

Large variance Gaussian priors in Bayesian nonparametric estimation: a maxiset approach

Florent Autin, Dominique Picard, Vincent Rivoirard

▶ To cite this version:

Florent Autin, Dominique Picard, Vincent Rivoirard. Large variance Gaussian priors in Bayesian nonparametric estimation: a maxiset approach. Mathematical Methods of Statistics, 2006, 15 (4), pp.349-373. hal-00634287

HAL Id: hal-00634287 https://hal.science/hal-00634287

Submitted on 20 Oct 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.

Large variance Gaussian priors in Bayesian nonparametric estimation: a maxiset approach^{*}

Florent Autin, Dominique Picard and Vincent Rivoirard

CNRS- Universités de Paris X-Nanterre, Paris VII and Paris XI

August 18, 2005

Abstract

In this paper we compare wavelet Bayesian rules taking into account the sparsity of the signal with priors which are combinations of a Dirac mass with a standard distribution properly normalized. To perform these comparisons, we take the maxiset point of view: i. e. we consider the set of functions which are well estimated (at a prescribed rate) by each procedure. We especially consider the standard cases of Gaussian and heavy-tailed priors. We show that if heavy-tailed priors have extremely good maxiset behavior compared to traditional Gaussian priors, considering large variance Gaussian priors (LVGP) leads to equally successful maxiset behavior. Moreover, these LVGP can be constructed in an adaptive way. We also show, using comparative simulations results that large variance Gaussian priors have very good numerical performances, confirming the maxiset prediction, and providing the advantage of a high simplicity from the computational point of view.

1 Introduction

Bayesian techniques have now become very popular to estimate signals decomposed on wavelet bases. Many authors have built Bayes estimates showing, from the practical point of view, impressive properties especially to estimate inhomogeneous signals. Most of the simulations show that these procedures seriously outperform classical procedures and in particular thresholding procedures. See for instance, [Chipman et al., 1997], [Abramovich et al., 1998], [Clyde et al., 1998], [Johnstone and Silverman, 1998], [Vidakovic, 1998], [Clyde and George, 1998], or [Clyde and George, 2000] who discussed the choice of the Bayes model to capture the sparsity of the signal to be estimated and the choice of the Bayes rule (and among others, posterior mean or median). We also refer the reader to the very complete review paper of [Antoniadis et al., 2001] which provide descriptions and comparisons of various Bayesian wavelet shrinkage and wavelet thresholding estimators.

^{*}Key Words and Phrases: minimax, maxiset, nonparametric estimation, Bayesian methods. AMS 2000 Subject Classification: 62G05, 62G07, 62G20.

To capture the sparsity of the signal, the most common models introduce priors on the wavelet coefficients of the following form:

$$\beta_{jk} \sim \pi_{j,\epsilon} \gamma_{j,\epsilon} + (1 - \pi_{j,\epsilon}) \delta(0), \tag{1}$$

where $0 \leq \pi_{j,\epsilon} \leq 1$, $\delta(0)$ is a point mass at zero and the β_{jk} 's are independent. The nonzero part of the prior $\gamma_{j,\epsilon}$ is assumed to be the dilation of a fixed symmetric, positive, unimodal and continuous density γ :

$$\gamma_{j,\epsilon}(\beta_{jk}) = \frac{1}{\tau_{j,\epsilon}} \gamma\left(\frac{\beta_{jk}}{\tau_{j,\epsilon}}\right),$$

where the dilation parameter $\tau_{j,\epsilon}$ is positive. The parameter $\pi_{j,\epsilon}$ can be interpreted as the proportion of non negligible coefficients. We also introduce the parameter

$$w_{j,\epsilon} = \frac{\pi_{j,\epsilon}}{1 - \pi_{j,\epsilon}}.$$

When the signal is sparse, most of the $w_{j,\epsilon}$'s are small. These priors or very close forms have extensively been used by the authors cited above and especially by [Abramovich et al., 2004], [Johnstone and Silverman, 2004a] and [Johnstone and Silverman, 2004b]. To complete the definition of the prior model, we have to fix the hyperparameters $\tau_{j,\epsilon}$ and $w_{j,\epsilon}$. Finally the density γ will play a very important role. The most popular choice for γ is the normal density. It is also the density giving rise to the easiest procedures from a computational point of view. However heavy-tailed priors have proved also to work extremely well.

From the minimax point of view, recent works have studied these Bayes procedures and it has been proved that Bayes rules can achieve optimal rates of convergence. [Abramovich et al., 2004] investigated theoretical performance of the procedures introduced by [Abramovich et al., 1998], considering priors of the form quoted above with some particular choice of the hyperparameters. For the mean squared error, they proved that the non adaptive posterior mean and posterior median achieve optimal rates up to a logarithmic factor on the Besov spaces $\mathcal{B}_{p,q}^s$ when $p \geq 2$. When p < 2, these estimators show less impressive properties since they only behave as linear estimates. As [Abramovich et al., 2004], [Johnstone and Silverman, 2004a] and [Johnstone and Silverman, 2004b] investigated minimax properties of Bayes rules, with priors based on heavy-tailed distributions and they considered an empirical Bayes setting. In this case, the posterior mean and median turn out to be optimal for the whole scale of Besov spaces. Other more sophisticated results concerning minimax properties of Bayes rules have been established by [Zhang, 2002].

Hence, summarizing the results cited above, the minimax results seem to indicate that Bayes procedures have comparable results to thresholding estimates at least on the range of Besov spaces, but also seem to show a preference to heavy-tailed priors.

The main goal of the paper is to push a little further this type of comparison on Bayesian procedures by adopting the maxiset point of view. In particular, since Gaussian priors have very interesting properties from the computational point of view, one of our motivations was to answer the following question: Are Gaussian priors always outperformed by heavy-tailed priors ? And quite happily, one of our results will be to show that if some Bayesian procedures using Gaussian priors behave quite unwell (in terms of maxisets as it was the case in terms of minimax rates) compared to those with heavy tails, it is nevertheless possible to attain a very good maxiset behavior, among procedures based on Gaussian priors. We prove that this can only be achieved under the condition that the hyperparameter $\tau_{j,\epsilon}$ is "large". Under this assumption, the density $\gamma_{j,\epsilon}$ is then more spread around 0, mimicking in some ways the behavior of a distribution with heavy-tails. Moreover, we prove that these procedures can be built in an adaptive way: their construction does not depend on the specified regularity or sparsity of the function at hand.

As these Bayesian procedures with large variance Gaussian priors have not been much studied in the literature yet, we investigated their behavior also from a practical point of view and show a comparative simulations study with many standard and Bayesian procedures of the literature. As can be seen in our last section, such estimators turn out to have excellent numerical performances.

Let us only recall here that the maxiset point of view consists in determining the set of all functions which can be estimated at a specified rate of convergence for a specified procedure. Exhibiting maxisets of different estimation rules allow to say that a procedure is more powerful than another one if its maxiset is larger.

The results that have been obtained up to now, using the maxiset point of view, are very promising since they generally show that the maxisets of well-known procedures are spaces which are well understandable and easily interpretable sets. They have the advantage of being generally less pessimistic and seem also to enjoy the important advantage of giving theoretical claims which are often closer to the practical (simulations) situation, than other theoretical results (such as minimax rates).

The following section details the model and Bayesian rules we are going to consider in the paper. The third section recalls the definition of maxisets, and briefly details some results obtained in the area, to allow a comparison with the results obtained later for Bayesian rules. The forth section investigates maxisets of standard Bayesian rules: first the 'small variance' Gaussian priors, then the heavy-tailed priors. The fifth section is devoted to large variance Gaussian priors, and the last section details the simulations results.

2 Model and Bayesian rules

For sake of simplicity, we will consider a white noise setting: $X_{\epsilon}(.)$ is a random measure satisfying on [0, 1] the following equation:

$$X_{\epsilon}(dt) = f(t)dt + \epsilon W(dt)$$

where $0 < \epsilon < 1$ is the noise level and f is a function defined on [0,1], W(.) is a Brownian motion on [0,1]. As usual, to connect with the standard framework of sequences of experiments we put $\epsilon = n^{-1/2}$.

 $\{\psi_{jk}(\cdot), j \geq -1, k \in \mathbb{Z}\}\$ is a compactly supported wavelet basis of $\mathbb{L}_2([0,1])$, such that any

 $f \in \mathbb{L}_2([0,1])$ can be represented as:

$$f = \sum_{j \ge -1} \sum_{k} \beta_{jk} \psi_{jk}$$

where $\beta_{jk} = (f, \psi_{jk})_{\mathbb{L}_2}$. As usual, ψ_{-1k} denote the translations of the scaling function, and ψ_{jk} , for $j \geq 0$ are the dilations and translations of the wavelet function. The model is reduced to a sequence space model if we put: $y_{jk} = X_{\epsilon}(\psi_{jk}) = \int f\psi_{jk} + \epsilon Z_{jk}$ where Z_{jk} are i.i.d $\mathcal{N}(0, 1)$. Let us note that at each level $j \geq 0$, the number of non-zero wavelet coefficients is smaller than or equal to $2^j + l_{\psi} - 1$, where l_{ψ} is the maximal size of the supports of the scaling function and the wavelet. So, there exists a constant S_{ψ} such that at each level $j \geq -1$, there are less than or equal to $S_{\psi} \times 2^j$ coefficients to be estimated. In the sequel, we shall not distinguish between f and $\beta = (\beta_{jk})_{jk}$ its sequence of wavelet coefficients.

As explained in the introduction, we consider the following priors: The β_{jk} 's are independent random variables with the following distribution,

$$\beta_{jk} \sim \pi_{j,\epsilon} \gamma_{j,\epsilon} + (1 - \pi_{j,\epsilon}) \delta(0),$$
 (2)

$$\gamma_{j,\epsilon}(\beta_{jk}) = \frac{1}{\tau_{j,\epsilon}} \gamma\left(\frac{\beta_{jk}}{\tau_{j,\epsilon}}\right)$$

$$w_{j,\epsilon} = \frac{\pi_{j,\epsilon}}{1 - \pi_{j,\epsilon}}$$
(3)

where $0 \le \pi_{j,\epsilon} \le 1$, $\delta(0)$ is the Dirac mass at 0, γ is a fixed symmetric, positive, unimodal and continuous density, $\tau_{j,\epsilon}$ is positive.

