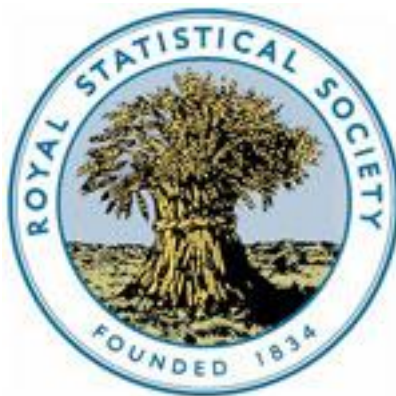


WILEY



Sparse Additive Models

Author(s): Pradeep Ravikumar, John Lafferty, Han Liu and Larry Wasserman

Source: *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, Vol. 71, No. 5 (Nov., 2009), pp. 1009-1030

Published by: Wiley for the Royal Statistical Society

Stable URL: <http://www.jstor.org/stable/40541567>

Accessed: 09-05-2016 15:46 UTC

REFERENCES

Linked references are available on JSTOR for this article:

http://www.jstor.org/stable/40541567?seq=1&cid=pdf-reference#references_tab_contents

You may need to log in to JSTOR to access the linked references.

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at

<http://about.jstor.org/terms>

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.



Wiley, Royal Statistical Society are collaborating with JSTOR to digitize, preserve and extend access to *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*

Sparse additive models

Pradeep Ravikumar,

University of California, Berkeley, USA

and John Lafferty, Han Liu and Larry Wasserman

Carnegie Mellon University, Pittsburgh, USA

[Received April 2008. Final revision March 2009]

Summary. We present a new class of methods for high dimensional non-parametric regression and classification called sparse additive models. Our methods combine ideas from sparse linear modelling and additive non-parametric regression. We derive an algorithm for fitting the models that is practical and effective even when the number of covariates is larger than the sample size. Sparse additive models are essentially a functional version of the grouped lasso of Yuan and Lin. They are also closely related to the COSSO model of Lin and Zhang but decouple smoothing and sparsity, enabling the use of arbitrary non-parametric smoothers. We give an analysis of the theoretical properties of sparse additive models and present empirical results on synthetic and real data, showing that they can be effective in fitting sparse non-parametric models in high dimensional data.

Keywords: Additive models; Lasso; Non-parametric regression; Sparsity

1. Introduction

Substantial progress has been made recently on the problem of fitting high dimensional linear regression models of the form $Y_i = X_i^T \beta + \varepsilon_i$, for $i = 1, \dots, n$. Here Y_i is a real-valued response, X_i is a predictor and ε_i is a mean 0 error term. Finding an estimate of β when $p > n$ that is both statistically well behaved and computationally efficient has proved challenging; however, under the assumption that the vector β is sparse, the lasso estimator (Tibshirani, 1996) has been remarkably successful. The lasso estimator $\hat{\beta}$ minimizes the l_1 -penalized sum of squares

$$\sum_i (Y_i - X_i^T \beta)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

with the l_1 -penalty $\|\beta\|_1$ encouraging sparse solutions, where many components $\hat{\beta}_j$ are 0. The good empirical success of this estimator has been recently backed up by results confirming that it has strong theoretical properties; see Bunea *et al.* (2007), Greenshtein and Ritov (2004), Zhao and Yu (2007), Meinshausen and Yu (2006) and Wainwright (2006).

The non-parametric regression model $Y_i = m(X_i) + \varepsilon_i$, where m is a general smooth function, relaxes the strong assumptions that are made by a linear model but is much more challenging in high dimensions. Hastie and Tibshirani (1999) introduced the class of additive models of the form

Address for correspondence: Larry Wasserman, Department of Statistics, 232 Baker Hall, Carnegie Mellon University, Pittsburgh, PA 15213-3890, USA.
E-mail: larry@stat.cmu.edu

$$Y_i = \sum_{j=1}^p f_j(X_{ij}) + \varepsilon_i. \quad (1)$$

This additive combination of univariate functions—one for each covariate X_j —is less general than joint multivariate non-parametric models but can be more interpretable and easier to fit; in particular, an additive model can be estimated by using a co-ordinate descent Gauss–Seidel procedure, called backfitting. Unfortunately, additive models only have good statistical and computational behaviour when the number of variables p is not large relative to the sample size n , so their usefulness is limited in the high dimensional setting.

