(r, d)$ . The convergence rate is now  $\log(n)n^{-\frac{2r}{2r+1}} + \log(p)/n$ . This rate (down to a  $\log(n)$  term) is the same as the one exhibited by Koltchinskii and Yuan (2010) in the context of multiple kernel learning ( $n^{-\frac{2r}{2r+1}} + \log(p)/n$ ). Suzuki and Sugiyama (2012) even obtain faster rates which correspond to smaller functional spaces. However, the results presented by both Koltchinskii and Yuan (2010) and Suzuki and Sugiyama (2012) are obtained under stringent conditions on the design, which are not necessary to prove Theorem 2.2.

A natural extension is to consider sparsity on both regressors and their expansion, instead of sparse regressors and nested expansion as before. That is, we no longer consider the first  $m_j$  dictionary functions for the expansion of regressor  $j$ . To this aim, we slightly extend the previous notation. Let  $K \in \mathbb{N}^*$  be the length of the dictionary. A model is now denoted by  $\mathbf{m} = (\mathbf{m}_1, \dots, \mathbf{m}_p)$  and for any  $j \in \{1, \dots, p\}$ ,  $\mathbf{m}_j = (m_{j1}, \dots, m_{jK})$  is a  $K$ -sized vector whose entries are 1 whenever the corresponding dictionary function is present in the model and 0 otherwise. Introduce the notation

$$S(\mathbf{m}) = \{\mathbf{m}_j \neq \mathbf{0}, j \in \{1, \dots, p\}\}, \quad S(\mathbf{m}_j) = \{m_{jk} \neq 0, k \in \{1, \dots, K\}\}.$$

The prior distribution on the models space  $\mathcal{M}$  is now

$$\eta_\alpha: \mathbf{m} \mapsto \frac{1 - \alpha \frac{1-\alpha^{K+1}}{1-\alpha}}{1 - \left(\alpha \frac{1-\alpha^{K+1}}{1-\alpha}\right)^{p+1}} \binom{p}{|S(\mathbf{m})|}^{-1} \prod_{j \in S(\mathbf{m})} \binom{K}{|S(\mathbf{m}_j)|}^{-1} \alpha^{|S(\mathbf{m}_j)|},$$

for any  $\alpha \in (0, 1/2)$ .

**Theorem 2.3.** *Let  $\hat{\psi}$  and  $\hat{\psi}^\Lambda$  be realizations of the Gibbs estimators defined by (2.2)–(2.3), respectively. Let A1 and A2 hold. Set  $w = 8C \max(L, C)$  and  $\delta = n\ell/[w + 4(\sigma^2 + C^2)]$ , for  $\ell \in (0, 1)$ , and let  $\varepsilon \in (0, 1)$ . Then with  $\mathbb{P}$ -probability*

at least  $1 - 2\varepsilon$ ,

$$\left. \begin{aligned} R(\hat{\psi}) - R(\psi^*) \\ R(\hat{\psi}^A) - R(\psi^*) \end{aligned} \right\} \leq \mathcal{D} \inf_{\mathbf{m} \in \mathcal{M}} \inf_{\theta \in \mathcal{B}_{\mathbf{m}}^1(0, C)} \left\{ R(\psi_\theta) - R(\psi^*) \right. \\ \left. + |S(\mathbf{m})| \frac{\log(p/|S(\mathbf{m})|)}{n} + \frac{\log(nK)}{n} \sum_{j \in S(\mathbf{m})} |S(\mathbf{m}_j)| + \frac{\log(1/\varepsilon)}{n} \right\},$$

where  $\mathcal{D}$  depends upon  $w$ ,  $\sigma$ ,  $C$ ,  $\ell$  and  $\alpha$  defined above.

### 3. MCMC implementation

In this section, we describe an implementation of the method outlined in the previous section. Our goal is to sample from the Gibbs posterior distribution  $\rho_\delta$ . We use a version of the so-called Subspace Carlin and Chib (SCC) developed by Petralias (2010); Petralias and Dellaportas (2012) which originates in the Shotgun Stochastic Search algorithm (see Hans, Dobra and West (2007)). The key idea of the algorithm lies in a stochastic search heuristic that restricts moves in neighborhoods of the visited models. Let  $T \in \mathbb{N}^*$  and denote by  $\{\theta(t), \mathbf{m}(t)\}_{t=0}^T$  the Markov chain of interest, with  $\theta(t) \in \Theta_{\mathbf{m}(t)}$ . Define  $i: t \mapsto \{+, -, =\}$ , the three possible moves performed by the algorithm: an addition, a deletion or an adjustment of a regressor. Let  $\{\mathbf{e}_1, \dots, \mathbf{e}_p\}$  be the canonical base of  $\mathbb{R}^p$ . For any model  $\mathbf{m}(t) = (m_1(t), \dots, m_p(t)) \in \mathcal{M}$ , define its neighborhood  $\{\mathcal{V}^+[\mathbf{m}(t)], \mathcal{V}^-[\mathbf{m}(t)], \mathcal{V}^=[\mathbf{m}(t)]\}$ , where

$$\begin{aligned} \mathcal{V}^+[\mathbf{m}(t)] &= \{\mathbf{k} \in \mathcal{M}: \mathbf{k} = \mathbf{m}(t) + x\mathbf{e}_j, x \in \mathbb{N}^*, j \in \{1, \dots, p\} \setminus S[\mathbf{m}(t)]\}, \\ \mathcal{V}^-[\mathbf{m}(t)] &= \{\mathbf{k} \in \mathcal{M}: \mathbf{k} = \mathbf{m}(t) - m_j(t)\mathbf{e}_j, j \in S[\mathbf{m}(t)]\}, \end{aligned}$$

and

$$\mathcal{V}^=[\mathbf{m}(t)] = \{\mathbf{k} \in \mathcal{M}: S(\mathbf{k}) = S[\mathbf{m}(t)]\}.$$

A move  $i(t)$  is chosen with probability  $q[i(t)]$ . By convention, if  $S[\mathbf{m}(t)] = p$  (respectively  $S[\mathbf{m}(t)] = 1$ ) the probability of performing an addition move (respectively a deletion move) is zero. Note  $\xi: \{+, -\} \mapsto \{-, +\}$  and let  $D_{\mathbf{m}}$  be the design matrix in model  $\mathbf{m} \in \mathcal{M}$ . Denote by  $\text{LSE}_{\mathbf{m}}$  the least square estimate  $\text{LSE}_{\mathbf{m}} = (D'_{\mathbf{m}}D_{\mathbf{m}})^{-1}D'_{\mathbf{m}}\mathbf{Y}$  (with  $\mathbf{Y} = (Y_1, \dots, Y_n)$ ) in model  $\mathbf{m} \in \mathcal{M}$ . For ease of notation, let  $\mathcal{J}$  denote the identity matrix. Finally, denote by  $\phi(\cdot; \mu, \Gamma)$  the density of a Gaussian distribution  $\mathcal{N}(\mu, \Gamma)$  with mean  $\mu$  and covariance matrix  $\Gamma$ . A description of the full algorithm is presented in Algorithm 1.

The estimates  $\hat{\Psi}$  and  $\hat{\Psi}^A$  are obtained as

$$\hat{\Psi} = \sum_{j=1}^p \sum_{k=1}^K \theta_{jk}(T) \varphi_k,$$

and for some burnin  $b \in \{1, \dots, T-1\}$ ,

$$\hat{\Psi}^A = \sum_{j=1}^p \sum_{k=1}^K \left( \frac{1}{T-b} \sum_{\ell=b+1}^T \theta_{jk}(\ell) \right) \varphi_k.$$

**Algorithm 1** A Subspace Carlin and Chib-based algorithm

- 
- 1: Initialize  $(\theta(0), \mathbf{m}(0))$ .
  - 2: **for**  $t = 1$  to  $T$  **do**
  - 3:   Choose a move  $i(t)$  with probability  $q[i(t)]$ .
  - 4:   For any  $\mathbf{k} \in \mathcal{V}^{i(t)}[\mathbf{m}(t-1)]$ , generate  $\theta_{\mathbf{k}}$  from the proposal density  $\phi(\cdot; \text{LSE}_{\mathbf{k}}, \sigma^2 \mathcal{J})$ .
  - 5:   Propose a model  $\mathbf{k} \in \mathcal{V}^{i(t)}[\mathbf{m}(t-1)]$  with probability

$$\gamma(\mathbf{m}(t-1), \mathbf{k}) = \frac{A_{\mathbf{k}}}{\sum_{\mathbf{j} \in \mathcal{V}^{i(t)}[\mathbf{m}(t-1)]} A_{\mathbf{j}}},$$

where

$$A_{\mathbf{j}} = \frac{\rho_{\delta}(\theta_{\mathbf{j}})}{\phi(\theta_{\mathbf{j}}; \text{LSE}_{\mathbf{j}}, \sigma^2 \mathcal{J})}.$$