2.1 Gaussian priors

Let us consider the case where γ is the Gaussian density, which is the most classical choice. In this case, we easily derive the Bayes rules of β_{jk} associated with the l^1 -loss and the l^2 -loss, respectively the 'a posteriori' median and mean:

$$\breve{\beta}_{jk} = \operatorname{Med}(\beta_{jk}|y_{jk}) = \operatorname{sign}(y_{jk}) \max(0, \xi_{jk}), \tag{4}$$

$$\tilde{\beta}_{jk} = \mathbb{E}(\beta_{jk}|y_{jk}) = \frac{b_j}{1 + \eta_{jk}} y_{jk}, \qquad (5)$$

where,

$$\begin{aligned} \xi_{jk} &= b_j |y_{jk}| - \epsilon \sqrt{b_j} \Phi^{-1} \left(\frac{1 + \min(\eta_{jk}, 1)}{2} \right), \\ b_j &= \frac{\tau_{j,\epsilon}^2}{\epsilon^2 + \tau_{j,\epsilon}^2}, \\ \eta_{jk} &= \frac{1}{w_{j,\epsilon}} \frac{\sqrt{\epsilon^2 + \tau_{j,\epsilon}^2}}{\epsilon} \exp\left(-\frac{\tau_{j,\epsilon}^2 y_{jk}^2}{2\epsilon^2(\epsilon^2 + \tau_{j,\epsilon}^2)} \right), \end{aligned}$$

and Φ is the normal cumulative distributive function.

To study the properties of such rules, it is interesting to make use of their shrinkage properties. Let us recall that $\hat{\beta}$ is said to be a shrinkage rule if $y_{jk} \longrightarrow \hat{\beta}_{jk}$ is antisymmetric, increasing on $(-\infty, +\infty)$ and

$$0 \le \hat{eta}_{jk} \le y_{jk}, \quad \forall \ y_{jk} \ge 0.$$

Both rules quoted above obviously are shrinkage rules. We also note that $\check{\beta}_{jk}$ is zero whenever y_{jk} falls in an implicitly defined interval $[-\lambda_{j,\epsilon}, \lambda_{j,\epsilon}]$.

We will first consider the following very classical form for the hyperparameters:

$$\tau_{j,\epsilon}^2 = c_1 2^{-\alpha j}, \quad \pi_{j,\epsilon} = \min(1, c_2 2^{-bj}),$$
(6)

where c_1 , c_2 , α and b are positive constants. This particular form was suggested by [Abramovich et al., 1998] and then used by [Abramovich et al., 2004]. A nice interpretation was provided by these authors who explained how α , b, c_1 and c_2 can be derived for applications.

Our second part will be concerned with large variance rules. In this case, we will consider hyperparameters of the form

$$\tau_{j,\epsilon} = \tau(\epsilon) \text{ and } w_{j,\epsilon} = w(\epsilon)$$
 (7)

with specified conditions on the functions τ and w.

Remark 1. An alternative for eliciting these hyperparameters consists in using empirical Bayes methods and EM algorithm (see [Clyde and George, 1998], [Clyde and George, 2000] or [Johnstone and Silverman, 1998]).

2.2 Heavy-tailed priors

For sake of comparison, we will also consider priors where the density γ is no longer Gaussian. We assume that there exist two positive constants M and M_1 such that

$$\sup_{\beta \ge M_1} \left| \frac{d}{d\beta} \log \gamma(\beta) \right| = M < \infty.$$
(8)

The hypothesis (8) means that the tails of γ have to be exponential or heavier. Indeed, under (8), we have:

 $\forall u \ge M_1, \quad \gamma(u) \ge \gamma(M_1) \exp(-M(u - M_1)).$

In the minimax approach of [Johnstone and Silverman, 2004a] and [Johnstone and Silverman, 2004b], the priors also verified (8). To complete the prior model, we assume that:

$$\tau_{j,\epsilon} = \epsilon, \ w_{j,\epsilon} = w(\epsilon) \to 0, \ \text{as } \epsilon \to 0$$
(9)

and w a positive continuous function. Using these assumptions, the following proposition describes the properties of the posterior median and mean:

Proposition 1. Under the conditions (8) and (9) the estimates $\check{\beta}_{jk}^{HT} = Med(\beta_{jk}|y_{jk})$ and $\tilde{\beta}_{jk}^{HT} = \mathbb{E}(\beta_{jk}|y_{jk})$ are shrinkage rules. Moreover, $\check{\beta}_{jk}^{HT}$ is a thresholding rule: there exists \check{t}_{ϵ} such that

$$\check{\beta}_{jk}^{HT} = 0 \iff |y_{jk}| \le \check{t}_{\epsilon},$$

where the threshold \check{t}_{ϵ} verifies for ϵ small enough, $\check{t}_{\epsilon} \geq \epsilon \sqrt{2 \log(1/w(\epsilon))}$ and

$$\lim_{\epsilon \to 0} \frac{\check{t}_{\epsilon}}{\epsilon \sqrt{2\log(1/w(\epsilon))}} = 1$$

Proof: The first point has been established by [Johnstone and Silverman, 2004a] and [Johnstone and Silverman, 2004b]. The second point is an immediate consequence of Proposition 3 of [Rivoirard, 2003].

3 Maxisets and associated functional spaces

Let us first briefly recall the definition of maximum sets. We consider a sequence of models $\mathcal{E}_n = \{P_{\theta}^n, \theta \in \Theta\}$, where the P_{θ}^n 's are probability distributions on the measurable spaces Ω_n , and Θ is the set of parameters. We also consider a sequence of estimates \hat{q}_n of a quantity $q(\theta)$ associated with this sequence of models, a loss function $\rho(\hat{q}_n, q(\theta))$, and a rate of convergence α_n tending to 0. Then, we define the **maxiset** associated with the sequence \hat{q}_n , the loss function ρ , the rate α_n and the constant T as the following set:

$$MS(\hat{q}_n, \rho, \alpha_n)(T) = \{\theta \in \Theta, \sup_{n \in \mathcal{B}} \mathbb{E}^n_{\theta} \rho(\hat{q}_n, q(\theta))(\alpha_n)^{-1} \le T\}$$

The focus in this domain has mainly been on the nonparametric situation. Let us briefly mention the differences with the minimax point of view. In this latter, we fix a set of functions and look at the worst performances of estimators. Here, instead of a priori fixing a (functional) set such as a Hölder, Sobolev or Besov ball, we choose to settle the problem in a wider context: The parameter set Θ can be very large, such as the set of bounded, measurable functions. Then, the maxiset is associated with the procedure in a more genuine way since it only depends on the model and the estimation rule at hand.

As explained more in details later in this section, there already exist very interpretable results about maxisets. For instance, it has been established in [Kerkyacharian and Picard, 1993] that the maxisets of linear kernel methods are in fact Besov spaces under fairly reasonable conditions on the kernel, whereas the maxisets of thresholding estimates (see [Cohen et al., 2001]) are Lorentz spaces reflecting extremely well the practical observation that wavelet thresholding performs well when the number of wavelet coefficients is small. It has also been observed (see [Kerkyacharian and Picard, 2002]) that there is a deep connection between oracle inequalities and maxisets, in the sense that verifying an oracle inequality is equivalent to proving that the maxiset of the procedure automatically contains a minimal set associated to the oracle.

Although both settings seem quite different, still there is a deep parallel between maxisets and minimax theory. For instance, facing a particular situation, the standard procedure to prove that a set B is the maxiset usually consists (exactly as in minimax theory) in two steps: first showing that $B \subset MS(\hat{q}_n, \rho, \alpha_n)(T)$, but this is generally obtained using similar arguments as for proving upper bound inequalities in minimax setting since it is simply needed to prove that if $\theta \in B$ then $\mathbb{E}^n_{\theta}\rho(\hat{q}_n, q(\theta)) \leq T\alpha_n$. The gain of the maxiset setting is probably that the second inclusion $MS(\hat{q}_n, \rho, \alpha_n)(T) \subset B$ is often much simpler than proving lower bound for minimax rates over complicated spaces.

3.1 Functional spaces

In this paper, for simplicity, we shall restrict to the case where ρ is the square of the \mathbb{L}_2 norm, even though a large majority of the results can be extended to more general losses. For this study, we need to introduce the following classes of functions which are of typical use in maxiset theory.

3.1.1 Besov and weak Besov spaces

Here, we give definitions of the Besov and weak Besov spaces depending on the wavelet basis. However, as is established in [Meyer, 1990] and [Cohen et al., 2001], most of them also have different definitions proving that this dependence in the basis is not crucial at all.

Definition 1. Let s > 0 and R > 0. A function $f = \sum_{j=-1}^{+\infty} \sum_k \beta_{jk} \psi_{jk} \in \mathbb{L}_2([0,1])$ belongs to the **Besov ball** $\mathcal{B}^s_{p,\infty}(R)$, if and only if:

$$\left[\sup_{j\geq -1} 2^{j(s+\frac{1}{2}-\frac{1}{p})p} \sum_{k} |\beta_{jk}|^{p}\right]^{1/p} \leq R.$$

Note that, when p = 2, f belongs to $\mathcal{B}_{2,\infty}^s$ if and only if:

$$\sup_{J \ge -1} 2^{2Js} \sum_{j \ge J} \sum_{k} \beta_{jk}^2 < +\infty.$$
 (10)

This characterization is often used in the sequel. Recall that the class of Besov spaces $\mathcal{B}_{p,\infty}^s$ provides a useful tool to classify wavelet decomposed signals in function of their regularity and sparsity properties. See [Donoho et al., 1995], [Donoho and Johnstone, 1994] or [Johnstone, 1994]. Roughly speaking, regularity increases when s increases whereas sparsity increases when p decreases. Especially, the spaces with indices p < 2 are of particular interest since they describe very wide classes of inhomogeneous but sparse functions. To model sparsity, a very convenient and natural tool consists in introducing the following particular class of Lorentz spaces that are in addition directly connected to the estimation procedures considered in this paper.