In this paper we investigate sparse additive models (SPAMs), which extend the advantages of sparse linear models to the additive non-parametric setting. The underlying model is the same as in equation (1), but we impose a sparsity constraint on the index set $\{j: f_j \neq 0\}$ of functions f_j that are not identically zero. Lin and Zhang (2006) have proposed COSSO, an extension of the lasso to this setting, for the case where the component functions f_j belong to a reproducing kernel Hilbert space. They penalized the sum of the reproducing kernel Hilbert space norms of the component functions. Yuan (2007) proposed an extension of the non-negative garrotte to this setting. As with the parametric non-negative garrotte, the success of this method depends on the initial estimates of component functions f_j .

In Section 3, we formulate an optimization problem in the population setting that induces sparsity. Then we derive a sample version of the solution. The SPAM estimation procedure that we introduce allows the use of arbitrary non-parametric smoothing techniques, effectively resulting in a combination of the lasso and backfitting. The algorithm extends to classification problems by using generalized additive models. As we explain later, SPAMs can also be thought of as a functional version of the grouped lasso (Antoniadis and Fan, 2001; Yuan and Lin, 2006).

The main results of this paper include the formulation of a convex optimization problem for estimating a SPAM, an efficient backfitting algorithm for constructing the estimator and theoretical results that analyse the effectiveness of the estimator in the high dimensional setting. Our theoretical results are of two different types. First, we show that, under suitable choices of the design parameters, the SPAM backfitting algorithm recovers the correct sparsity pattern asymptotically; this is a property that we call *sparsistency*, as a shorthand for ‘sparsity pattern consistency’. Second, we show that the estimator is *persistent*, in the sense of Greenshtein and Ritov (2004), which is a form of risk consistency.

In the following section we establish notation and assumptions. In Section 3 we formulate SPAMs as an optimization problem and derive a scalable backfitting algorithm. Examples showing the use of our sparse backfitting estimator on high dimensional data are included in Section 5. In Section 6.1 we formulate the sparsistency result, when orthogonal function regression is used for smoothing. In Section 6.2 we give the persistence result. Section 7 contains a discussion of the results and possible extensions. Proofs are contained in Appendix A.

The statements of the theorems in this paper were given, without proof, in Ravikumar *et al.* (2008). The backfitting algorithm was also presented there. Related results were obtained in Meier *et al.* (2008) and Koltchinskii and Yuan (2008).

2. Notation and assumptions

We assume that we are given independent data $(X_1, Y_1), \dots, (X_n, Y_n)$ where $X_i = (X_{i1}, \dots, X_{ij}, \dots, X_{ip})^T \in [0, 1]^p$ and

$$Y_i = m(X_i) + \varepsilon_i \quad (2)$$

with $\varepsilon_i \sim N(0, \sigma^2)$ independent of X_i and

$$m(x) = \sum_{j=1}^p f_j(x_j). \quad (3)$$

Let μ denote the distribution of X , and let μ_j denote the marginal distribution of X_j for each $j = 1, \dots, p$. For a function f_j on $[0, 1]$ denote its $L_2(\mu_j)$ norm by

$$\|f_j\|_{\mu_j} = \sqrt{\int_0^1 f_j^2(x) d\mu_j(x)} = \sqrt{\mathbb{E}\{f_j(X_j)^2\}}. \quad (4)$$

When the variable X_j is clear from the context, we remove the dependence on μ_j in the notation $\|\cdot\|_{\mu_j}$ and simply write $\|f_j\|$.

For $j \in \{1, \dots, p\}$, let \mathcal{H}_j denote the Hilbert subspace $L_2(\mu_j)$ of measurable functions $f_j(x_j)$ of the single scalar variable x_j with zero mean, $\mathbb{E}\{f_j(X_j)\} = 0$. Thus, \mathcal{H}_j has the inner product

$$\langle f_j, f'_j \rangle = \mathbb{E}\{f_j(X_j) f'_j(X_j)\} \quad (5)$$

and $\|f_j\| = \sqrt{\mathbb{E}\{f_j(X_j)^2\}} < \infty$. Let $\mathcal{H} = \mathcal{H}_1 \oplus \mathcal{H}_2 \oplus \dots \oplus \mathcal{H}_p$ denote the Hilbert space of functions of (x_1, \dots, x_p) that have the additive form: $m(x) = \sum_j f_j(x_j)$, with $f_j \in \mathcal{H}_j$, $j = 1, \dots, p$.