- 6:   **if**  $i(t) \in \{+, -\}$  **then**
- 7:     For any  $\mathbf{h} \in \mathcal{V}^{\xi(i(t))}[\mathbf{k}]$ , generate  $\theta_{\mathbf{h}}$  from the proposal density  $\phi(\cdot; \text{LSE}_{\mathbf{h}}, \sigma^2 \mathcal{J})$ . Note that  $\mathbf{m}(t-1) \in \mathcal{V}^{\xi(i(t))}[\mathbf{k}]$ .
- 8:     Accept model  $\mathbf{k}$ , *i.e.*, set  $\mathbf{m}(t) = \mathbf{k}$  and  $\theta(t) = \theta_{\mathbf{k}}$ , with probability

$$\begin{aligned} \alpha &= \min \left( 1, \frac{A_{\mathbf{k}} q[i(t)] \gamma(\mathbf{k}, \mathbf{m}(t-1))}{A_{\mathbf{m}(t-1)} q[\xi(i(t))] \gamma(\mathbf{m}(t-1), \mathbf{k})} \right) \\ &= \min \left( 1, \frac{q[i(t)] \sum_{\mathbf{h} \in \mathcal{V}^{i(t)}[\mathbf{m}(t-1)]} A_{\mathbf{h}}}{q[\xi(i(t))] \sum_{\mathbf{h} \in \mathcal{V}^{\xi(i(t))}[\mathbf{k}]} A_{\mathbf{h}}} \right). \end{aligned}$$

Otherwise, set  $\mathbf{m}(t) = \mathbf{m}(t-1)$  and  $\theta(t) = \theta_{\mathbf{m}(t-1)}$ .

- 9:   **else**
- 10:    Generate  $\theta_{\mathbf{m}(t-1)}$  from the proposal density  $\phi(\cdot; \text{LSE}_{\mathbf{m}(t-1)}, \sigma^2 \mathcal{J})$ .
- 11:    Accept model  $\mathbf{k}$ , *i.e.*, set  $\mathbf{m}(t) = \mathbf{k}$  and  $\theta(t) = \theta_{\mathbf{k}}$ , with probability

$$\alpha = \min \left( 1, \frac{A_{\mathbf{k}} \gamma(\mathbf{k}, \mathbf{m}(t-1))}{A_{\mathbf{m}(t-1)} \gamma(\mathbf{m}(t-1), \mathbf{k})} \right).$$

Otherwise, set  $\mathbf{m}(t) = \mathbf{m}(t-1)$  and  $\theta(t) = \theta_{\mathbf{m}(t-1)}$ .

- 12:   **end if**
  - 13: **end for**
- 

The transition kernel of the chain defined above is reversible with respect to  $\rho_{\delta} \otimes \eta_{\alpha}$ , hence this procedure ensures that  $\{\theta(t)\}_{t=1}^T$  is a Markov Chain with stationary distribution  $\rho_{\delta}$ .

#### 4. Numerical studies

In this section we validate the effectiveness of our method on simulated data. All our numerical studies have been performed with the software R (see R Core Team (2012)). The method is available on the CRAN website (<http://www.cran.r-project.org/web/packages/pacbpred/index.html>), under the name `pacbpred` (see Guedj (2012)).

Some comments are in order here about how to calibrate the constants  $C$ ,  $\sigma^2$ ,  $\delta$  and  $\alpha$ . Clearly, a too small value for  $C$  will stuck the algorithm, preventing the chain to escape from the initial model. Indeed, most proposed models will be discarded since the acceptance ratio will frequently take the

TABLE 1  
Each number is the mean (standard deviation) of the RSS over 10 independent runs

MCMC	$p = 50$ 3000 it.	$p = 200$ 10000 it.	$p = 400$ 20000 it.
Model 1	0.0318 (0.0047)	0.0320 (0.0029)	0.0335 (0.0056)
Model 2	0.0411 (0.0061)	0.1746 (0.0639)	0.2201 (0.0992)
Model 3	0.0665 (0.0421)	0.1151 (0.0399)	0.1597 (0.0579)

value 0. Conversely, a large value for  $C$  deteriorates the quality of the bound in [Theorem 2.1](#), [Theorem 2.2](#), [Theorem 2.3](#) and [Theorem 5.1](#). However, this only influences the theoretical bound, as its contribution to the acceptance ratio is limited to  $\log(2C)$ . We thereby proceeded with typically large values of  $C$  (such as  $C = 10^6$ ). As the parameter  $\sigma^2$  is the variance of the proposal distribution  $\phi$ , the practitioner should tune it in accordance with the noise level of the data. The parameter requiring the finest calibration is  $\delta$ : the convergence of the algorithm is sensitive to its choice. Dalalyan and Tsybakov (2008, 2012) exhibit the theoretical value  $\delta = n/4\sigma^2$ . This value leads to very good numerical performances, as it has been also noticed by Dalalyan and Tsybakov (2008, 2012); Alquier and Biau (2011). The choice for  $\alpha$  is guided by a similar reasoning to the one for  $C$ . Its contribution to the acceptance ratio is limited to a  $\log(1/\alpha)$  term. The value  $\alpha = 0.25$  was used in the simulations for its apparent good properties. Although it would be computationally costly, a finer calibration through methods such as cross-validation is possible.

Finally and as a general rule, we strongly encourage practitioners to run several chains of unequal lengths and to adjust the number of iterations needed by observing if the empirical risk is stabilized.

**Model 1.**  $n = 200$  and  $S^* = \{1, 2, 3, 4\}$ . This model is similar to Meier, van de Geer and Bühlmann (2009, Section 3, Example 1) and is given by

$$Y_i = \psi_1^*(X_{i1}) + \psi_2^*(X_{i2}) + \psi_3^*(X_{i3}) + \psi_4^*(X_{i4}) + \xi_i,$$

with

$$\begin{aligned} \psi_1^*: x \mapsto -\sin(2x), \quad \psi_2^*: x \mapsto x^3, \quad \psi_3^*: x \mapsto x, \\ \psi_4^*: x \mapsto e^{-x} - e/2, \quad \xi_i \sim \mathcal{N}(0, 0.1), \quad i \in \{1, \dots, n\}. \end{aligned}$$

The covariates are sampled from independent uniform distributions over  $(-1, 1)$ .

**Model 2.**  $n = 200$  and  $S^* = \{1, 2, 3, 4\}$ . As above but correlated. The covariates are sampled from a multivariate gaussian distribution with covariance matrix  $\Sigma_{ij} = 2^{-|i-j|-2}$ ,  $i, j \in \{1, \dots, p\}$ .

**Model 3.**  $n = 200$  and  $S^* = \{1, 2, 3, 4\}$ . This model is similar to Meier, van de Geer and Bühlmann (2009, Section 3, Example 3) and is given by

$$Y_i = 5\psi_1^*(X_{i1}) + 3\psi_2^*(X_{i2}) + 4\psi_3^*(X_{i3}) + 6\psi_4^*(X_{i4}) + \xi_i,$$

with

$$\begin{aligned} \psi_1^* &: x \mapsto x, & \psi_2^* &: x \mapsto 4(x^2 - x - 1), & \psi_3^* &: x \mapsto \frac{\sin(2\pi x)}{2 - \sin(2\pi x)}, \\ \psi_4^* &: x \mapsto 0.1 \sin(2\pi x) + 0.2 \cos(2\pi x) + 0.3 \sin^2(2\pi x) + 0.4 \cos^3(2\pi x) \\ & & & + 0.5 \sin^3(2\pi x), & \xi_i & \sim \mathcal{N}(0, 0.5), \quad i \in \{1, \dots, n\}. \end{aligned}$$

The covariates are sampled from independent uniform distributions over  $(-1, 1)$ .

The results of the simulations are summarized in [Table 1](#) and illustrated by [Figure 1](#) and [Figure 2](#). The reconstruction of the true regression function  $\psi^*$  is achieved even in very high-dimensional situations, pulling up our method at the level of the gold standard Lasso.