Definition 2. Let 0 < r < 2 and R > 0. A function $f = \sum_{j=-1}^{+\infty} \sum_k \beta_{jk} \psi_{jk} \in \mathbb{L}_2([0,1])$ belongs to the weak Besov ball $W_r(R)$ if and only if:

$$\left[\sup_{\lambda>0}\lambda^{r-2}\sum_{j\geq -1}\sum_{k}\beta_{jk}^{2}I\{|\beta_{jk}|\leq\lambda\}\right]^{1/2}\leq R.$$

It is not difficult to prove (see [Cohen et al., 2001]) that

$$f \in W_r \Leftrightarrow \sup_{\lambda > 0} \lambda^r \sum_j I\{|\beta_{jk}| > \lambda\} < \infty,$$

We have, in particular,

$$\sup_{\lambda>0} \lambda^{r} \sum_{j,k} I\{|\beta_{jk}| > \lambda\} \le \frac{2^{2-r}}{1-2^{-r}} \sup_{\lambda>0} \lambda^{r-2} \sum_{j\geq -1} \sum_{k} \beta_{jk}^{2} I\{|\beta_{jk}| \le \lambda\},$$
(11)

which shows the natural relationship between sparsity and weak Besov spaces and the connection with the regular Besov spaces introduced above. If \subsetneq denotes the strict inclusion between two functional spaces, the previous embeddings are not difficult to show (see for instance [Meyer, 1990], [Kerkyacharian and Picard, 2002] or [Rivoirard, 2004a]):

$$\mathcal{B}_{p,\infty}^s \subsetneq \mathcal{B}_{2,\infty}^s \subsetneq W_{\frac{2}{1+2s}}, \qquad \text{if} \quad s > 0, \ p > 2, \tag{12}$$

$$\mathcal{B}^s_{p,\infty} \subsetneq W_{\frac{2}{1+2s}}, \qquad \text{if} \quad s > 0, \ p < 2.$$
(13)

3.2 First connections between the spaces and maxiset results

In the present setting of white noise model, [Rivoirard, 2004a] proved that, the maxisets of linear estimates for polynomial rates of convergence of the form $\epsilon^{4s/(1+2s)}$ are Besov spaces $\mathcal{B}_{2,\infty}^s$. A similar result in the context of kernel estimates was established in [Kerkyacharian and Picard, 1993]. If we introduce the classical hard and soft thresholding rules:

$$\hat{f}^T = \sum_{-1 \le j < j_{\epsilon}} \sum_{k} y_{jk} I\{|y_{jk}| > mt_{\epsilon}\} \psi_{jk},$$
$$\hat{f}^S = \sum_{-1 \le j < j_{\epsilon}} \sum_{k} \left(1 - \frac{mt_{\epsilon}}{|y_{jk}|}\right) I\{|y_{jk}| > mt_{\epsilon}\} y_{jk} \psi_{jk},$$

with m a positive constant, $j_{\epsilon} \in \mathbb{N}$ such that

$$t_{\epsilon} = \epsilon \sqrt{\log(1/\epsilon)} \tag{14}$$

$$2^{-j_{\epsilon}} \leq t_{\epsilon}^2 < 2^{1-j_{\epsilon}}, \tag{15}$$

(which will be denoted in the sequel by $2^{j_{\epsilon}} \sim t_{\epsilon}^{-2}$). Under mild conditions, [Kerkyacharian and Picard, 2000] proved:

$$MS(\hat{f}^{T}, \|.\|_{2}^{2}, (\epsilon\sqrt{\log(1/\epsilon}))^{4s/(1+2s)}) = \mathcal{B}_{2,\infty}^{s/(2s+1)} \cap W_{\frac{2}{2s+1}}$$

A similar result is obtained for the soft thresholding rule \hat{f}^S .

Remark 2. The embeddings mentioned above ((12) and (13)) give clear informations about the respective performances of linear procedures and thresholding rules, which have been extensively confirmed by practical results. In particular, one can observe that the spaces $\mathcal{B}_{p,\infty}^s$ for p < 2 are never included into the maxisets of the linear procedures $(\mathcal{B}_{2,\infty}^s)$, while they are included into the maxisets of the linear procedures $(\mathcal{B}_{2,\infty}^s)$ under fairly wide conditions.

Notation: If \subset denotes the inclusion between two spaces and for \mathcal{A} , a given space, the following notations:

$$\begin{array}{rcl} MS(\hat{f}_{\epsilon}, \|.\|_{2}^{2}, \lambda_{\epsilon}) & \subset & \mathcal{A} \\ (resp.) & \mathcal{A} & \subset & MS(\hat{f}_{\epsilon}, \|.\|_{2}^{2}, \lambda_{\epsilon}) \end{array}$$

will mean in the sequel

$$\forall M \exists M', MS(\hat{f}_{\epsilon}, \|.\|_{2}^{2}, \lambda_{\epsilon})(M) \subset \mathcal{A}(M')$$

$$(resp.) \quad \forall M' \exists M, \mathcal{A}(M') \subset MS(\hat{f}_{\epsilon}, \|.\|_{2}^{2}, \lambda_{\epsilon})(M),$$

where M and M' respectively denote the radii of balls of $MS(\hat{f}_{\epsilon}, \|.\|_2^2, \lambda_{\epsilon})$ and \mathcal{A} .

4 Maxisets results for 'heavy-tailed' and 'small variance Gaussian' priors

4.1 Maxisets results for small variance Gaussian priors

Let us consider now, the Bayesian rules with Gaussian priors as explained in section 2.1, and especially those verifying conditions (6), as introduced in [Abramovich et al., 1998] and studied in [Abramovich et al., 2004].

Theorem 1. With the previous choice for the hyperparameters, for s > 0 and $\beta^0 \in \{\breve{\beta}, \tilde{\beta}\}$,

- $\alpha > 2s+1$ implies $\mathcal{B}_{p,\infty}^s \not\subset MS(\beta^0, \|.\|_2^2, t_{\epsilon}^{4s/(1+2s)})$ for any $1 \le p \le \infty$,
- $\alpha = 2s + 1$ implies $\mathcal{B}_{p,\infty}^s \not\subset MS(\beta^0, \|.\|_2^2, t_{\epsilon}^{4s/(1+2s)})$ if p < 2.

Remark 3. Theorem 1 is established for the rate $t_{\epsilon}^{4s/(1+2s)}$ but it can be generalized for any rate of convergence of the form $\epsilon^{4s/(1+2s)}(\log(1/\epsilon))^m$, with $m \ge 0$. The results established in Theorem 1 (if we for example refer to remark 2) proved that the performances obtained by these rules are obviously outperformed by thresholding rules. It is worthwhile to notice in addition, that their behavior are (just like linear procedures) highly non robust regarding the tuning constant α . The behavior of these rules turns out to be very comparable to linear rule as is confirmed in Appendix where more details about the maxisets of these procedure are given.

The proof of Theorem 1 is based on the following result:

Proposition 2. If $\beta \in MS(\beta^0, \|.\|_2^2, t_{\epsilon}^{4s/(1+2s)})$, then there exists a constant C such that, for ϵ small enough:

$$\sum_{j,k} \beta_{jk}^2 I\{\tau_{j,\epsilon}^2 \le \epsilon^2\} I\{|\beta_{jk}| > t_\epsilon\} \le C t_\epsilon^{\frac{4s}{1+2s}}.$$
(16)

Proof of Proposition 2: Here we shall distinguish the cases of the posterior mean and median. The posterior median can be written as follows:

$$\breve{\beta}_{jk} = \operatorname{sign}(y_{jk})(b_j|y_{jk}| - g(\epsilon, \tau_{j,\epsilon}, y_{jk})),$$

with $0 \leq g(\epsilon, \tau_{j,\epsilon}, y_{jk}) \leq b_j |y_{jk}|$. Let us assume that $b_j |y_{jk} - \beta_{jk}| \leq (1 - b_j) |\beta_{jk}|/2$ and $\tau_{j,\epsilon}^2 \leq \epsilon^2$, so $b_j \leq 1/2$. First, let us suppose that $y_{jk} \geq 0$ so $\beta_{jk} \geq 0$. If $\beta_{jk} \geq 0$, then

$$\begin{split} |\breve{\beta_{jk}} - \beta_{jk}| &= |b_j(y_{jk} - \beta_{jk}) - (1 - b_j)\beta_{jk} - g(\epsilon, \tau_{j,\epsilon}, y_{jk})| \\ &= (1 - b_j)\beta_{jk} - b_j(y_{jk} - \beta_{jk}) + g(\epsilon, \tau_{j,\epsilon}, y_{jk}) \\ &\geq \frac{1}{2}(1 - b_j)\beta_{jk} \\ &\geq \frac{1}{4}\beta_{jk}. \end{split}$$

If $\beta_{jk} \leq 0$, then

$$|\check{\beta_{jk}} - \beta_{jk}| \ge \frac{1}{4}|\beta_{jk}|.$$

The case $y_{jk} \leq 0$ is handled by using similar arguments and the particular form of the posterior median. So, we obtain:

$$\mathbb{E}(\beta_{jk} - \beta_{jk})^2 I\{\tau_{j,\epsilon}^2 \le \epsilon^2\} \ge \frac{1}{16} \beta_{jk}^2 \mathbb{P}(b_j | y_{jk} - \beta_{jk} | \le (1 - b_j) |\beta_{jk}|/2) I\{\tau_{j,\epsilon}^2 \le \epsilon^2\} \\ \ge \frac{1}{16} \beta_{jk}^2 \mathbb{P}(|y_{jk} - \beta_{jk}| \le |\beta_{jk}|/2) I\{\tau_{j,\epsilon}^2 \le \epsilon^2\} \\ \ge \frac{1}{16} \beta_{jk}^2 (1 - \mathbb{P}(|y_{jk} - \beta_{jk}| > |\beta_{jk}|/2)) I\{\tau_{j,\epsilon}^2 \le \epsilon^2\}.$$

Using the large deviations inequalities for the Gaussian variables, we obtain for ϵ small enough:

$$\mathbb{E}(\tilde{\beta_{jk}} - \beta_{jk})^2 I\{\tau_{j,\epsilon}^2 \le \epsilon^2\} I\{|\beta_{jk}| > t_\epsilon\} \ge \frac{1}{16} \beta_{jk}^2 (1 - \mathbb{P}(|y_{jk} - \beta_{jk}| > t_\epsilon/2)) I\{\tau_{j,\epsilon}^2 \le \epsilon^2\} I\{|\beta_{jk}| > t_\epsilon\} \\ \ge \frac{1}{32} \beta_{jk}^2 I\{\tau_{j,\epsilon}^2 \le \epsilon^2\} I\{|\beta_{jk}| > t_\epsilon\}.$$

This implies (16).