Let $\{\psi_{jk}, k = 0, 1, \dots\}$ denote a uniformly bounded, orthonormal basis with respect to $L^2[0, 1]$. Unless stated otherwise, we assume that $f_j \in \mathcal{T}_j$ where

$$\mathcal{T}_j = \left\{ f_j \in \mathcal{H}_j : f_j(x_j) = \sum_{k=0}^{\infty} \beta_{jk} \psi_{jk}(x_j), \sum_{k=0}^{\infty} \beta_{jk}^2 k^{2\nu_j} \leq C^2 \right\} \quad (6)$$

for some $0 < C < \infty$. We shall take $\nu_j = 2$ although the extension to other levels of smoothness is straightforward. It is also possible to adapt to ν_j although we do not pursue that direction here.

Let $\Lambda_{\min}(A)$ and $\Lambda_{\max}(A)$ denote the minimum and maximum eigenvalues of a square matrix A . If $v = (v_1, \dots, v_k)^T$ is a vector, we use the norms

$$\|v\| = \sqrt{\sum_{j=1}^k v_j^2}, \quad \|v\|_1 = \sum_{j=1}^k |v_j|, \quad \|v\|_{\infty} = \max_j |v_j|. \quad (7)$$

3. Sparse backfitting

The outline of the derivation of our algorithm is as follows. We first formulate a population level optimization problem and show that the minimizing functions can be obtained by iterating through a series of soft thresholded univariate conditional expectations. We then plug in smoothed estimates of these univariate conditional expectations, to derive our sparse backfitting algorithm.

3.1. Population sparse additive models

For simplicity, assume that $\mathbb{E}(Y_i) = 0$. The standard additive model optimization problem in $L_2(\mu)$ (the population setting) is

$$\min_{f_j \in \mathcal{H}_j, 1 \leq j \leq p} \left[\mathbb{E} \left\{ Y - \sum_{j=1}^p f_j(X_j) \right\}^2 \right] \quad (8)$$

where the expectation is taken with respect to X and the noise ε . Now consider the following modification of this problem that introduces a scaling parameter for each function, and that imposes additional constraints:

$$\min_{\beta \in \mathbb{R}^p, g_j \in \mathcal{H}_j} \left[\mathbb{E} \left\{ Y - \sum_{j=1}^p \beta_j g_j(X_j) \right\}^2 \right] \quad (9)$$

subject to

$$\sum_{j=1}^p |\beta_j| \leq L, \quad (10)$$

$$\mathbb{E}(g_j^2) = 1, \quad j = 1, \dots, p, \quad (11)$$

noting that g_j is a function whereas $\beta = (\beta_1, \dots, \beta_p)^T$ is a vector. The constraint that β lies in the l_1 -ball $\{\beta : \|\beta\|_1 \leq L\}$ encourages sparsity of the estimated β , just as for the parametric lasso (Tibshirani, 1996). It is convenient to absorb the scaling constants β_j into the functions f_j , and to re-express the minimization in the following equivalent Lagrangian form:

$$\mathcal{L}(f, \lambda) = \frac{1}{2} \mathbb{E} \left\{ Y - \sum_{j=1}^p f_j(X_j) \right\}^2 + \lambda \sum_{j=1}^p \sqrt{\mathbb{E}\{f_j^2(X_j)\}}. \quad (12)$$

Theorem 1. The minimizers $f_j \in \mathcal{H}_j$ of equation (12) satisfy

$$f_j = \left[1 - \frac{\lambda}{\sqrt{\mathbb{E}(P_j^2)}} \right]_+ P_j \quad \text{almost surely} \quad (13)$$

where $[\cdot]_+$ denotes the positive part, and $P_j = \mathbb{E}(R_j | X_j)$ denotes the projection of the residual $R_j = Y - \sum_{k \neq j} f_k(X_k)$ onto \mathcal{H}_j .

An outline of the proof of this theorem appears in Ravikumar *et al.* (2008). A formal proof is given in Appendix A. At the population level, the f_j s can be found by a co-ordinate descent procedure that fixes $(f_k : k \neq j)$ and fits f_j by equation (13), and then iterates over j .