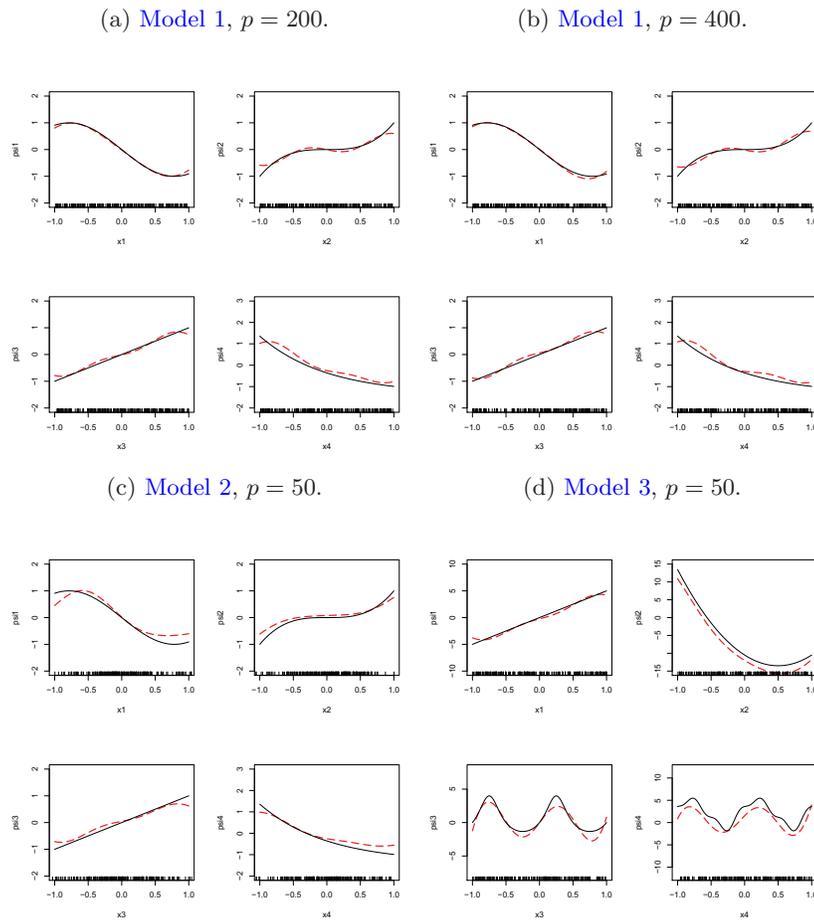


FIG 1. Estimates (red dashed lines) for  $\psi_1^*$ ,  $\psi_2^*$ ,  $\psi_3^*$  and  $\psi_4^*$  (solid black lines). Other estimates (for  $\psi_j^*$ ,  $j \notin \{1, 2, 3, 4\}$ ) are mostly zero.

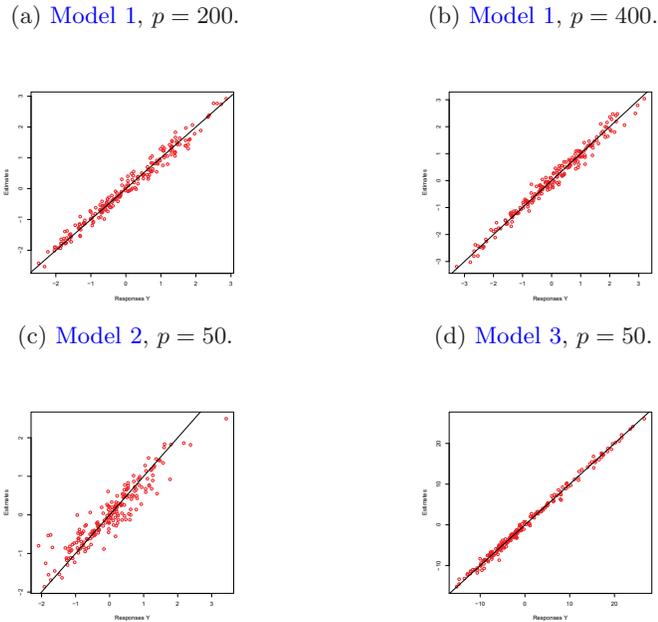


FIG 2. plot of the responses  $Y_1, \dots, Y_n$  against their estimates. The more points on the first bisectrix (solid black line), the better the estimation.

## 5. Proofs

To start the chain of proofs leading to [Theorem 2.1](#), [Theorem 2.2](#) and [Theorem 2.3](#), we recall and prove some lemmas to establish [Theorem 5.1](#) which consists in a general PAC-Bayesian inequality in the spirit of [Catoni \(2004, Theorem 5.5.1\)](#) for classification or [Catoni \(2004, Lemma 5.8.2\)](#) for regression. Note also that [Dalalyan and Tsybakov \(2012, Theorem 1\)](#) provides a similar inequality in the deterministic design case. A salient fact on [Theorem 5.1](#) is that the validity of the oracle inequalities only involves the distribution of the noise variable  $\xi_1$ , and that distribution is independent of the sample size  $n$ .

The proofs of the following two classical results are omitted. [Lemma 5.1](#) is a version of Bernstein's inequality which originates in [Massart \(2007, Proposition 2.19\)](#), whereas [Lemma 5.2](#) appears in [Catoni \(2004, Equation 5.2.1\)](#).

For  $x \in \mathbb{R}$ , denote  $(x)_+ = \max(x, 0)$ . Let  $\mu_1, \mu_2$  be two probabilities. The Kullback-Leibler divergence of  $\mu_1$  with respect to  $\mu_2$  is denoted  $\mathcal{KL}(\mu_1, \mu_2)$  and is

$$\mathcal{KL}(\mu_1, \mu_2) = \begin{cases} \int \log \left( \frac{d\mu_1}{d\mu_2} \right) d\mu_1 & \text{if } \mu_1 \ll \mu_2, \\ \infty & \text{otherwise.} \end{cases}$$

Finally, for any measurable space  $(A, \mathcal{A})$  and any probability  $\pi$  on  $(A, \mathcal{A})$ , denote by  $\mathcal{M}_{+, \pi}^1(A, \mathcal{A})$  the set of probabilities on  $(A, \mathcal{A})$  absolutely continuous with respect to  $\pi$ .

**Lemma 5.1.** *Let  $(T_i)_{i=1}^n$  be independent real-valued variables. Assume that there exist two positive constants  $v$  and  $w$  such that, for any integer  $k \geq 2$ ,*

$$\sum_{i=1}^n \mathbb{E}[(T_i)_+^k] \leq \frac{k!}{2} v w^{k-2}.$$

*Then for any  $\gamma \in (0, \frac{1}{w})$ ,*

$$\mathbb{E} \left[ \exp \left( \gamma \sum_{i=1}^n (T_i - \mathbb{E} T_i) \right) \right] \leq \exp \left( \frac{v \gamma^2}{2(1-w\gamma)} \right).$$

**Lemma 5.2.** *Let  $(A, \mathcal{A})$  be a measurable space. For any probability  $\mu$  on  $(A, \mathcal{A})$  and any measurable function  $h : A \rightarrow \mathbb{R}$  such that  $\int (\exp \circ h) d\mu < \infty$ ,*

$$\log \int (\exp \circ h) d\mu = \sup_{m \in \mathcal{M}_{+, \pi}^1(A, \mathcal{A})} \int h d m - \mathcal{KL}(m, \mu),$$

*with the convention  $\infty - \infty = -\infty$ . Moreover, as soon as  $h$  is upper-bounded on the support of  $\mu$ , the supremum with respect to  $m$  on the right-hand side is reached for the Gibbs distribution  $g$  given by*

$$\frac{dg}{d\mu}(a) = \frac{\exp(h(a))}{\int (\exp \circ h) d\mu}, \quad a \in A.$$

**Theorem 5.1** is valid in the general regression framework. In the proofs of [Lemma 5.3](#), [Lemma 5.4](#), [Lemma 5.5](#) and [Theorem 5.1](#), we consider a general regression function  $\psi^*$ . Denote by  $(\Theta, \mathcal{T})$  a space of functions equipped with a countable generated  $\sigma$ -algebra, and let  $\pi$  be a probability on  $(\Theta, \mathcal{T})$ , referred to as the prior. [Lemma 5.3](#), [Lemma 5.4](#), [Lemma 5.5](#) and [Theorem 5.1](#) follow from the work of [Catoni \(2004\)](#); [Dalalyan and Tsybakov \(2008, 2012\)](#); [Alquier \(2008\)](#); [Alquier and Biau \(2011\)](#). Let  $\delta > 0$  and consider the so-called *posterior* Gibbs transition density  $\rho_\delta$  with respect to  $\pi$ , defined as

$$\rho_\delta(\{\mathbf{x}_i, y_i\}_{i=1}^n, \psi) = \frac{\exp[-\delta r_n(\{\mathbf{x}_i, y_i\}_{i=1}^n, \psi)]}{\int \exp[-\delta r_n(\{\mathbf{x}_i, y_i\}_{i=1}^n, \psi)] \pi(d\psi)}. \quad (5.1)$$

In the following three lemmas, denote by  $\rho$  a *so-called* posterior probability absolutely continuous with respect to  $\pi$ . Let  $\psi$  be a realization of a random variable  $\Psi$  sampled from  $\rho$ .