For the posterior mean, we have:

$$\mathbb{E}(\tilde{\beta_{jk}} - \beta_{jk})^{2} = \mathbb{E}\left(\frac{b_{j}}{1 + \eta_{jk}}(y_{jk} - \beta_{jk}) - (1 - \frac{b_{j}}{1 + \eta_{jk}})\beta_{jk}\right)^{2} \\
\geq \frac{1}{4}\mathbb{E}\left((1 - \frac{b_{j}}{1 + \eta_{jk}})\beta_{jk}\right)^{2}I\left\{\frac{b_{j}}{1 + \eta_{jk}}|y_{jk} - \beta_{jk}| \le (1 - \frac{b_{j}}{1 + \eta_{jk}})|\beta_{jk}|/2\right\}.$$

So, we obtain:

$$\mathbb{E}(\tilde{\beta_{jk}} - \beta_{jk})^{2} I\{\tau_{j,\epsilon}^{2} \leq \epsilon^{2}\} \geq \frac{1}{16} \beta_{jk}^{2} \mathbb{P}(|y_{jk} - \beta_{jk}| \leq |\beta_{jk}|/2) I\{\tau_{j,\epsilon}^{2} \leq \epsilon^{2}\} \\
\geq \frac{1}{16} \beta_{jk}^{2} (1 - \mathbb{P}(|y_{jk} - \beta_{jk}| > |\beta_{jk}|/2)) I\{\tau_{j,\epsilon}^{2} \leq \epsilon^{2}\}.$$

Finally, using similar arguments as those used for the posterior median, we obtain (16). Proposition 2 is proved. \Box

Proof of Theorem 1: Let us first investigate the case $\alpha > 2s + 1$. Let us take β such that all the β_{jk} 's are zero, except 2^j coefficients at each level j that are equal to $2^{-j(s+\frac{1}{2})}$. Then, $\beta \in \mathcal{B}_{p,\infty}^s$. Since $\tau_{j,\epsilon}^2 = c_1 2^{-j\alpha}$, if we put $2^{J_{\alpha}} \sim c_1^{\frac{1}{\alpha}} \epsilon^{-\frac{2}{\alpha}}$ and $2^{J_s} \sim t_{\epsilon}^{\frac{-2}{2s+1}}$, we observe that asymptotically $J_{\alpha} < J_s$. So, for ϵ small enough:

$$\sum_{j,k} \beta_{jk}^2 I\{\tau_{j,\epsilon}^2 \le \epsilon^2\} I\{|\beta_{jk}| > t_\epsilon\} = \sum_{\substack{J_\alpha \le j < J_s \\ \ge c\epsilon^{\frac{4s}{\alpha}},}} 2^{-2js}$$

with c a positive constant. Using Proposition 2, β does not belong to $MS(\beta^0, \|.\|_2^2, t_{\epsilon}^{4s/(1+2s)})$.

Let us then investigate the case $\alpha = 2s + 1$. Let us take β such that all the β_{jk} 's are zero, except 1 coefficient at each level j that is equal to $2^{-j(s+\frac{1}{2}-\frac{1}{p})}$. Then, $\beta \in \mathcal{B}_{p,\infty}^s$. Similarly, we put $2^{J_{\alpha}} \sim c_1^{\frac{1}{\alpha}} \epsilon^{-\frac{2}{\alpha}}$ and $2^{\tilde{J}_s} \sim t_{\epsilon}^{-1/(s+\frac{1}{2}-\frac{1}{p})}$, we observe that asymptotically $J_{\alpha} < \tilde{J}_s$. So, for ϵ small enough:

$$\sum_{j,k} \beta_{jk}^2 I\{\tau_{j,\epsilon}^2 \le \epsilon^2\} I\{|\beta_{jk}| > t_\epsilon\} = \sum_{\substack{J_\alpha \le j < \tilde{J}_s \\ \ge \tilde{c} \epsilon^{4(s+\frac{1}{2}-\frac{1}{p})/\alpha},}$$

with \tilde{c} a positive constant. Using Proposition 2, β does not belong to $MS(\beta^0, \|.\|_2^2, t_{\epsilon}^{4s/(1+2s)})$, since p < 2.

4.2 Heavy-tailed priors

Let us consider now the case of priors verifying the condition (8) and (9). If we set,

$$\check{f}_{\epsilon}^{HT} = \sum_{j < j_{\epsilon}} \sum_{k} \check{\beta}_{jk}^{HT} \psi_{jk}, \quad \check{\beta}_{jk}^{HT} = \operatorname{Med}(\beta_{jk} | y_{jk}), \tag{17}$$

and

$$\tilde{f}_{\epsilon}^{HT} = \sum_{j < j_{\epsilon}} \sum_{k} \tilde{\beta}_{jk}^{HT} \psi_{jk}, \quad \tilde{\beta}_{jk}^{HT} = \mathbb{E}(\beta_{jk} | y_{jk}), \tag{18}$$

where j_{ϵ} is such that $2^{j_{\epsilon}} \sim t_{\epsilon}^{-2}$, using results of Proposition 1, we expect these procedures to mimic classical thresholding rules from the maxiset point of view, at least when the posterior median is considered. Indeed Theorems 2, 3, 4 and 5 established by [Rivoirard, 2003] lead to the following result.

Theorem 2. Let s > 0. We suppose that there exist two positive constants ρ_1 and ρ_2 such that for $\epsilon > 0$ small enough,

$$\epsilon^{\rho_1} \le w(\epsilon) \le \epsilon^{\rho_2}$$

Then, we have:

$$MS(f_{\epsilon}^{0}, \|.\|_{2}^{2}, (\epsilon\sqrt{\log(1/\epsilon}))^{4s/(1+2s)}) = \mathcal{B}_{2,\infty}^{s/(2s+1)} \cap W_{\frac{2}{2s+1}},$$

where $f_{\epsilon}^0 \in {\{\tilde{f}_{\epsilon}^{HT}, \check{f}_{\epsilon}^{HT}\}}$, as soon as $\rho_2 \geq 16$ for the posterior median and $\rho_2 \geq 64$ for the posterior mean.

So, the performances achieved by adaptive Bayesian procedures based on heavy-tailed prior densities are similar to those of classical non linear procedures in the maxiset framework. In particular, they obviously outperform the previous small variance Bayesian procedures from the maxiset point of view.

5 Gaussian priors with large variance

The previous subsection has shown the power of the Bayes procedures built from heavy-tailed prior models in the maxiset setting. The goal of this section is to answer the following questions. Are heavy-tailed priors unavoidable? Is it possible to build Gaussian priors leading to procedures with maxiset properties comparable to the heavy-tailed methods discussed above ? Moreover, can we imagine to contruct these procedures in such a way that it automatically adapts to the regularity of the function (adaptivity property). In other words, if γ is the Gaussian density, does there exist an adaptive choice of the hyperparameters $\pi_{j,\epsilon}$ and $w_{j,\epsilon}$ such that

$$MS(f_{\epsilon}^{0}, \|.\|_{2}^{2}, (\epsilon\sqrt{\log(1/\epsilon)})^{4s/(1+2s)}) = \mathcal{B}_{2,\infty}^{s/(2s+1)} \cap W_{\frac{2}{2s+1}}$$

This is a very important issue since calculation using Gaussian priors are mostly direct and obviously much easier than heavy-tailed priors. The answers are provided by the following theorem 3.

Let us consider the following estimates:

$$\check{f}_{\epsilon}^{LV} = \sum_{j < j_{\epsilon}} \sum_{k} \check{\beta}_{jk} \psi_{jk}, \quad \check{\beta}_{jk} = \operatorname{Med}(\beta_{jk} | y_{jk}), \tag{19}$$

and

$$\tilde{f}_{\epsilon}^{LV} = \sum_{j < j_{\epsilon}} \sum_{k} \tilde{\beta}_{jk} \psi_{jk}, \quad \tilde{\beta}_{jk} = \mathbb{E}(\beta_{jk} | y_{jk}), \tag{20}$$

(Recall that the posterior mean and median are given in (5) and (4)), with the following choice of hyperparameters

$$\tau_{j,\epsilon} = \tau(\epsilon) \text{ and } w_{j,\epsilon} = w(\epsilon)$$
 (21)

Theorem 3. We consider the prior model (1), where γ is the Gaussian density. We assume that $\tau_{j,\epsilon} = \tau(\epsilon)$ and $w_{j,\epsilon} = w(\epsilon)$ are independent of j with w a continuous positive function. We consider \check{f}_{ϵ} and \tilde{f}_{ϵ} introduced in (19) and (20). If

$$1 + \epsilon^{-2} \tau(\epsilon)^2 = t_\epsilon^{-1}$$

and there exist q_1 and q_2 such that for ϵ small enough

$$\epsilon^{q_1} \le w(\epsilon) \le \epsilon^{q_2},$$

we have:

$$MS(f_{\epsilon}^{0}, \|.\|_{2}^{2}, (\epsilon\sqrt{\log(1/\epsilon)})^{4s/(1+2s)}) = \mathcal{B}_{2,\infty}^{s/(2s+1)} \cap W_{\frac{2}{2s+1}},$$

where $f_{\epsilon}^0 \in \{\tilde{f}_{\epsilon}, \check{f}_{\epsilon}\}$ as soon as $q_2 > 63/2$ for the posterior median and $q_2 \geq 65/2$ for the posterior mean.

Unlike the previous choice ($\tau_{j,\epsilon}^2 = \epsilon^2$ or $\tau_{j,\epsilon}^2 = 2^{-j\alpha}$), here we impose a "larger" variance. It is the key point of the proof of Theorem 3. In a sense, we re-create the heavy tails by increasing the variance. The proof of Theorem 3 essentially relies on the following proposition.