3.2. Data version of sparse additive models

To obtain a sample version of the population solution, we insert sample estimates into the population algorithm, as in standard backfitting (Hastie and Tibshirani, 1999). Thus, we estimate the projection $P_j = \mathbb{E}(R_j | X_j)$ by smoothing the residuals:

$$\hat{P}_j = S_j R_j \quad (14)$$

where S_j is a linear smoother, such as a local linear or kernel smoother. Let

$$\hat{s}_j = \frac{1}{\sqrt{n}} \|\hat{P}_j\| = \sqrt{\text{mean}(\hat{P}_j^2)} \quad (15)$$

be the estimate of $\sqrt{\mathbb{E}(P_j^2)}$. Using these plug-in estimates in the co-ordinate descent procedure yields the SPAM backfitting algorithm that is given in Table 1.

This algorithm can be seen as a functional version of the co-ordinate descent algorithm for solving the lasso. In particular, if we solve the lasso by iteratively minimizing with respect to a single co-ordinate, each iteration is given by soft thresholding; Table 2. Convergence properties of variants of this simple algorithm have been recently treated by Daubechies *et al.* (2004, 2007). Our sparse backfitting algorithm is a direct generalization of this algorithm, and it reduces to it in the case where the smoothers are local linear smoothers with large bandwidths, i.e., as the bandwidth approaches ∞ , the local linear smoother approaches a global linear fit, yielding the estimator $\hat{P}_j(i) = \hat{\beta}_j X_{ij}$. When the variables are standardized,

Table 1. SPAM backfitting algorithm†

Input: data (X_i, Y_i) , regularization parameter λ
Initialize $\hat{f}_j = 0$, for $j = 1, \dots, p$
Iterate until convergence, *for each* $j = 1, \dots, p$

Step 1: compute the residual, $R_j = Y - \sum_{k \neq j} \hat{f}_k(X_k)$
 Step 2: estimate $P_j = \mathbb{E}(R_j | X_j)$ by smoothing, $\hat{P}_j = S_j R_j$
 Step 3: estimate the norm, $\hat{s}_j^2 = (1/n) \sum_{i=1}^n \hat{P}_j^2(i)$
 Step 4: soft threshold, $\hat{f}_j = [1 - \lambda/\hat{s}_j]_+ \hat{P}_j$
 Step 5: centre, $\hat{f}_j \leftarrow \hat{f}_j - \text{mean}(\hat{f}_j)$.

Output: component functions \hat{f}_j and estimator $\hat{m}(X_i) = \sum_j \hat{f}_j(X_{ij})$

†The first two steps in the iterative algorithm are the usual backfitting procedure; the remaining steps carry out functional soft thresholding.

Table 2. Co-ordinate descent lasso†

Input: data (X_i, Y_i) , regularization parameter λ
Initialize $\hat{\beta}_j = 0$, for $j = 1, \dots, p$
Iterate until convergence, *for each* $j = 1, \dots, p$

Step 1: compute the residual, $R_j = Y - \sum_{k \neq j} \hat{\beta}_k X_k$
 Step 2: project residual onto X_j , $P_j = X_j^T R_j$
 Step 3: soft threshold, $\hat{\beta}_j = [1 - \lambda/|P_j|]_+ P_j$

Output: estimator $\hat{m}(X_i) = \sum_j \hat{\beta}_j X_{ij}$

†The SPAM backfitting algorithm is a functional version of the co-ordinate descent algorithm for the lasso, which computes $\hat{\beta} = \arg \min(\frac{1}{2} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1)$.

$$\hat{s}_j = \sqrt{\left(\frac{1}{n} \sum_{i=1}^n \hat{\beta}_j^2 X_{ij}^2 \right)} = |\hat{\beta}_j|$$

so the soft thresholding in step 4 of the SPAM backfitting algorithm is the same as the soft thresholding in step 3 in the co-ordinate descent lasso algorithm.

3.3. Basis functions

It is useful to express the model in terms of basis functions. Recall that $B_j = (\psi_{jk} : k = 1, 2, \dots)$ is an orthonormal basis for \mathcal{T}_j and that $\sup_x |\psi_{jk}(x)| \leq B$ for some B . Then

$$f_j(x_j) = \sum_{k=1}^{\infty} \beta_{jk} \psi_{jk}(x_j) \quad (16)$$

where $\beta_{jk} = \int f_j(x_j) \psi_{jk}(x_j) dx_j$.