**Lemma 5.3.** *Let [A1](#) and [A2](#) hold. Set  $w = 8C \max(L, C)$ ,  $\delta \in (0, n/[w + 4(\sigma^2 + C^2)])$  and  $\varepsilon \in (0, 1)$ . Then with  $\mathbb{P}$ -probability at least  $1 - \varepsilon$*

$$R(\psi) - R(\psi^*) \leq \frac{1}{1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \left( R_n(\psi) - R_n(\psi^*) + \frac{\log \frac{d\rho}{d\pi}(\psi) + \log \frac{1}{\varepsilon}}{\delta} \right).$$

*Proof.* Apply [Lemma 5.1](#) to the variables  $T_i$  defined as follow: for any  $\psi \in$ ,

$$T_i = -(Y_i - \psi(\mathbf{X}_i))^2 + (Y_i - \psi^*(\mathbf{X}_i))^2, \quad i \in \{1, \dots, n\}. \quad (5.2)$$

First, let us note that

$$\begin{aligned} R(\psi) - R(\psi^*) &= \mathbb{E}[(Y_1 - \psi(\mathbf{X}_1))^2] - \mathbb{E}[(Y_1 - \psi^*(\mathbf{X}_1))^2] \\ &= \mathbb{E}[(2Y_1 - \psi(\mathbf{X}_1) - \psi^*(\mathbf{X}_1))(\psi^*(\mathbf{X}_1) - \psi(\mathbf{X}_1))] \\ &= \mathbb{E}[(\psi^*(\mathbf{X}_1) - \psi(\mathbf{X}_1)) \mathbb{E}[(2W_1 + \psi^*(\mathbf{X}_1) - \psi(\mathbf{X}_1)) | \mathbf{X}_1]] \\ &= 2 \mathbb{E}[(\psi^*(\mathbf{X}_1) - \psi(\mathbf{X}_1)) \mathbb{E}[\xi_1 | \mathbf{X}_1]] + \mathbb{E}[\psi^*(\mathbf{X}_1) - \psi(\mathbf{X}_1)]^2. \end{aligned}$$

As  $\mathbb{E}[\xi_1 | \mathbf{X}_1] = 0$ ,

$$R(\psi) - R(\psi^*) = \mathbb{E}[\psi^*(\mathbf{X}) - \psi(\mathbf{X})]^2. \quad (5.3)$$

By [\(5.2\)](#), the random variables  $(T_i)_{i=1}^n$  are independent. Using [Lemma 5.1](#), we get

$$\begin{aligned} \sum_{i=1}^n \mathbb{E} T_i^2 &= \sum_{i=1}^n \mathbb{E} [(2Y_i - \psi(\mathbf{X}_i) - \psi^*(\mathbf{X}_i))^2 (\psi(\mathbf{X}_i) - \psi^*(\mathbf{X}_i))^2] \\ &= \sum_{i=1}^n \mathbb{E} \mathbb{E} [(2W_i + \psi^*(\mathbf{X}_i) - \psi(\mathbf{X}_i))^2 (\psi(\mathbf{X}_i) - \psi^*(\mathbf{X}_i))^2 | \mathbf{X}_i]. \end{aligned}$$

Next, using that  $|a + b|^k \leq 2^{k-1}(|a| + |b|)$  for any  $a, b \in \mathbb{R}$  and  $k \in \mathbb{N}^*$ , we get

$$\begin{aligned} \sum_{i=1}^n \mathbb{E} T_i^2 &\leq 2 \sum_{i=1}^n \mathbb{E} [(\psi(\mathbf{X}_i) - \psi^*(\mathbf{X}_i))^2 \mathbb{E} [(4W_i^2 + 4C^2) | \mathbf{X}_i]] \\ &\leq 8(\sigma^2 + C^2) \sum_{i=1}^n \mathbb{E} [(\psi(\mathbf{X}_i) - \psi^*(\mathbf{X}_i))^2] \\ &= 8n(\sigma^2 + C^2) (R(\psi) - R(\psi^*)) \stackrel{\text{def}}{=} v, \end{aligned} \quad (5.4)$$

where we have used [\(5.3\)](#) in the last equation. It follows that for any integer  $k \geq 3$ ,

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}[(T_i)_+^k] &= \sum_{i=1}^n \mathbb{E} \mathbb{E}[(T_i)_+^k | \mathbf{X}_i] \\ &\leq \sum_{i=1}^n \mathbb{E} \mathbb{E} [|2Y_i - \psi(\mathbf{X}_i) - \psi^*(\mathbf{X}_i)|^k |\psi(\mathbf{X}_i) - \psi^*(\mathbf{X}_i)|^k | \mathbf{X}_i] \\ &= \sum_{i=1}^n \mathbb{E} \mathbb{E} [|2W_i + \psi^*(\mathbf{X}_i) - \psi(\mathbf{X}_i)|^k |\psi(\mathbf{X}_i) - \psi^*(\mathbf{X}_i)|^k | \mathbf{X}_i] \\ &\leq 2^{k-1} \sum_{i=1}^n \mathbb{E} \mathbb{E} [(2^k |\xi_i|^k + |\psi^*(\mathbf{X}_i) - \psi(\mathbf{X}_i)|^k) |\psi(\mathbf{X}_i) - \psi^*(\mathbf{X}_i)|^k | \mathbf{X}_i]. \end{aligned}$$

Using that  $|\psi(\mathbf{x}_i) - \psi^*(\mathbf{x}_i)|^k \leq (2C)^{k-2} |\psi(\mathbf{x}_i) - \psi^*(\mathbf{x}_i)|^2$  and (5.3), we get

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}[(T_i)_+^k] &\leq 2^{k-1} \sum_{i=1}^n (2^{k-1} k! \sigma^2 L^{k-2} + (2C)^k) (2C)^{k-2} [R(\psi) - R(\psi^*)] \\ &= \frac{k!}{2} v (2C)^{k-2} \left( \frac{2^{2k-4} \sigma^2 L^{k-2} + \frac{2}{k!} 2^{2k-4} C^k}{\sigma^2 + C^2} \right). \end{aligned}$$

Recalling that  $C > \max(1, \sigma)$  gives

$$\begin{aligned} \frac{2^{2k-4} \sigma^2 L^{k-2} + \frac{2}{k!} 2^{2k-4} C^k}{\sigma^2 + C^2} &\leq \frac{4^{k-2} \sigma^2 L^{k-2}}{2\sigma^2} + \frac{\frac{2}{k!} 4^{k-2} C^k}{C^2} \\ &\leq \frac{1}{2} (4L)^{k-2} + \frac{1}{2} (4C)^{k-2} = [4 \max(L, C)]^{k-2}. \end{aligned}$$

Hence

$$\sum_{i=1}^n \mathbb{E}[(T_i)_+^k] \leq \frac{k!}{2} v w^{k-2}, \quad \text{with } w \stackrel{\text{def}}{=} 8C \max(L, C). \quad (5.5)$$

Applying Lemma 5.1, we obtain, for any real  $\delta \in (0, \frac{n}{w})$ , with  $\gamma = \frac{\delta}{n}$ ,

$$\mathbb{E} \exp[\delta(R_n(\psi^*) - R_n(\psi) + R(\psi) - R(\psi^*))] \leq \exp\left(\frac{v\delta^2}{2n^2(1 - \frac{w\delta}{n})}\right),$$

that is, that for any real number  $\varepsilon \in (0, 1)$ ,

$$\begin{aligned} \mathbb{E} \exp \left[ \delta [R_n(\psi^*) - R_n(\psi)] + \delta [R(\psi) - R(\psi^*)] \left( 1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} \right) \right. \\ \left. - \log \frac{1}{\varepsilon} \right] \leq \varepsilon. \quad (5.6) \end{aligned}$$

Next, we use a standard PAC-Bayesian approach (as developed in Audibert (2004a); Catoni (2004, 2007); Alquier (2008)). For any prior probability  $\pi$  on  $(\Theta, \mathcal{T})$ ,

$$\begin{aligned} \int \mathbb{E} \exp \left[ \delta [R(\psi) - R(\psi^*)] \left( 1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} \right) \right. \\ \left. + \delta [R_n(\psi^*) - R_n(\psi)] - \log \frac{1}{\varepsilon} \right] \pi(d\psi) \leq \varepsilon. \end{aligned}$$

By the Fubini-Tonelli theorem

$$\begin{aligned} \mathbb{E} \int \exp \left[ \delta [R(\psi) - R(\psi^*)] \left( 1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} \right) \right. \\ \left. + \delta [R_n(\psi^*) - R_n(\psi)] - \log \frac{1}{\varepsilon} \right] \pi(d\psi) \leq \varepsilon. \end{aligned}$$

Therefore, for any data-dependent posterior probability measure  $\rho$  absolutely continuous with respect to  $\pi$ , adopting the convention  $\infty \times 0 = 0$ ,

$$\begin{aligned} \mathbb{E} \int \exp \left[ \delta [R(\psi) - R(\psi^*)] \left( 1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} \right) \right. \\ \left. + \delta [R_n(\psi^*) - R_n(\psi)] - \log \frac{d\rho}{d\pi}(\psi) - \log \frac{1}{\varepsilon} \right] \rho(d\psi) \leq \varepsilon. \end{aligned} \quad (5.7)$$

Recalling that  $\mathbb{E}$  stands for the expectation computed with respect to  $\mathbb{P}$ , the integration symbol may be omitted and we get

$$\begin{aligned} \mathbb{E} \exp \left[ \delta [R(\psi) - R(\psi^*)] \left( 1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} \right) \right. \\ \left. + \delta [R_n(\psi^*) - R_n(\psi)] - \log \frac{d\rho}{d\pi}(\psi) - \log \frac{1}{\varepsilon} \right] \leq \varepsilon. \end{aligned}$$

Using the elementary inequality  $\exp(\delta x) \geq \mathbb{1}_{\mathbb{R}_+}(x)$ , we get, with  $\mathbb{P}$ -probability at most  $\varepsilon$

$$\begin{aligned} \left( 1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} \right) [R(\psi) - R(\psi^*)] \geq R_n(\psi) - R_n(\psi^*) \\ + \frac{\log \frac{d\rho}{d\pi}(\psi) + \log \frac{1}{\varepsilon}}{\delta}. \end{aligned}$$