Proposition 3. Let s > 0 and $\varpi_{jk}(\epsilon)$ a sequence of random weights lying in [0,1]. We assume that there exist positive constants c, m and $K(\varpi)$ such that for any $\epsilon > 0$,

$$\hat{\beta}(\epsilon) = (\varpi_{jk}(\epsilon)y_{jk})_{jk}$$

is a shrinkage rule verifying for any ϵ ,

$$\varpi_{jk}(\epsilon) = 0, \quad a.e. \quad \forall \ j \ge j_{\epsilon} \ with \ 2^{j_{\epsilon}} \sim t_{\epsilon}^{-2}, \quad \forall \ k,$$
(22)

$$|y_{jk}| \le mt_{\epsilon} \Rightarrow \varpi_{jk}(\epsilon) \le ct_{\epsilon}, \quad a.e. \quad \forall j < j_{\epsilon}, \ \forall k,$$
(23)

$$(1 - \varpi_{jk}(\epsilon)) \le K(\varpi) \left(\frac{t_{\varepsilon}}{|y_{jk}|} + t_{\epsilon} \right), \quad a.e. \quad \forall j < j_{\epsilon}, \ \forall k.$$

$$(24)$$

and let

$$\hat{f}_{\epsilon} = \sum_{j < j_{\epsilon}} \sum_{k} \varpi_{jk}(\epsilon) y_{jk} \psi_{jk}.$$

Let $f \in \mathcal{B}_{2,\infty}^{\frac{s}{1+2s}} \cap W_{\frac{2}{1+2s}}$ and let us note

$$\|f\|_{B^{\frac{s}{1+2s}}_{2,\infty}}^2 = \sup_{J \ge -1} 2^{2Js} \sum_{j \ge J} \sum_k \beta_{jk}^2 < \infty,$$

and

$$\|f\|_{W_{\frac{2}{1+2s}}}^2 = \sup_{\lambda>0} \lambda^{r-2} \sum_{j\geq -1} \sum_k \beta_{jk}^2 I\{|\beta_{jk}| \le \lambda\} < \infty.$$

Then, as soon as $m \ge 8$, we have the following inequality:

$$\mathbb{E}\|\hat{f}_{\epsilon} - f\|_{2}^{2} \leq \left[4c^{2}S_{\psi} + 4(1 + K(\varpi)^{2})\|f\|_{2}^{2} + 4\sqrt{3}S_{\psi} + 2(2^{\frac{4s}{1+2s}} + 2^{\frac{-4s}{1+2s}})m^{\frac{4s}{1+2s}}\|f\|_{W_{\frac{2}{1+2s}}}^{2} + \frac{8m^{-2/1+2s}}{(1-2^{-2/1+2s})}(1 + 8K(\varpi)^{2})\|f\|_{W_{\frac{2}{1+2s}}}^{2} + \|f\|_{B_{2,\infty}^{\frac{1}{1+2s}}}^{2}\right]t_{\epsilon}^{\frac{4s}{1+2s}},$$

and

$$\mathcal{B}_{2,\infty}^{\vec{1+2s}} \cap W_{\frac{2}{1+2s}} \subset MS(\hat{f}_{\epsilon}, \|.\|_{2}^{2}, t_{\epsilon}^{4s/(1+2s)}).$$

Proof of Proposition 3: Using (22), we have

$$\mathbb{E}\|\hat{f}_{\epsilon} - f\|_{2}^{2} = \mathbb{E}\|\sum_{j < j_{\epsilon}, k} (\varpi_{jk}(\epsilon)y_{jk} - \beta_{jk})\psi_{j,k}\|_{2}^{2} + \sum_{j \ge j_{\epsilon}, k} \beta_{jk}^{2}$$

The second term is a bias term bounded by $t_{\epsilon}^{\frac{4s}{1+2s}} ||f||_{B_{2,\infty}^{\frac{s}{1+2s}}}^2$. We split $\mathbb{E} \sum_{j < j_{\epsilon},k} (\varpi_{jk}(\epsilon)y_{jk} - \beta_{jk})^2$ into 2(A+B) with

$$A = \mathbb{E} \sum_{j < j_{\epsilon}, k} [\varpi_{jk}(\epsilon)^{2} (y_{jk} - \beta_{jk})^{2} + (1 - \varpi_{jk}(\epsilon))^{2} \beta_{jk}^{2}] I\{|y_{jk}| \le mt_{\epsilon}\},\$$

$$B = \mathbb{E} \sum_{j < j_{\epsilon}, k} [\varpi_{jk}(\epsilon)^{2} (y_{jk} - \beta_{jk})^{2} + (1 - \varpi_{jk}(\epsilon))^{2} \beta_{jk}^{2}] I\{|y_{jk}| > mt_{\epsilon}\}.$$

Again, we split A into $A_1 + A_2$, and using (23)

$$A_1 = \mathbb{E} \sum_{j < j_{\epsilon}, k} \varpi_{jk}(\epsilon)^2 (y_{jk} - \beta_{jk})^2 I\{|y_{jk}| \le mt_{\epsilon}\}$$

$$\le c^2 S_{\psi} 2^{j_{\epsilon}} t_{\epsilon}^2 \epsilon^2$$

$$\le 2c^2 S_{\psi} t_{\epsilon}^2.$$

$$\begin{aligned} A_{2} &= \mathbb{E} \sum_{j < j_{\epsilon}, k} (1 - \varpi_{jk}(\epsilon))^{2} \beta_{jk}^{2} I\{|y_{jk}| \leq mt_{\epsilon}\} \\ &\leq \mathbb{E} \sum_{j < j_{\epsilon}, k} \beta_{jk}^{2} I\{|y_{jk}| \leq mt_{\epsilon}\} [I\{|\beta_{jk}| \leq 2mt_{\epsilon}\} + I\{|\beta_{jk}| > 2mt_{\epsilon}\}] \\ &\leq (2mt_{\epsilon})^{4s/1+2s} \|f\|_{W_{\frac{2}{1+2s}}}^{2} + \sum_{j < j_{\epsilon}, k} \beta_{jk}^{2} \mathbb{P}(|\beta_{jk} - y_{jk}| \geq mt_{\epsilon}) \\ &\leq (2mt_{\epsilon})^{4s/1+2s} \|f\|_{W_{\frac{2}{1+2s}}}^{2} + \|f\|_{2}^{2} \epsilon^{m^{2}/2} \\ &\leq (2mt_{\epsilon})^{4s/1+2s} \|f\|_{W_{\frac{2}{1+2s}}}^{2} + \|f\|_{2}^{2} t_{\epsilon}^{2}. \end{aligned}$$

We have used here the concentration property of the Gaussian distribution and the fact that $m^2 \ge 4$.

$$B := B_1 + B_2$$

= $\mathbb{E} \sum_{j < j_{\epsilon}, k} [\varpi_{jk}(\epsilon)^2 (y_{jk} - \beta_{jk})^2 + (1 - \varpi_{jk}(\epsilon))^2 \beta_{jk}^2] I\{|y_{jk}| > mt_{\epsilon}\} [I\{|\beta_{jk}| \le mt_{\epsilon}/2\} + I\{|\beta_{jk}| > mt_{\epsilon}/2\}].$

For B_1 we use the Schwartz inequality:

$$\mathbb{E}(y_{jk} - \beta_{jk})^2 I\{|y_{jk} - \beta_{jk}| > mt_{\epsilon}/2\} \le (\mathbb{P}(|y_{jk} - \beta_{jk}| > mt_{\epsilon}/2))^{1/2} (\mathbb{E}(y_{jk} - \beta_{jk})^4)^{1/2}.$$

Now, observing that $\mathbb{E}(y_{jk} - \beta_{jk})^4 = 3\epsilon^4$ and that $\mathbb{P}(|y_{jk} - \beta_{jk}| > mt_{\epsilon}/2) \le \epsilon^{\frac{m^2}{8}}$, we have for $m^2 \ge 32$:

$$B_{1} \leq \sqrt{3} \sum_{j < j_{\epsilon}, k} \epsilon^{2} I\{|\beta_{jk}| \leq mt_{\epsilon}/2\} \epsilon^{\frac{m^{2}}{16}} + \sum_{j < j_{\epsilon}, k} \beta_{jk}^{2} I\{|\beta_{jk}| \leq mt_{\epsilon}/2\}$$
$$\leq 2\sqrt{3} S_{\psi} t_{\epsilon}^{2} + \left(\frac{m}{2} t_{\epsilon}\right)^{4s/1+2s} \|f\|_{W_{\frac{s}{1+2s}}}^{2}.$$

For B_2 , we use (11) to obtain

$$\begin{split} B_2 &= \mathbb{E} \sum_{j < j_{\epsilon}, k} [\varpi_{jk}(\epsilon)^2 (y_{jk} - \beta_{jk})^2 + (1 - \varpi_{jk}(\epsilon))^2 \beta_{jk}^2] I\{|y_{jk}| > mt_{\epsilon}\} I\{|\beta_{jk}| > mt_{\epsilon}/2\} \\ &\leq \sum_{j < j_{\epsilon}, k} [\epsilon^2 I\{|\beta_{jk}| > mt_{\epsilon}/2\} + B_3 \\ &\leq \frac{4m^{-2/1+2s}}{(1 - 2^{-2/1+2s})} \|f\|_{W_{\frac{2}{1+2s}}}^2 t_{\epsilon}^{4s/1+2s} + B_3. \end{split}$$
$$B_3 &:= \sum_{j < j_{\epsilon}, k} \mathbb{E} (1 - \varpi_{jk}(\epsilon))^2 \beta_{jk}^2 I\{|y_{jk}| > mt_{\epsilon}\} I\{|\beta_{jk}| > mt_{\epsilon}/2\} [I\{|y_{jk}| \ge |\beta_{jk}|/2\} + I\{|y_{jk}| < |\beta_{jk}|/2\}] \\ &:= B_3' + B_3''_3. \end{split}$$

$$B''_{3} \leq \sum_{j < j_{\epsilon}, k} \beta_{jk}^{2} \mathbb{P}(|y_{jk} - \beta_{jk}| \ge mt_{\varepsilon}/4)$$
$$\leq ||f||_{2}^{2} t_{\epsilon}^{2},$$

since $m^2 \ge 64$. We have used in the line above the concentration property of the Gaussian distribution. Now using (24) and (11), we get,

$$B'_{3} \leq \sum_{j < j_{\epsilon},k} \mathbb{E}\beta_{jk}^{2} (1 - \varpi_{jk}(\epsilon))^{2} I\{|y_{jk}| \geq |\beta_{jk}|/2\} I\{|\beta_{jk}| > mt_{\epsilon}/2\} I\{|y_{jk}| \geq mt_{\epsilon}\}]$$

$$\leq \sum_{j < j_{\epsilon},k} \mathbb{E}\beta_{jk}^{2} K(\varpi)^{2} \left(\frac{t_{\varepsilon}}{|y_{jk}|} + t_{\epsilon}\right)^{2} I\{|y_{jk}| \geq |\beta_{jk}|/2\} I\{|\beta_{jk}| > mt_{\epsilon}/2\})$$

$$\leq K(\varpi)^{2} \frac{32m^{-2/1+2s}}{1 - 2^{-2/1+2s}} \|f\|_{W_{\frac{2}{1+2s}}}^{2} t_{\epsilon}^{4s/1+2s} + 2K(\varpi)^{2} \|f\|_{2}^{2} t_{\epsilon}^{2}.$$