Let us also define

$$\tilde{f}_j(x_j) = \sum_{k=1}^d \beta_{jk} \psi_{jk}(x_j) \quad (17)$$

where $d = d_n$ is a truncation parameter. For the Sobolev space \mathcal{T}_j of order 2 we have that $\|f_j - \tilde{f}_j\|^2 = O(1/d^4)$. Let $S = \{j : f_j \neq 0\}$. Assuming the sparsity condition $|S| = O(1)$ it follows that $\|m - \tilde{m}\|^2 = O(1/d^4)$ where $\tilde{m} = \sum_j \tilde{f}_j$. The usual choice is $d \asymp n^{1/5}$, yielding truncation bias $\|m - \tilde{m}\|^2 = O(n^{-4/5})$.

In this setting, the smoother can be taken to be the least squares projection onto the truncated set of basis functions $\{\psi_{j1}, \dots, \psi_{jd}\}$; this is also called orthogonal series smoothing. Let Ψ_j denote the $n \times d_n$ matrix that is given by $\Psi_j(i, l) = \psi_{j,l}(X_{ij})$. The smoothing matrix is the projection matrix $\mathcal{S}_j = \Psi_j(\Psi_j^T \Psi_j)^{-1} \Psi_j^T$. In this case, the backfitting algorithm in Table 1 is a co-ordinate descent algorithm for minimizing

$$\frac{1}{2n} \left\| Y - \sum_{j=1}^p \Psi_j \beta_j \right\|_2^2 + \lambda \sum_{j=1}^p \sqrt{\left(\frac{1}{n} \beta_j^T \Psi_j^T \Psi_j \beta_j \right)}$$

which is the sample version of equation (12). This is the Lagrangian of a second-order cone program, and standard convexity theory implies the existence of a minimizer. In Section 6.1 we prove theoretical properties of SPAMs by assuming that this particular smoother is being used.

3.4. Connection with the grouped lasso

The SPAM model can be thought of as a functional version of the grouped lasso (Yuan and Lin, 2006) as we now explain. Consider the following linear regression model with multiple factors:

$$Y = \sum_{j=1}^{p_n} X_j \beta_j + \varepsilon = X\beta + \varepsilon, \quad (18)$$

where Y is an $n \times 1$ response vector, ε is an $n \times 1$ vector of independent and identically distributed mean 0 noise, X_j is an $n \times d_j$ matrix corresponding to the j th factor and β_j is the corresponding $d_j \times 1$ coefficient vector. Assume for convenience (in this subsection only) that each X_j is orthogonal, so that $X_j^T X_j = I_{d_j}$, where I_{d_j} is the $d_j \times d_j$ identity matrix. We use $X = (X_1, \dots, X_{p_n})$ to denote the full design matrix and use $\beta = (\beta_1^T, \dots, \beta_{p_n}^T)^T$ to denote the parameter.

The *grouped lasso* estimator is defined as the solution of the following convex optimization problem:

$$\hat{\beta}(\lambda_n) = \arg \min_{\beta} \left(\|Y - X\beta\|_2^2 + \lambda_n \sum_{j=1}^{p_n} \sqrt{d_j} \|\beta_j\| \right) \quad (19)$$

where $\sqrt{d_j}$ scales the j th term to compensate for different group sizes.

It is obvious that, when $d_j = 1$ for $j = 1, \dots, p_n$, the grouped lasso becomes the standard lasso. From the Karush–Kuhn–Tucker optimality conditions, a necessary and sufficient condition for $\hat{\beta} = (\hat{\beta}_1^T, \dots, \hat{\beta}_{p_n}^T)^T$ to be the grouped lasso solution is

$$\begin{aligned} -X_j^T(Y - X\hat{\beta}) + \frac{\lambda \sqrt{d_j} \hat{\beta}_j}{\|\hat{\beta}_j\|} &= \mathbf{0}, & \forall \hat{\beta}_j \neq \mathbf{0}, \\ \|X_j^T(Y - X\hat{\beta})\| &\leq \lambda \sqrt{d_j}, & \forall \hat{\beta}_j = \mathbf{0}. \end{aligned} \quad (20)$$