Taking  $\delta < n/[w + 4(\sigma^2 + C^2)]$  implies

$$1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} > 0,$$

and with  $\mathbb{P}$ -probability at least  $1 - \varepsilon$ ,

$$R(\psi) - R(\psi^*) \leq \frac{1}{1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \left( R_n(\psi) - R_n(\psi^*) + \frac{\log \frac{d\rho}{d\pi}(\psi) + \log \frac{1}{\varepsilon}}{\delta} \right).$$

□

**Lemma 5.4.** *Let A1 and A2 hold. Set  $w = 8C \max(L, C)$ ,  $\delta \in (0, n/[w + 4(\sigma^2 + C^2)])$  and  $\varepsilon \in (0, 1)$ . Then with  $\mathbb{P}$ -probability at least  $1 - \varepsilon$*

$$\begin{aligned} \int R_n(\psi) \rho(d\psi) - R_n(\psi^*) \leq \left[ 1 + \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} \right] \left[ \int R(\psi) \rho(d\psi) \right. \\ \left. - R(\psi^*) \right] + \frac{\mathcal{KL}(\rho, \pi) + \log \frac{1}{\varepsilon}}{\delta}. \end{aligned} \quad (5.8)$$

*Proof.* Set  $\psi \in \mathbb{F}$  and  $Z_i = (Y_i - \psi(\mathbf{X}_i))^2 - (Y_i - \psi^*(\mathbf{X}_i))^2$ ,  $i \in \{1, \dots, n\}$ . Since  $Z_i = -T_i$  where  $T_i$  is defined in (5.2), using the same arguments that lead to

(5.6), we get that for any  $\delta \in (0, n/w)$  and  $\varepsilon \in (0, 1)$

$$\begin{aligned} \mathbb{E} \int \exp \left[ -\delta [R(\psi) - R(\psi^*)] \left( 1 + \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} \right) \right. \\ \left. + \delta [R_n(\psi) - R_n(\psi^*)] - \log \frac{d\rho}{d\pi}(\psi) - \log \frac{1}{\varepsilon} \right] \rho(d\psi) \leq \varepsilon. \end{aligned}$$

Using Jensen's inequality, we get

$$\begin{aligned} \mathbb{E} \exp \left[ - \int \left\{ \delta [R(\psi) - R(\psi^*)] \left( 1 + \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} \right) \right. \right. \\ \left. \left. + \delta [R_n(\psi) - R_n(\psi^*)] - \log \frac{d\rho}{d\pi}(\psi) - \log \frac{1}{\varepsilon} \right\} \rho(d\psi) \right] \leq \varepsilon. \end{aligned}$$

Since  $\exp(\delta x) \geq \mathbb{1}_{\mathbb{R}_+}(x)$ , we obtain with  $\mathbb{P}$ -probability at most  $\varepsilon$

$$\begin{aligned} \left[ - \int R(\psi) \rho(d\psi) + R(\psi^*) \right] \left( 1 + \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} \right) + \int R_n(\psi) \rho(d\psi) \\ - R_n(\psi^*) - \frac{\mathcal{KL}(\rho, \pi) + \log \frac{1}{\varepsilon}}{\delta} \geq 0. \end{aligned}$$

Taking  $\delta < n/[w + 4(\sigma^2 + C^2)]$  yields (5.8).  $\square$

**Lemma 5.5.** *Let A1 and A2 hold. Set  $w = 8C \max(L, C)$ ,  $\delta \in (0, n/[w + 4(\sigma^2 + C^2)])$  and  $\varepsilon \in (0, 1)$ . Then with  $\mathbb{P}$ -probability at least  $1 - \varepsilon$*

$$\begin{aligned} \int R(\psi) \rho(d\psi) - R(\psi^*) \leq \frac{1}{1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \left( \int R_n(\psi) \rho(d\psi) - R_n(\psi^*) \right. \\ \left. + \frac{\mathcal{KL}(\rho, \pi) + \log \frac{1}{\varepsilon}}{\delta} \right). \end{aligned}$$

*Proof.* Recall (5.7). By Jensen's inequality,

$$\begin{aligned} \mathbb{E} \exp \left[ \delta \left( \int R(\psi) \rho(d\psi) - R(\psi^*) \right) \left[ 1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta} \right] \right. \\ \left. + \delta \left( R_n(\psi^*) - \int R_n(\psi) \rho(d\psi) \right) - \mathcal{KL}(\rho, \pi) - \log \frac{1}{\varepsilon} \right] \leq \varepsilon. \end{aligned}$$

Using  $\exp(\delta x) \geq \mathbb{1}_{\mathbb{R}_+}(x)$  yields the expected result.  $\square$

**Theorem 5.1.** *Let  $\hat{\psi}$  and  $\hat{\psi}^\Lambda$  be realizations of the Gibbs estimators defined by (2.2)–(2.3), respectively. Let A1 and A2 hold. Set  $w = 8C \max(L, C)$  and  $\delta = n\ell/[w + 4(\sigma^2 + C^2)]$ , for  $\ell \in (0, 1)$ , and let  $\varepsilon \in (0, 1)$ . Then with probability*

at least  $1 - 2\varepsilon$ ,

$$\left. \begin{array}{l} R(\hat{\psi}) - R(\psi^*) \\ R(\hat{\psi}^A) - R(\psi^*) \end{array} \right\} \leq \mathcal{D} \inf_{\rho \in \mathcal{M}_{+, \pi}^1(\Theta, \mathcal{T})} \left\{ \int R(\psi) \rho(d\psi) - R(\psi^*) + \frac{\mathcal{KL}(\rho, \pi) + \log \frac{1}{\varepsilon}}{n} \right\}, \quad (5.9)$$

where  $\mathcal{D}$  is a constant depending only upon  $w$ ,  $\sigma$ ,  $C$  and  $\ell$ .

*Proof.* Recall that the randomized Gibbs estimator  $\hat{\Psi}$  is sampled from  $\rho_\delta$ . Denote by  $\hat{\psi}$  a realization of the variable  $\hat{\Psi}$ . By [Lemma 5.3](#), with  $\mathbb{P}$ -probability at least  $1 - \varepsilon$ ,

$$R(\hat{\psi}) - R(\psi^*) \leq \frac{1}{1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \left( R_n(\hat{\psi}) - R_n(\psi^*) + \frac{\log \frac{d\rho_\delta}{d\pi}(\hat{\psi}) + \log \frac{1}{\varepsilon}}{\delta} \right).$$

Note that

$$\begin{aligned} \log \frac{d\rho_\delta}{d\pi}(\hat{\psi}) &= \log \frac{\exp[-\delta R_n(\hat{\psi})]}{\int \exp[-\delta R_n(\psi)] \pi(d\psi)} \\ &= -\delta R_n(\hat{\psi}) - \log \int \exp[-\delta R_n(\psi)] \pi(d\psi). \end{aligned}$$

Thus, with  $\mathbb{P}$ -probability at least  $1 - \varepsilon$ ,

$$R(\hat{\psi}) - R(\psi^*) \leq \frac{1}{1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \left( -R_n(\psi^*) - \frac{1}{\delta} \log \int \exp[-\delta R_n(\psi)] \pi(d\psi) + \frac{1}{\delta} \log \frac{1}{\varepsilon} \right).$$

By [Lemma 5.2](#), with  $\mathbb{P}$ -probability at least  $1 - \varepsilon$ ,

$$R(\hat{\psi}) - R(\psi^*) \leq \frac{1}{1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \inf_{\rho \in \mathcal{M}_{+, \pi}^1(\Theta, \mathcal{T})} \left( \int R_n(\psi) \rho(d\psi) - R_n(\psi^*) + \frac{\mathcal{KL}(\rho, \pi) + \log \frac{1}{\varepsilon}}{\delta} \right).$$

Finally, by [Lemma 5.4](#), with  $\mathbb{P}$ -probability at least  $1 - 2\varepsilon$ ,

$$R(\hat{\psi}) - R(\psi^*) \leq \frac{1 + \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}}{1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \inf_{\rho \in \mathcal{M}_{+, \pi}^1(\Theta, \mathcal{T})} \left\{ \int R(\psi) \rho(d\psi) - R(\psi^*) + \frac{2}{1 + \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \frac{\mathcal{KL}(\rho, \pi) + \log \frac{1}{\varepsilon}}{\delta} \right\}.$$

Apply [Lemma 5.5](#) with the Gibbs posterior probability defined by [\(5.1\)](#). With  $\mathbb{P}$ -probability at least  $1 - \varepsilon$ ,

$$\int R(\psi)\rho_\delta(d\psi) - R(\psi^*) \leq \frac{1}{1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \left( \int R_n(\psi)\rho_\delta(d\psi) - R_n(\psi^*) + \frac{\mathcal{KL}(\rho_\delta, \pi) + \log \frac{1}{\varepsilon}}{\delta} \right).$$