Proof of Theorem 3: We shall prove that under our assumption the LVGP rules verify Assumptions (22), (23) and (24). First assumption is obviously checked. Note that we already remarked in subsection (2.1) that they are shrinkage rules. Now, let us fix $m \ge 8$ and let us assume that $|y_{jk}| \le mt_{\epsilon}$. Then,

$$\eta_{jk} = \frac{1}{w(\epsilon)} \frac{\sqrt{\epsilon^2 + \tau(\epsilon)^2}}{\epsilon} \exp\left(-\frac{\tau(\epsilon)^2 y_{jk}^2}{2\epsilon^2(\epsilon^2 + \tau(\epsilon)^2)}\right)$$
$$\geq \frac{1}{w(\epsilon)} t_{\epsilon}^{-1/2} \exp\left(-\frac{m^2 t_{\epsilon}^2}{2\epsilon^2}\right)$$
$$\geq \epsilon^{\frac{m^2}{2} - \frac{1}{2}} \frac{1}{w(\epsilon)} (\log(1/\epsilon))^{-1/4}.$$

- If $q_2 > \frac{m^2 1}{2}$, for ϵ small enough, $\eta_{jk} \ge 1$ and $\check{\beta}_{jk} = 0$.
- If $q_2 \ge \frac{m^2+1}{2}$, for ϵ small enough, $\eta_{jk} \ge t_{\epsilon}^{-1}$ and $\frac{b_j}{1+\eta_{jk}} \le t_{\epsilon}$.

So, Assumption (23) is checked for both rules. Now, let us prove Assumption (24). Let us fix a constant $M \ge \sqrt{6+4q_1}$. We assume $|y_{jk}| > Mt_{\epsilon}$. Then, for ϵ small enough,

$$\eta_{jk} = \frac{1}{w(\epsilon)} \frac{\sqrt{\epsilon^2 + \tau(\epsilon)^2}}{\epsilon} \exp\left(-\frac{\tau(\epsilon)^2 y_{jk}^2}{2\epsilon^2(\epsilon^2 + \tau(\epsilon)^2)}\right)$$

$$\leq \frac{1}{w(\epsilon)} \frac{\sqrt{\epsilon^2 + \tau(\epsilon)^2}}{\epsilon} \epsilon^{\frac{M^2}{4}}$$

$$\leq \frac{1}{w(\epsilon)} t_{\epsilon}^{-1/2} \epsilon^{\frac{M^2}{4}}$$

$$\leq t_{\epsilon}.$$

Let us first consider the posterior median. Using the previous inequality, we have for ϵ small enough, and for any $j < j_{\epsilon}$ and any k,

$$\epsilon \sqrt{b_j} \Phi^{-1}\left(\frac{1+\min(\eta_{jk},1)}{2}\right) \le t_{\epsilon}.$$

So,

$$\begin{aligned} |y_{jk} - \breve{\beta_{jk}}| &= |y_{jk} - \breve{\beta_{jk}}|I\{|y_{jk}| > Mt_{\epsilon}\} + |y_{jk} - \breve{\beta_{jk}}|I\{|y_{jk}| \le Mt_{\epsilon}\} \\ &\leq ((1 - b_j)|y_{jk}| + t_{\epsilon})I\{|y_{jk}| > Mt_{\epsilon}\} + 2|y_{jk}|I\{|y_{jk}| \le Mt_{\epsilon}\} \\ &\leq t_{\epsilon}|y_{jk}| + (1 + 2M)t_{\epsilon}, \end{aligned}$$

which implies (24) for the posterior median. Now, let us deal with the posterior mean. For ϵ

small enough, and for any $j < j_{\epsilon}$ and any k,

$$\begin{aligned} |y_{jk} - \tilde{\beta_{jk}}| &= |y_{jk} - \tilde{\beta_{jk}}|I\{|y_{jk}| > Mt_{\epsilon}\} + |y_{jk} - \tilde{\beta_{jk}}|I\{|y_{jk}| \le Mt_{\epsilon}\} \\ &\leq \left(1 - \frac{b_j}{1 + \eta_{jk}}\right)|y_{jk}|I\{|y_{jk}| > Mt_{\epsilon}\} + 2|y_{jk}|I\{|y_{jk}| \le Mt_{\epsilon}\} \\ &\leq (1 - b_j + \eta_{jk})|y_{jk}|I\{|y_{jk}| > Mt_{\epsilon}\} + 2|y_{jk}|I\{|y_{jk}| \le Mt_{\epsilon}\} \\ &\leq 2t_{\epsilon}|y_{jk}| + 2Mt_{\epsilon}, \end{aligned}$$

which implies (24) for the posterior mean.

Assumptions (22),(23) and (24) are checked for both rules, which finally proves that their maxiset contains $\mathcal{B}_{2,\infty}^{\frac{s}{1+2s}} \cap W_{\frac{2}{1+2s}}$ for the rate $t_{\epsilon}^{4s/(1+2s)} = (\epsilon \sqrt{\log(1/\epsilon)})^{4s/(1+2s)}$. We prove now the reverse inclusion:

$$MS(f_{\epsilon}^{0}, \|.\|_{2}^{2}, (\epsilon\sqrt{\log(1/\epsilon)})^{4s/(1+2s)}) \subset \mathcal{B}_{2,\infty}^{s/(2s+1)} \cap W_{\frac{2}{2s+1}}.$$

Observe that $\beta_{jk}^0 = 0$ when $j \ge j_{\epsilon}$, which implies,

$$\sum_{j>j_{\epsilon},k} \beta_{jk}^2 \le \mathbb{E} \|f_{\epsilon}^0 - f\|_2^2 \le ct_{\epsilon}^{\frac{4s}{1+2s}} \le c2^{-j_{\epsilon}\frac{2s}{1+2s}}.$$

Letting ϵ vary, we obtain the characterization (10), which proves that:

$$MS(f_{\epsilon}^{0}, \|.\|_{2}^{2}, (\epsilon \sqrt{\log(1/\epsilon)})^{4s/(1+2s)}) \subset \mathcal{B}_{2,\infty}^{s/(2s+1)}$$

If we remember that if $|y_{jk}| \leq mt_{\epsilon}$ then $0 \leq \beta_{jk}^0/y_{jk} \leq ct_{\epsilon}$ (Assumption (23)), we have for $f \in MS(f_{\epsilon}^0, \|.\|_2^2, (\epsilon \sqrt{\log(1/\epsilon)})^{4s/(1+2s)})(M)$:

$$\begin{aligned} (1 - ct_{\epsilon})^{2} \sum_{j,k} \beta_{jk}^{2} I\{|\beta_{jk}| \leq mt_{\epsilon}\} \\ &= 2(1 - ct_{\epsilon})^{2} \sum_{j,k} \beta_{jk}^{2} \left[\mathbb{P}(y_{jk} - \beta_{jk} < 0) I\{\beta_{jk} \geq 0\} + \mathbb{P}(y_{jk} - \beta_{jk} > 0) I\{\beta_{jk} < 0\} \right] I\{|\beta_{jk}| \leq mt_{\epsilon}\} \\ &\leq 2\mathbb{E} \sum_{j,k} \left[(\beta_{jk} - \beta_{jk}^{0})^{2} I\{\beta_{jk} \geq 0\} + (\beta_{jk} - \beta_{jk}^{0})^{2} I\{\beta_{jk} < 0\} \right] I\{|\beta_{jk}| \leq mt_{\epsilon}\} \\ &\leq 2\mathbb{E} \sum_{j,k} (\beta_{jk} - \beta_{jk}^{0})^{2} \\ &\leq 2M \left(\epsilon \sqrt{\log(1/\epsilon)} \right)^{4s/(1+2s)}. \end{aligned}$$

We deduce that

$$\sup_{\lambda>0} \lambda^{-\frac{4s}{2s+1}} \sum_{j\geq -1} \sum_{k} \beta_{jk}^2 I\{|\beta_{jk}| \leq \lambda\} < \infty,$$

and f belongs to $W_{\frac{2}{2s+1}}$.

6 Simulations

Dealing with the prior model (1), we compare in this section the performances of both LVGP rules described in the previous section, in (19) and (20), with many other procedures: the thresholding rules of [Donoho and Johnstone, 1994] called VisuShrink and of [Nason, 1996] called GlobalSure, the ParetoThresh (with p=1.3) proposed by [Rivoirard, 2004b] built using Pareto priors and hyperparameters as well as the Bayesian procedures of [Abramovich et al., 1998] denoted as BayesThresh and those proposed by [Johnstone and Silverman, 2004b] and implemented by [Antoniadis et al., 2000] built with the heavy-tailed Laplace prior with scale factor $\alpha = 0.5$ (LaplaceBayesMedian, LaplaceBayesMean) and with the heavy-tailed quasi-Cauchy prior(CauchyBayesMedian, CauchyBayesMean). For this purpose, we use the mean-squared error in the following regression model.

6.1 Model and discrete wavelet transform

Let us consider the standard regression problem:

$$g_i = f(\frac{i}{n}) + \sigma \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, 1), \quad 1 \le i \le n,$$
 (25)

where n = 1024. We introduce the discrete wavelet transform (denoted DWT) of the vector $f_0 = (f(\frac{i}{n}), \quad 1 \le i \le n)^T$:

$$d := \mathcal{W}f^0.$$

The DWT matrix \mathcal{W} is orthogonal. Therefore, we can reconstruct f_0 by the relation

$$f_0 = \mathcal{W}^T d.$$

These transformations performed by Mallat's fast algorithm require only O(n) operations [Mallat, 1998]. The DWT provides n discrete wavelet coefficients d_{jk} , $-1 \leq j \leq N-1$, $k \in \mathcal{I}_j$. They are related to the wavelet coefficients β_{jk} of f by the simple relation

$$d_{jk} \approx \beta_{jk} \times \sqrt{n}.$$

Using the DWT, the regression model (25) is reduced to the following one:

$$y_{jk} = d_{jk} + \sigma z_{jk}, \quad -1 \le j \le N - 1, \quad k \in \mathcal{I}_j,$$

where

$$y := (y_{jk})_{j,k} = \mathcal{W}g$$

and

$$z := (z_{jk})_{j,k} = \mathcal{W}\epsilon.$$

Since \mathcal{W} is orthogonal, z is a vector of independent $\mathcal{N}(0,1)$ variables. Now, instead of estimating f, we estimate the d_{jk} 's.