On the basis of this stationary condition, an iterative blockwise co-ordinate descent algorithm can be derived; as shown by Yuan and Lin (2006), a solution to equation (20) satisfies

$$\hat{\beta}_j = \left[1 - \frac{\lambda \sqrt{d_j}}{\|S_j\|} \right]_+ S_j \quad (21)$$

where $S_j = X_j^T(Y - X\beta_{\setminus j})$, with $\beta_{\setminus j} = (\beta_1^T, \dots, \beta_{j-1}^T, \mathbf{0}^T, \beta_{j+1}^T, \dots, \beta_{p_n}^T)$. By iteratively applying equation (21), the grouped lasso solution can be obtained.

As discussed in Section 1, the COSSO model of Lin and Zhang (2006) replaces the lasso constraint on $\sum_j |\beta_j|$ with a reproducing kernel Hilbert space constraint. The advantage of our formulation is that it decouples smoothness ($g_j \in \mathcal{T}_j$) and sparsity ($\sum_j |\beta_j| \leq L$). This leads to a

simple algorithm that can be carried out with any non-parametric smoother and scales easily to high dimensions.

4. Choosing the regularization parameter

We choose λ by minimizing an estimate of the risk. Let ν_j be the effective degrees of freedom for the smoother on the j th variable, i.e. $\nu_j = \text{tr}(\mathcal{S}_j)$ where \mathcal{S}_j is the smoothing matrix for the j th dimension. Also let $\hat{\sigma}^2$ be an estimate of the variance. Define the total effective degrees of freedom as

$$\text{df}(\lambda) = \sum_j \nu_j I(\|\hat{f}_j\| \neq 0). \quad (22)$$

Two estimates of risk are

$$C_p = \frac{1}{n} \sum_{i=1}^n \left\{ Y_i - \sum_{j=1}^p \hat{f}_j(X_j) \right\}^2 + \frac{2\hat{\sigma}^2}{n} \text{df}(\lambda) \quad (23)$$

and

$$\text{GCV}(\lambda) = \frac{(1/n) \sum_{i=1}^n \left\{ Y_i - \sum_j \hat{f}_j(X_{ij}) \right\}^2}{\{1 - \text{df}(\lambda)/n\}^2}. \quad (24)$$

The first is C_p and the second is generalized cross-validation but with degrees of freedom defined by $\text{df}(\lambda)$. A proof that these are valid estimates of risk is not currently available; thus, these should be regarded as heuristics.

On the basis of the results in Wasserman and Roeder (2007) about the lasso, it seems likely that choosing λ by risk estimation can lead to overfitting. One can further clean the estimate by testing $H_0: f_j = 0$ for all j such that $\hat{f}_j \neq 0$. For example, the tests in Fan and Jiang (2005) could be used.

5. Examples

To illustrate the method, we consider a few examples.

5.1. Synthetic data

We generated $n = 100$ observations for an additive model with $p = 100$ and four relevant variables,

$$Y_i = \sum_{j=1}^4 f_j(X_{ij}) + \varepsilon_i,$$

where $\varepsilon_i \sim N(0, 1)$; the relevant component functions are given by

$$\begin{aligned} f_1(x) &= -\sin(1.5x), \\ f_2(x) &= x^3 + 1.5(x - 0.5)^2, \\ f_3(x) &= -\phi(x, 0.5, 0.8^2), \\ f_4(x) &= \sin\{\exp(-0.5x)\} \end{aligned}$$

where $\phi(\cdot, 0.5, 0.8^2)$ is the Gaussian probability distribution function with mean 0.5 and standard deviation 0.8. The data therefore have 96 irrelevant dimensions. The covariates are sampled independent and identically distributed from $\text{uniform}(-2.5, 2.5)$. All the component functions are standardized, i.e.