Note that

$$\begin{aligned} \mathcal{KL}(\rho_\delta, \pi) &= \int \log \frac{\exp[-\delta R_n(\psi)]}{\int \exp[-\delta R_n(\psi)]\pi(d\psi)} \rho_\delta(d\psi) \\ &= -\delta \int R_n(\psi)\rho_\delta(d\psi) - \log \left( \int \exp[-\delta R_n(\psi)]\pi(d\psi) \right). \end{aligned}$$

By [Lemma 5.2](#), with  $\mathbb{P}^{\otimes n}$ -probability at least  $1 - \varepsilon$

$$\int R(\psi)\rho_\delta(d\psi) - R(\psi^*) \leq \frac{1}{1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \inf_{\rho \in \mathcal{M}_{+, \pi}^1(\Theta, \mathcal{T})} \left\{ \int R_n(\psi)\rho(d\psi) - R_n(\psi^*) + \frac{\mathcal{KL}(\rho, \pi) + \log \frac{1}{\varepsilon}}{\delta} \right\}.$$

By [Lemma 5.4](#), with  $\mathbb{P}^{\otimes n}$ -probability at least  $1 - 2\varepsilon$

$$\int R(\psi)\rho_\delta(d\psi) - R(\psi^*) \leq \frac{1 + \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}}{1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \inf_{\rho \in \mathcal{M}_{+, \pi}^1(\Theta, \mathcal{T})} \left\{ \int R(\psi)\rho(d\psi) - R(\psi^*) + \frac{2}{1 + \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} \frac{\mathcal{KL}(\rho, \pi) + \log \frac{1}{\varepsilon}}{\delta} \right\}.$$

As  $R$  is a convex function, applying Jensen's inequality gives

$$\int R(\psi)\rho_\delta(d\psi) \geq R(\hat{\psi}^A).$$

Finally, note that

$$\frac{1 + \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}}{1 - \frac{4\delta(\sigma^2 + C^2)}{n - w\delta}} = 1 + \frac{8\ell(\sigma^2 + C^2)}{(1 - \ell)(w + 4\sigma^2 + 4C^2)}.$$

□

*Proof of [Theorem 2.1](#).* Let  $\rho \in \mathcal{M}_{+, \pi}^1(\Theta, \mathcal{T})$ . For any  $A \in \mathcal{T}$ , note that  $\rho(A) = \sum_{\mathbf{m} \in \mathcal{M}} \rho_{\mathbf{m}}(A)$  where  $\rho_{\mathbf{m}}(\cdot) = \rho(\cdot \cap \Theta_{\mathbf{m}})$ , the trace of  $\rho$  on  $\Theta_{\mathbf{m}}$ . By [Theorem 5.1](#),

with  $\mathbb{P}$ -probability at least  $1 - 2\varepsilon$

$$R(\hat{\psi}) - R(\psi^*) \leq \mathcal{D} \inf_{\mathbf{m} \in \mathcal{M}} \inf_{\rho \in \mathcal{M}_{+, \pi}^1(\Theta, \mathcal{T})} \left\{ \int R(\psi) \rho_{\mathbf{m}}(d\psi) - R(\psi^*) + \frac{\mathcal{KL}(\rho_{\mathbf{m}}, \pi) + \log \frac{1}{\varepsilon}}{n} \right\}. \quad (5.10)$$

Note that for any  $\rho \in \mathcal{M}_{+, \pi}^1(\Theta, \mathcal{T})$  and any  $\mathbf{m} \in \mathcal{M}$ ,

$$\begin{aligned} \mathcal{KL}(\rho_{\mathbf{m}}, \pi) &= \int \log \left( \frac{d\rho_{\mathbf{m}}}{d\pi_{\mathbf{m}}} \right) d\rho_{\mathbf{m}} + \int \log \left( \frac{d\pi_{\mathbf{m}}}{d\pi} \right) d\rho_{\mathbf{m}} \\ &= \mathcal{KL}(\rho_{\mathbf{m}}, \pi_{\mathbf{m}}) + \log(1/\alpha) \sum_{j \in S(\mathbf{m})} m_j + \log \binom{p}{|S(\mathbf{m})|} + \log \left( \frac{1 - \left(\frac{\alpha}{1-\alpha}\right)^{p+1}}{1 - \frac{\alpha}{1-\alpha}} \right). \end{aligned}$$

Next, using the elementary inequality  $\log \binom{n}{k} \leq k \log(ne/k)$  and that  $\frac{\alpha}{1-\alpha} < 1$ ,

$$\begin{aligned} \mathcal{KL}(\rho_{\mathbf{m}}, \pi) &\leq \mathcal{KL}(\rho_{\mathbf{m}}, \pi_{\mathbf{m}}) + \log(1/\alpha) \sum_{j \in S(\mathbf{m})} m_j + |S(\mathbf{m})| \log \left( \frac{pe}{|S(\mathbf{m})|} \right) \\ &\quad + \log \left( \frac{1-\alpha}{1-2\alpha} \right). \end{aligned}$$

We restrict the set of all probabilities absolutely continuous with respect to  $\pi_{\mathbf{m}}$  to uniform probabilities on the ball  $\mathcal{B}_{\mathbf{m}}^1(\mathbf{x}, \zeta)$ , with  $\mathbf{x} \in \mathcal{B}_{\mathbf{m}}^1(0, C)$  and  $0 < \zeta \leq C - \|\theta\|_1$ . Such a probability is denoted by  $\mu_{\mathbf{x}, \zeta}$ . With  $\mathbb{P}$ -probability at least  $1 - 2\varepsilon$ , it yields that

$$\begin{aligned} R(\hat{\psi}) - R(\psi^*) &\leq \mathcal{D} \inf_{\mathbf{m} \in \mathcal{M}} \inf_{\theta \in \mathcal{B}_{\mathbf{m}}^1(0, C)} \inf_{\mu_{\theta, \zeta}, 0 < \zeta \leq C - \|\theta\|_1} \left\{ \int R(\psi_{\bar{\theta}}) \mu_{\theta, \zeta}(d\bar{\theta}) - \right. \\ &\quad \left. R(\psi^*) + \frac{1}{n} \left[ \mathcal{KL}(\mu_{\theta, \zeta}, \pi_{\mathbf{m}}) + \log \frac{1}{\varepsilon} + |S(\mathbf{m})| \log \left( \frac{p}{|S(\mathbf{m})|} \right) + \sum_{j \in S(\mathbf{m})} m_j \right] \right\}. \end{aligned}$$

Next, note that

$$\mathcal{KL}(\mu_{\theta, \zeta}, \pi_{\mathbf{m}}) = \log \left( \frac{V_{\mathbf{m}}(C)}{V_{\mathbf{m}}(\zeta)} \right) = \log \left( \frac{C}{\zeta} \right) \sum_{j \in S(\mathbf{m})} m_j.$$

Note also that

$$\begin{aligned} \int R(\psi_{\bar{\theta}}) \mu_{\theta, \zeta}(d\bar{\theta}) &= \int \mathbb{E} [Y_1 - \psi_{\bar{\theta}}(\mathbf{X}_1)]^2 \mu_{\theta, \zeta}(d\bar{\theta}) \\ &= \int \mathbb{E} [Y_1 - \psi_{\theta}(\mathbf{X}_1) + \psi_{\theta}(\mathbf{X}_1) - \psi_{\bar{\theta}}(\mathbf{X}_1)]^2 \mu_{\theta, \zeta}(d\bar{\theta}), \end{aligned}$$

and

$$\begin{aligned} & \int \mathbb{E} [Y_1 - \psi_\theta(\mathbf{X}_1) + \psi_\theta(\mathbf{X}_1) - \psi_{\bar{\theta}}(\mathbf{X}_1)]^2 \mu_{\theta, \zeta}(\mathrm{d}\bar{\theta}) \\ &= \int R(\psi_\theta) \mu_{\theta, \zeta}(\mathrm{d}\bar{\theta}) + \int \mathbb{E} [\psi_\theta(\mathbf{X}_1) - \psi_{\bar{\theta}}(\mathbf{X}_1)]^2 \mu_{\theta, \zeta}(\mathrm{d}\bar{\theta}) \\ & \quad + 2 \int \mathbb{E} \{ [Y_1 - \psi_\theta(\mathbf{X}_1)] [\psi_\theta(\mathbf{X}_1) - \psi_{\bar{\theta}}(\mathbf{X}_1)] \} \mu_{\theta, \zeta}(\mathrm{d}\bar{\theta}). \end{aligned}$$