In the sequel, we suppose that σ is known. Nevertheless, it could robustly be estimated by the median absolute deviation of the $(d_{N-1,k})_{k \in \mathcal{I}_{N-1}}$ divided by 0.6745 (see [Donoho and Johnstone, 1994]).

To implement the LVGP rules, we reconstruct the d_{jk} 's, as posterior median and the posterior mean of a prior having the following form:

$$d_{jk} \sim \frac{\omega_n}{1+\omega_n} \gamma_{j,n} + \frac{1}{1+\omega_n} \,\delta(0),$$

where $\omega_n = \omega^* = 10(\frac{\sigma}{\sqrt{n}})^q$ $(q > 0), \delta(0)$ is a point mass at zero, γ is the Gaussian density and

$$\gamma_{j,n}(d_{jk}) = \frac{1}{\tau_n} \gamma(\frac{d_{jk}}{\tau_n}),$$

with τ_n is such that $\frac{n\tau_n^2}{\sigma^2 + n\tau_n^2} = 0,999.$

Dealing with this prior model, we respectively denote *GaussMedian* and *GaussMean*, the LVGP rules described in (19) and (20).

The Symmlet 8 wavelet basis (as described on page 198 of [Daubechies, 1992]) is used for all the methods of reconstruction. In Table 1 we measure the performances of both estimators by using the four test functions: "Blocks", "Bumps", "Heavisine" and "Doppler" thanks to the mean-squared error defined by:

MSE
$$(\hat{f}) = \frac{1}{n} \sum_{i=1}^{n} \left(\hat{f}(\frac{i}{n}) - f(\frac{i}{n}) \right)^2$$
.

Remark: Recall that the test functions have been chosen by [Donoho et al., 1995] to represent a large variety of inhomogeneous signals.

6.2 Simulations and discussion

Table 1 shows the average mean-squared error (denoted AMSE) using 100 replications for VisuShrink, GlobalSure, ParetoThresh, BayesThresh, GaussMedian, GaussMean (for q = 1), LaplaceBayesMedian, LaplaceBayesMean, CauchyBayesMedian and CauchyBayesMean, with different values for the root signal to noise ration (RSNR).

The results provided below can be summarized as follows:

• According to Table 1, we remark that "purely Bayesian" procedures (BayesThresh, Gauss-Median, GaussMean, CauchyBayesMedian, CauchyBayesMean, LaplaceBayesMedian and LaplaceBayesMean) are preferable to "purely deterministic" ones (VisuShrink and GlobalSure) under the AMSE approach for inhomogeneous signals. Paretothresh appears as a good compromise between purely Bayesian rules and deterministic ones.

- We observe that Bayesian rules using the posterior mean (GaussMean, LaplaceBayesMean and CauchyBayesMean) have better performances than those using the posterior median (GaussMedian, LaplaceBayesMedian and CauchyBayesMedian).
- CauchyBayesMean provides the best behaviors here since its AMSEs are globally the smallest (11 times on 12).
- GaussMean shows performances which are rather close to CauchyBayesMean. It outperforms BayesThresh 11 times on 12. This confirms our maxiset previous results, and shows that GaussMean is an excellent choice if we take into account the performances as well as the computation time.

RSNR=5	Blocks	Bumps	Heavisine	Doppler
VisuShrink	2.08	2.99	0.17	0.77
GlobalSure	0.82	0.92	0.18	0.59
ParetoThresh	0.73	0.85	0.15	0.36
BayesThresh	0.67	0.74	0.15	0.30
GaussMedian	0.72	0.76	0.20	0.30
GaussMean	0.62	0.68	0.19	0.29
LaplaceBayesMedian	0.59	0.69	0.14	0.30
LaplaceBayesMean	0.56	0.65	0.13	0.28
CauchyBayesMedian	0.60	0.67	0.14	0.29
CauchyBayesMean	0.55	0.63	0.13	0.27
RSNR=7	Blocks	Bumps	Heavisine	Doppler
VisuShrink	1.29	1.77	0.12	0.47
GlobalSure	0.42	0.48	0.12	0.21
ParetoThresh	0.40	0.46	0.09	0.21
BayesThresh	0.38	0.45	0.10	0.16
GaussMedian	0.41	0.42	0.12	0.15
GaussMean	0.35	0.38	0.11	0.15
LaplaceBayesMedian	0.33	0.37	0.09	0.17
LaplaceBayesMean	0.31	0.36	0.08	0.16
CauchyBayesMedian	0.32	0.36	0.09	0.17
CauchyBayesMean	0.29	0.34	0.08	0.15
RSNR=10	Blocks	Bumps	Heavisine	Doppler
VisuShrink	0.77	1.04	0.08	0.27
GlobalSure	0.25	0.29	0.08	0.11
ParetoThresh	0.21	0.25	0.06	0.12
BayesThresh	0.22	0.25	0.06	0.09
GaussMedian	0.21	0.23	0.06	0.08
GaussMean	0.18	0.20	0.06	0.07
LaplaceBayesMedian	0.17	0.20	0.05	0.09
LaplaceBayesMean	0.17	0.19	0.05	0.09
CauchyBayesMedian	0.17	0.19	0.05	0.09
CauchyBayesMean	0.16	0.18	0.05	0.09

Table 1: AMSEs pour VisuShrink, GlobalSure, ParetoThresh, BayesThresh, GaussMedian, GaussMean, LaplaceBayesMedian, LaplaceBayesMean, CauchyBayesMedian and Cauchy-BayesMean, with various test functions and various values of the RSNR.

In the sequel, we present some simulations of the Bayesian rules using Gaussian priors (1 and 2) and heavy-tailed priors (3 and 4) when RSNR=5.



Figure 1: Original test functions and reconstructions using GaussMedian and GaussMean with q = 1 (RSNR=5).

In Figure 1, we note that in both Bayesian procedures some high-frequency artefacts appear. However, these artefacts disappear if we take large values of q. Figure 2 shows an example of reconstructions using GaussMedian and GaussMean when the RSNR is equal to 5 ($\sigma = 7/5$) for different values of q.



Figure 2: Reconstructions with GaussMedian (schemes a,b et c) and GaussMean (schemes d,e et f) for various values of q when RSNR=5; a: AMSE=0.37. b: AMSE=0.30. c: AMSE=0.33. d: AMSE=0.39. e: AMSE=0.29. f: AMSE=0.30.

As we can see in Figure 2, the artefacts are less numerous when q increases. But this improvement has a cost: in general the AMSE increases when q is close to 0 or strictly greater than 1. Consequently, the value q = 1 appears as a good compromise to obtain good reconstructions and good AMSE with the GaussMedian and GaussMean procedures.



Figure 3: Original test functions and reconstructions using LaplaceBayesMedian and LaplaceBayesMean (RSNR=5).



Figure 4: Original test functions and reconstructions using CauchyBayesMedian and Cauchy-BayesMean (RSNR=5).

7 Appendix: More on maxisets of 'small variance Gaussian priors'

In a minimax setting, [Abramovich et al., 2004] obtained the following result:

Theorem 4. Let β^0 be $\check{\beta}$ or $\tilde{\beta}$. With $\alpha = 2s + 1$ and any $0 \le b < 1$, there exist two positive constants C_1 and C_2 such that $\forall \epsilon > 0$,

$$C_1(\epsilon \sqrt{\log(1/\epsilon)})^{4s/(2s+1)} \le \sup_{\beta \in \mathcal{B}^s_{2,\infty}(M)} \mathbb{E} \|\beta^0 - \beta\|_2^2 \le C_2 \log(1/\epsilon) \epsilon^{4s/(2s+1)}.$$

So, posterior mean and median achieve the optimal rate up to an unavoidable logarithmic term. Now, let us consider the maxiset setting.

Theorem 5. For s > 0, $\alpha = 2s + 1$, any $0 \le b < 1$, and if β^0 is $\check{\beta}$ or $\tilde{\beta}$,

1. for the rate $e^{4s/(1+2s)}$,

$$MS(\beta^0, \|.\|_2^2, \epsilon^{4s/(1+2s)}) \subsetneq \mathcal{B}_{2,\infty}^s$$

2. For the rate $(\epsilon \sqrt{\log(1/\epsilon)})^{4s/(1+2s)}$,

$$MS(\beta^0, \|.\|_2^2, (\epsilon \sqrt{\log(1/\epsilon)})^{4s/(1+2s)}) \subset \mathcal{B}_{2,\infty}^{*s},$$

with

$$\mathcal{B}_{2,\infty}^{*s} = \left\{ f \in \mathbb{L}^2 : \quad \sup_{J>0} 2^{2Js} J^{-2s/(1+2s)} \sum_{j \ge J} \sum_k \beta_{jk}^2 < \infty \right\}.$$

3. For the rate $\epsilon^{4s/(1+2s)}\log(1/\epsilon)$,

$$\mathcal{B}_{2,\infty}^{s} \subset MS(\beta^{0}, \|.\|_{2}^{2}, \epsilon^{4s/(1+2s)}\log(1/\epsilon)).$$

Proof: Let us first prove the inclusion

$$MS(\beta^0, \|.\|_2^2, \epsilon^{4s/(1+2s)}) \subset \mathcal{B}_{2,\infty}^s.$$

For this, let us note $\lambda_{\epsilon} = (c_1^{-1}\epsilon^2)^{1/\alpha}$. We observe that if $2^{-j} \leq \lambda_{\epsilon}$ then $b_j \leq 1/2$ and

$$|\beta_{jk}^0| \le \frac{1}{2} |y_{jk}|.$$

Since $y_{jk} \times \beta_{jk}^0 \ge 0$, if $2^{-j} \le \lambda_{\epsilon}$,

$$\mathbb{E}\beta_{jk}^2 I\{\beta_{jk} \ge 0\} I\{y_{jk} < \beta_{jk}\} \le 4\mathbb{E}(\beta_{jk}^0 - \beta_{jk})^2 I\{\beta_{jk} \ge 0\} I\{y_{jk} < \beta_{jk}\},\$$

and

$$\mathbb{E}\beta_{jk}^2 I\{\beta_{jk} < 0\} I\{y_{jk} > \beta_{jk}\} \le 4\mathbb{E}(\beta_{jk}^0 - \beta_{jk})^2 I\{\beta_{jk} < 0\} I\{y_{jk} > \beta_{jk}\}.$$