Since  $\bar{\theta} \in \mathcal{B}_{\mathbf{m}}^1(\theta, \zeta)$ ,

$$\begin{aligned} & \int \mathbb{E} [\psi_\theta(\mathbf{X}_1) - \psi_{\bar{\theta}}(\mathbf{X}_1)]^2 \mu_{\theta, \zeta}(\mathrm{d}\bar{\theta}) \\ &= \int \mathbb{E} \left[ \sum_{j \in S(\mathbf{m})} \sum_{k=1}^{m_j} (\theta_{jk} - \bar{\theta}_{jk}) \varphi_k(X_{1j}) \right]^2 \mu_{\theta, \zeta}(\mathrm{d}\bar{\theta}) \\ & \leq \|\theta - \bar{\theta}\|_1^2 \max_k \|\varphi_k\|_\infty^2 \leq \zeta^2, \end{aligned}$$

and by the Fubini-Tonelli theorem,

$$\begin{aligned} & 2 \int \mathbb{E} \{ [Y_1 - \psi_\theta(\mathbf{X}_1)] [\psi_\theta(\mathbf{X}_1) - \psi_{\bar{\theta}}(\mathbf{X}_1)] \} \mu_{\theta, \zeta}(\mathrm{d}\bar{\theta}) \\ &= 2 \mathbb{E} \left[ [Y_1 - \psi_\theta(\mathbf{X}_1)] \int [\psi_\theta(\mathbf{X}_1) - \psi_{\bar{\theta}}(\mathbf{X}_1)] \mu_{\theta, \zeta}(\mathrm{d}\bar{\theta}) \right] = 0, \end{aligned}$$

since  $\int \psi_{\bar{\theta}}(\mathbf{X}_1) \mu_{\theta, \zeta}(\mathrm{d}\bar{\theta}) = \psi_\theta(\mathbf{X}_1)$ . Consequently, as

$$\int R(\psi_\theta) \mu_{\theta, \zeta}(\mathrm{d}\bar{\theta}) = R(\psi_\theta),$$

we get

$$\int R(\psi_{\bar{\theta}}) \mu_{\theta, \zeta}(\mathrm{d}\bar{\theta}) \leq R(\psi_\theta) + \zeta^2.$$

So with  $\mathbb{P}$ -probability at least  $1 - 2\varepsilon$ ,

$$\begin{aligned} R(\hat{\psi}) - R(\psi^*) &\leq \mathcal{D} \inf_{\mathbf{m} \in \mathcal{M}} \inf_{\theta \in \mathcal{B}_{\mathbf{m}}^1(0, C)} \inf_{\mu_{\theta, \zeta}, 0 < \zeta \leq C - \|\theta\|_1} \left\{ R(\psi_\theta) + \zeta^2 - R(\psi^*) \right. \\ & \left. + \frac{1}{n} \left[ \log(C/\zeta) \sum_{j \in S(\mathbf{m})} m_j + \log \frac{1}{\varepsilon} + |S(\mathbf{m})| \log \left( \frac{p}{|S(\mathbf{m})|} \right) + \sum_{j \in S(\mathbf{m})} m_j \right] \right\}. \end{aligned}$$

The function  $t \mapsto t^2 + \log(C/t) \sum_{j \in S(\mathbf{m})} m_j/n$  is convex. Its minimum is unique and is reached for  $t = [\sum_{j \in S(\mathbf{m})} m_j/(2n)]^{1/2}$ . With  $\mathbb{P}$ -probability at least  $1 - 2\varepsilon$ ,

$$R(\hat{\psi}) - R(\psi^*) \leq \mathcal{D} \inf_{\mathbf{m} \in \mathcal{M}} \inf_{\theta \in \mathcal{B}_{\mathbf{m}}^1(0, C)} \left\{ R(\psi_\theta) - R(\psi^*) + |S(\mathbf{m})| \frac{\log(p/|S(\mathbf{m})|)}{n} + \frac{\log(n)}{n} \sum_{j \in S(\mathbf{m})} m_j + \frac{\log(1/\varepsilon)}{n} \right\},$$

where  $\mathcal{D}$  is a constant depending only on  $w, \sigma, C, \ell$  and  $\alpha$ . As the same inequality holds for  $\hat{\psi}^A$ , this concludes the proof.  $\square$

*Proof of Theorem 2.2.* Recall Theorem 2.1. A3 gives

$$R(\psi_\theta) - R(\psi^*) = \int (\psi_\theta(\mathbf{x}) - \psi^*(\mathbf{x}))^2 d\mathcal{P}(\mathbf{x}) \leq B \int (\psi_\theta(\mathbf{x}) - \psi^*(\mathbf{x}))^2 d\mathbf{x}.$$

For any  $\mathbf{m} \in \mathcal{M}$ , define

$$\psi_{\mathbf{m}}^* = \sum_{j \in S^*} \sum_{k=1}^{m_j} \theta_{jk}^* \varphi_k.$$

To proceed, we need to check that the projection of  $\theta^*$  onto model  $\mathbf{m}$  lies in  $\mathcal{B}_{\mathbf{m}}^1(0, C)$ , i.e.,

$$\sum_{j \in S^*} \sum_{k=1}^{m_j} |\theta_{jk}^*| \leq C.$$

Using the Cauchy-Schwarz inequality, we get

$$\begin{aligned} \sum_{j \in S^*} \sum_{k=1}^{m_j} |\theta_{jk}^*| &= \sum_{j \in S^*} \sum_{k=1}^{m_j} k^{r_j} |\theta_{jk}^*| k^{-r_j} \\ &\leq \sum_{j \in S^*} \left[ \sqrt{\sum_{k=1}^{m_j} k^{2r_j} (\theta_{jk}^*)^2} \sqrt{\sum_{k=1}^{m_j} k^{-2r_j}} \right]. \end{aligned}$$

Since for any  $t \geq 1$ ,  $\sum_{k=1}^{m_j} k^{-2t} \leq \sum_{k=1}^{\infty} k^{-2t} = \pi^2/6$ , the previous inequality yields

$$\sum_{j \in S^*} \sum_{k=1}^{m_j} |\theta_{jk}^*| \leq \frac{\pi}{\sqrt{6}} \sum_{j \in S^*} \sqrt{d_j} \leq C.$$

Recalling (5.3) and A3, for a  $\mathbf{m} \in \mathcal{M}$  we may now write that

$$\begin{aligned} \inf_{\theta \in \Theta_{\mathbf{m}}} R(\psi_\theta) - R(\psi^*) &\leq R(\psi_{\mathbf{m}}^*) - R(\psi^*) \leq B \int (\psi^*(\mathbf{x}) - \psi_{\mathbf{m}}^*(\mathbf{x}))^2 d\mathbf{x} \\ &= B \int \left( \sum_{j \in S^*} \sum_{k=1+m_j}^{\infty} \theta_{jk}^* \varphi_k(\mathbf{x}) \right)^2 d\mathbf{x}. \end{aligned}$$

As  $\{\varphi_k\}_{k=1}^\infty$  forms an orthogonal basis,

$$\begin{aligned} B \int \left( \sum_{j \in S^*} \sum_{k=1+m_j}^\infty \theta_{jk}^* \varphi_k(\mathbf{x}) \right)^2 d\mathbf{x} &= B \sum_{j \in S^*} \sum_{k=1+m_j}^\infty (\theta_{jk}^*)^2 \\ &\leq B \sum_{j \in S^*} d_j (1+m_j)^{-2r_j}, \end{aligned}$$

where the normalizing numerical factors are included in the now generic constant  $B$ . As a consequence, with  $\mathbb{P}$ -probability at least  $1 - 2\varepsilon$ ,

$$\begin{aligned} R(\hat{\psi}) - R(\psi^*) &\leq \mathcal{D} \inf_{\mathbf{m} \in \mathcal{M}} \left\{ B \sum_{j \in S^*} \left\{ d_j (1+m_j)^{-2r_j} + \frac{m_j}{n} \log(n) \right\} \right. \\ &\quad \left. + |S^*| \frac{\log(p/|S^*|)}{n} + \frac{\log(1/\varepsilon)}{n} \right\}, \end{aligned}$$

where  $\mathcal{D}$  is the same constant as in [Theorem 2.1](#). For any  $r \geq 2$ , the function  $t \mapsto d_j(1+t)^{-2r_j} + \frac{\log(n)}{n}t$  is convex and admits a minimum in  $(\frac{\log(n)}{2r_j d_j n})^{-\frac{1}{2r_j+1}} - 1$ .

Accordingly, choosing  $m_j \sim (\frac{\log(n)}{2r_j d_j n})^{-\frac{1}{2r_j+1}} - 1$  yields that with  $\mathbb{P}$ -probability at least  $1 - 2\varepsilon$ ,

$$R(\hat{\psi}) - R(\psi^*) \leq \mathcal{D} \left\{ \sum_{j \in S^*} d_j^{\frac{1}{2r_j+1}} \left( \frac{\log(n)}{2nr_j} \right)^{\frac{2r_j}{2r_j+1}} + |S^*| \frac{\log\left(\frac{p}{|S^*|}\right)}{n} + \frac{\log(1/\varepsilon)}{n} \right\},$$

where  $\mathcal{D}$  is a constant depending only on  $\alpha, w, \sigma, C, \ell$  and  $B$ , and that ends the proof.  $\square$

*Proof of [Theorem 2.3](#).* The proof is similar to the proof of [Theorem 2.1](#). From [\(5.10\)](#) and for any  $\rho \in \mathcal{M}_{+, \pi}^1(\Theta, \mathcal{T})$  and any  $\mathbf{m} \in \mathcal{M}$ ,

$$\begin{aligned} \mathcal{KL}(\rho_{\mathbf{m}}, \pi) &= \mathcal{KL}(\rho_{\mathbf{m}}, \pi_{\mathbf{m}}) + \log(1/\alpha)|S(\mathbf{m})| + \log\left(\frac{p}{|S(\mathbf{m})|}\right) \\ &\quad + \log\left(\frac{1 - \left(\alpha \frac{1-\alpha^{K+1}}{1-\alpha}\right)^{p+1}}{1 - \alpha \frac{1-\alpha^{K+1}}{1-\alpha}}\right) + \sum_{j \in S(\mathbf{m})} \log\left(\frac{K}{|S(\mathbf{m}_j)|}\right). \end{aligned}$$

Using the elementary inequality  $\log\binom{n}{k} \leq k \log(ne/k)$  and that  $\alpha \frac{1-\alpha^{K+1}}{1-\alpha} \in (0, 1)$  since  $\alpha < 1/2$ ,

$$\begin{aligned} \mathcal{KL}(\rho_{\mathbf{m}}, \pi) &\leq \mathcal{KL}(\rho_{\mathbf{m}}, \pi_{\mathbf{m}}) + |S(\mathbf{m})| \left[ \log(1/\alpha) + \log\left(\frac{pe}{|S(\mathbf{m})|}\right) \right] \\ &\quad + \sum_{j \in S(\mathbf{m})} |S(\mathbf{m}_j)| \log\left(\frac{Ke}{|S(\mathbf{m}_j)|}\right) + \log\left(\frac{1-\alpha}{1-2\alpha}\right). \end{aligned}$$

Thus with  $\mathbb{P}$ -probability at least  $1 - 2\varepsilon$ ,

$$R(\hat{\psi}) - R(\psi^*) \leq \mathcal{D} \inf_{\mathbf{m} \in \mathcal{M}} \inf_{\theta \in \mathcal{B}_{\mathbf{m}}^1(0, C)} \inf_{\mu_{\theta, \zeta}, 0 < \zeta \leq C - \|\theta\|_1} \left\{ R(\psi_{\theta}) + \zeta^2 - R(\psi^*) \right. \\ \left. + \frac{1}{n} \left[ \log(C/\zeta) + \log(K) \right] \sum_{j \in S(\mathbf{m})} |S(\mathbf{m}_j)| + \log \frac{1}{\varepsilon} + |S(\mathbf{m})| \log \left( \frac{p}{|S(\mathbf{m})|} \right) \right\}.$$

Hence with  $\mathbb{P}$ -probability at least  $1 - 2\varepsilon$ ,

$$\left. \begin{array}{l} R(\hat{\psi}) - R(\psi^*) \\ R(\hat{\psi}^A) - R(\psi^*) \end{array} \right\} \leq \mathcal{D} \inf_{\mathbf{m} \in \mathcal{M}} \inf_{\theta \in \mathcal{B}_{\mathbf{m}}^1(0, C)} \left\{ R(\psi_{\theta}) - R(\psi^*) \right. \\ \left. + |S(\mathbf{m})| \frac{\log(p/|S(\mathbf{m})|)}{n} + \frac{\log(nK)}{n} \sum_{j \in S(\mathbf{m})} |S(\mathbf{m}_j)| + \frac{\log(1/\varepsilon)}{n} \right\},$$

where  $\mathcal{D}$  is a numerical constant depending upon  $w$ ,  $\sigma$ ,  $C$ ,  $\ell$  and  $\alpha$ .  $\square$

### Acknowledgements

The authors are grateful to Gérard Biau and Éric Moulines for their constant implication, and to Christophe Giraud and Taiji Suzuki for valuable insights and comments. They also thank an anonymous referee and an associate editor for providing constructive and helpful remarks.

### References

- ALQUIER, P. (2006). Transductive and Inductive Adaptive Inference for Regression and Density Estimation PhD thesis, Université Paris 6 - UPMC.
- ALQUIER, P. (2008). PAC-Bayesian Bounds for Randomized Empirical Risk Minimizers. *Mathematical Methods of Statistics* **17** 279–304. [arXiv:0712.1698v3](#) [MR2483458](#)
- ALQUIER, P. and BIAU, G. (2011). Sparse Single-Index Model. To appear in *Journal of Machine Learning Research*. [arXiv:1101.3229v2](#)
- ALQUIER, P. and LOUNICI, K. (2011). PAC-Bayesian Theorems for Sparse Regression Estimation with Exponential Weights. *Electronic Journal of Statistics* **5** 127–145. [MR2786484](#)
- AUDIBERT, J.-Y. (2004a). Aggregated estimators and empirical complexity for least square regression. *Annales de l'Institut Henri Poincaré: Probabilités et Statistiques* **40** 685–736. [MR2096215](#)
- AUDIBERT, J.-Y. (2004b). Théorie statistique de l'apprentissage: une approche PAC-Bayésienne PhD thesis, Université Paris 6 - UPMC.
- AUDIBERT, J.-Y. (2009). Fast learning rates in statistical inference through aggregation. *The Annals of Statistics* **37** 1591–1646. [MR2533466](#)



- HASTIE, T., TIBSHIRANI, R. and FRIEDMAN, J. (2009). *The Elements of Statistical Learning – Data mining, Inference, and Prediction*, Second ed. Springer. [MR2722294](#)
- HÄRDLE, W. K. (1990). *Applied nonparametric regression*. Cambridge University Press. [MR1161622](#)
- KOLTCHINSKII, V. and YUAN, M. (2010). Sparsity in multiple kernel learning. *The Annals of Statistics* **38** 3660–3695. [MR2766864](#)
- MARIN, J.-M. and ROBERT, C. P. (2007). *Bayesian Core: A Practical Approach to Computational Bayesian Statistics*. Springer. [MR2289769](#)
- MASSART, P. (2007). *Concentration Inequalities and Model Selection*. *École d'Été de Probabilités de Saint-Flour XXXIII – 2003*. Springer. [MR2319879](#)
- MCALLESTER, D. A. (1999). Some PAC-Bayesian Theorems. *Machine Learning* **37** 355–363. [MR1811587](#)
- MEIER, L., VAN DE GEER, S. A. and BÜHLMANN, P. (2009). High-dimensional additive modeling. *The Annals of Statistics* **37** 3779–3821. [arXiv:0806.4115](#) [MR2572443](#)
- MEINSHAUSEN, N. and YU, B. (2009). Lasso-type recovery of sparse representations for high-dimensional data. *The Annals of Statistics* **37** 246–270. [arXiv:0806.0145v2](#) [MR2488351](#)
- MEYN, S. and TWEEDIE, R. L. (2009). *Markov Chains and Stochastic Stability*, 2nd ed. Cambridge University Press. [MR2509253](#)
- PETRALIAS, A. (2010). Bayesian model determination and nonlinear threshold volatility models PhD thesis, Athens University of Economics and Business.
- PETRALIAS, A. and DELLAPORTAS, P. (2012). An MCMC model search algorithm for regression problems. *Journal of Statistical Computation and Simulation* **0** 1-19.
- R CORE TEAM (2012). R: A Language and Environment for Statistical Computing, Vienna, Austria ISBN 3-900051-07-0. <http://www.R-project.org/>
- RASKUTTI, G., WAINWRIGHT, M. J. and YU, B. (2012). Minimax-optimal rates for sparse additive models over kernel classes via convex programming. *Journal of Machine Learning Research* **13** 389-427. [MR2913704](#)
- RAVIKUMAR, P., LAFFERTY, J., LIU, H. and WASSERMAN, L. (2009). Sparse additive models. *Journal of the Royal Statistical Society, Series B* **71** 1009–1030. [arXiv:0711.4555v2](#) [MR2750255](#)
- RIGOLLET, P. (2006). Inégalités d'oracle, agrégation et adaptation PhD thesis, Université Paris 6 - UPMC.
- RIGOLLET, P. and TSYBAKOV, A. B. (2012). Sparse estimation by exponential weighting. *Statistical Science* **27** 558-575.
- SHAWE-TAYLOR, J. and WILLIAMSON, R. C. (1997). A PAC analysis of a Bayes estimator. In *Proceedings of the 10th annual conference on Computational Learning Theory* 2–9. ACM.
- STONE, C. J. (1985). Additive regression and other nonparametric models. *The Annals of Statistics* **13** 689–705. [MR0790566](#)
- SUZUKI, T. (2012). PAC-Bayesian Bound for Gaussian Process Regression and Multiple Kernel Additive Model. In *Proceedings of the 25th annual conference on Computational Learning Theory*.

- SUZUKI, T. and SUGIYAMA, M. (2012). Fast learning rates of Multiple kernel learning: trade-off between sparsity and smoothness. Submitted. [arXiv.org/abs/1203.0565v1](https://arxiv.org/abs/1203.0565v1)
- TIBSHIRANI, R. (1996). Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B* **58** 267–288. [MR1379242](#)
- TSYBAKOV, A. B. (2009). *Introduction to Nonparametric Estimation. Statistics*. Springer. [MR2724359](#)
- VAN DE GEER, S. A. (2008). High-dimensional generalized linear models and the Lasso. *The Annals of Statistics* **36** 614–645. [MR2396809](#)