So, since $\mathbb{P}(y_{jk} - \beta_{jk} < 0) = \mathbb{P}(y_{jk} - \beta_{jk} > 0) = 1/2$, if $f \in MS(\beta^0, \|.\|_2^2, \epsilon^{4s/(1+2s)})(M)$, we have:

$$\begin{split} &\sum_{j,k} \beta_{jk}^2 I\{2^{-j} \le \lambda_{\epsilon}\} \\ &= 2 \sum_{j,k} \beta_{jk}^2 \left[\mathbb{P}(y_{jk} - \beta_{jk} < 0) I\{\beta_{jk} \ge 0\} + \mathbb{P}(y_{jk} - \beta_{jk} > 0) I\{\beta_{jk} < 0\} \right] I\{2^{-j} \le \lambda_{\epsilon}\} \\ &\le 8 \mathbb{E} \sum_{j,k} \left[(\beta_{jk}^0 - \beta_{jk})^2 I\{\beta_{jk} \ge 0\} + (\beta_{jk}^0 - \beta_{jk})^2 I\{\beta_{jk} < 0\} \right] I\{2^{-j} \le \lambda_{\epsilon}\} \\ &\le 8 \mathbb{E} \sum_{j,k} (\beta_{jk}^0 - \beta_{jk})^2 \\ &\le 8 M \ \epsilon^{4s/(1+2s)}. \end{split}$$

Since $\alpha = 2s + 1$, we deduce

$$\sup_{J \ge -1} 2^{2Js} \sum_{j \ge J} \sum_{k} \beta_{jk}^2 \le 8M c_1^{2s/(1+2s)}$$

and f belongs to $\mathcal{B}_{2,\infty}^s$. To prove that the inclusion is strict, we just use Theorem 4. The second inclusion is easily obtained by using similar arguments. Finally, the proof of the last one is provided by Theorem 4.

As recalled in Section 3.2, for the rates $e^{4s/(1+2s)}$, the maxisets of linear estimates are exactly Besov spaces $\mathcal{B}^s_{2,\infty}$. So Theorem 5 shows that the Bayesian procedures built by [Abramovich et al., 2004] are outperformed by linear estimates for polynomial rates of convergence. Furthermore, these procedures cannot achieve the same performances as classical non linear procedures, since we have the following result.

Proposition 4. For any s > 0,

$$\mathcal{B}_{2,\infty}^{s/(2s+1)} \cap W_{\frac{2}{2s+1}} \not\subset \mathcal{B}_{2,\infty}^{*s}.$$

Proof: To prove this result, we build a *sparse* function belonging to $\mathcal{B}_{2,\infty}^{s/(2s+1)} \cap W_{\frac{2}{2s+1}}$ but not to $\mathcal{B}_{2,\infty}^{*s}$. Let us consider $f = \sum_{j,k} \beta_{jk} \psi_{jk}$, where at each level j, 2^{jn} wavelet coefficients take the value $2^{-j\beta}$, whereas the other ones are equal to 0, with $0 \le n \le 1$ and $\beta > n/2$ (so $f \in \mathbb{L}_2$). For any $J \ge 1$,

$$\begin{split} 2^{2Js}J^{-2s/(1+2s)}\sum_{j\geq J}\sum_k \beta_{jk}^2 &= 2^{2Js}J^{-2s/(1+2s)}\sum_{j\geq J} 2^{nj}2^{-2j\beta} \\ &\geq 2^{J(2s+n-2\beta)}J^{-2s/(1+2s)}. \end{split}$$

So,

$$n - 2\beta + 2s > 0 \Rightarrow f \notin \mathcal{B}_{2,\infty}^{*s}.$$
(26)

Similarly,

$$n - 2\beta + 2s/(1+2s) \le 0 \Rightarrow f \in \mathcal{B}_{2,\infty}^{s/(2s+1)}.$$
 (27)

And

$$\begin{split} \lambda^{-4s/(1+2s)} \sum_{j,k} \beta_{jk}^2 I\{|\beta_{jk}| \leq \lambda\} &= \lambda^{-4s/(1+2s)} \sum_j 2^{jn} 2^{-2j\beta} I\{2^{-j\beta} \leq \lambda\} \\ &\leq \lambda^{-4s/(1+2s)-n/\beta+2}. \end{split}$$

So,

$$n - 2\beta + 2ns \le 0 \Rightarrow f \in W_{\frac{2}{2s+1}}.$$
(28)

As soon as n < 1 (that yields that the signal is sparse), it is then possible to choose $\beta > n/2$ such that (26), (27) and (28) are checked. So, f belongs to $\mathcal{B}_{2,\infty}^{s/(2s+1)} \cap W_{\frac{2}{2s+1}}$ but not to $\mathcal{B}_{2,\infty}^{*s}$.

References

- [Abramovich et al., 2004] Abramovich, F., Amato, U., and Angelini, C. (2004). On optimality of Bayesian wavelet estimators. *Scand. J. Statist.*, 31(2):217–234.
- [Abramovich et al., 1998] Abramovich, F., Sapatinas, T., and Silverman, B. W. (1998). Wavelet thresholding via a Bayesian approach. J. R. Stat. Soc. Ser. B Stat. Methodol., 60(4):725–749.
- [Antoniadis et al., 2001] Antoniadis, A., Bigot, J., and Sapatinas, T. (2001). Wavelet estimators in nonparametric regression: a comparative simulation study. *Journal of Statistical Software*, 6(6):1–83.
- [Antoniadis et al., 2000] Antoniadis, A., Jansen, M., Johnstone, I., and Silverman, B. (2000). Ebayesthresh: Matlab software for empirical bayes thresholding. http://wwwlmc.imag.fr/lmc-sms/Anestis.Antoniadis/EBayesThresh/.
- [Chipman et al., 1997] Chipman, H. A., Kolaczyk, E. D., and McCulloch, R. E. (1997). Adaptive bayesian wavelet shrinkage. Journal of the American Statistical Association, 92:1413– 1421.
- [Clyde and George, 1998] Clyde, M. and George, E. I. (1998). Robust empirical bayes estimation in wavelets. *Technical Report*.
- [Clyde and George, 2000] Clyde, M. and George, E. I. (2000). Flexible empirical Bayes estimation for wavelets. J. R. Stat. Soc. Ser. B Stat. Methodol., 62(4):681–698.
- [Clyde et al., 1998] Clyde, M., Parmigiani, G., and Vidakovic, B. (1998). Multiple shrinkage and subset selection in wavelets. *Biometrika*, 85(2):391–401.
- [Cohen et al., 2001] Cohen, A., DeVore, R., Kerkyacharian, G., and Picard, D. (2001). Maximal spaces with given rate of convergence for thresholding algorithms. *Appl. Comput. Harmon. Anal.*, 11(2):167–191.
- [Daubechies, 1992] Daubechies, I. (1992). *Ten lectures on wavelets*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia.
- [Donoho and Johnstone, 1994] Donoho, D. L. and Johnstone, I. M. (1994). Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81(3):425–455.
- [Donoho et al., 1995] Donoho, D. L., Johnstone, I. M., Kerkyacharian, G., and Picard, D. (1995). Wavelet shrinkage: Asymptopia? Journal of the Royal Statistical Society, Series B, 57:301–369. With Discussion.
- [Johnstone, 1994] Johnstone, I. M. (1994). Minimax Bayes, asymptotic minimax and sparse wavelet priors. In Statistical decision theory and related topics, V (West Lafayette, IN, 1992), pages 303–326. Springer, New York.
- [Johnstone and Silverman, 1998] Johnstone, I. M. and Silverman, B. W. (1998). Empirical bayes approaches to mixture problems and wavelet regression. *Technical Report*.

- [Johnstone and Silverman, 2004a] Johnstone, I. M. and Silverman, B. W. (2004a). Empirical bayes selection of wavelet thresholds. *Technical Report*.
- [Johnstone and Silverman, 2004b] Johnstone, I. M. and Silverman, B. W. (2004b). Needles and hay in haystacks: Empirical bayes estimates of possibly sparse sequences. *Ann. Statist*, 32:1594–1649.
- [Kerkyacharian and Picard, 1993] Kerkyacharian, G. and Picard, D. (1993). Density estimation by kernel and wavelets methods: optimality of Besov spaces. *Statist. Probab. Lett.*, 18(4):327– 336.
- [Kerkyacharian and Picard, 2000] Kerkyacharian, G. and Picard, D. (2000). Thresholding algorithms, maxisets and well-concentrated bases. *Test*, 9(2):283–344.
- [Kerkyacharian and Picard, 2002] Kerkyacharian, G. and Picard, D. (2002). Minimax or maxisets? *Bernoulli*, 8(2):219–253.
- [Mallat, 1998] Mallat, S. (1998). A wavelet tour of signal processing. Academic Press Inc., San Diego.
- [Meyer, 1990] Meyer, Y. (1990). Ondelettes et opérateurs. I. Actualités Mathématiques. [Current Mathematical Topics]. Hermann, Paris.
- [Nason, 1996] Nason, G. P. (1996). Wavelet shrinkage using cross-validation. J. Roy. Statist. Syst. Sci. B, 23(6):1–11.
- [Rivoirard, 2003] Rivoirard, V. (2003). Bayesian modelization of sparse sequences and maxisets for bayes rules. *Technical Report. Submitted to Math. Methods Statist.*
- [Rivoirard, 2004a] Rivoirard, V. (2004a). Maxisets for linear procedures. *Statist. Probab. Lett.*, 67(3):267–275.
- [Rivoirard, 2004b] Rivoirard, V. (2004b). Thresholding procedures using priors based on pareto distributions. *Test*, 13(1):213–246.
- [Vidakovic, 1998] Vidakovic, B. (1998). Nonlinear wavelet shrinkage with Bayes rules and Bayes factors. J. Amer. Statist. Assoc., 93(441):173–179.
- [Zhang, 2002] Zhang, C.-H. (2002). General empirical bayes wavelet methods. Technical Report.