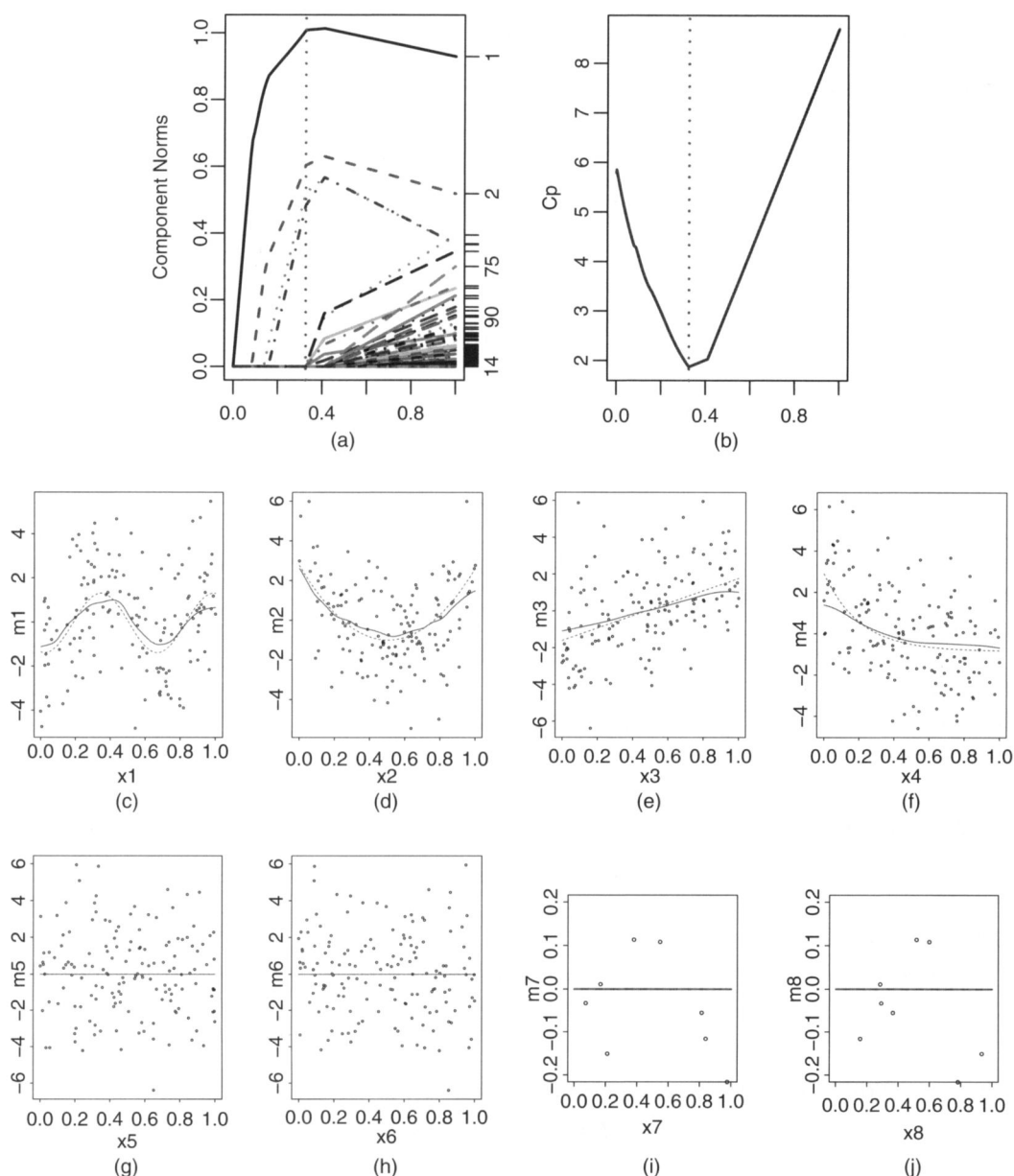


Fig. 1. Simulated data: (a) empirical l_2 -norm of the estimated components as plotted against the regularization parameter λ (the value on the x-axis is proportional to $\sum_j \|\hat{f}_j\|$); (b) C_p -scores against the amount of regularization (\cdot , value of λ which has the smallest C_p -score); estimated (—) versus true additive component functions (---) for (c)–(f) the first four relevant dimensions and (g)–(j) the first four irrelevant dimensions ((c) $I_1 = 97.05$; (d) $I_1 = 88.36$; (e) $I_1 = 90.65$; (f) $I_1 = 79.26$; (g)–(j) $I_1 = 0$)

$$\begin{aligned}\frac{1}{n} \sum_{i=1}^n f_j(X_{ij}) &= 0, \\ \frac{1}{n-1} \sum_{i=1}^n f_j^2(X_{ij}) &= 1.\end{aligned}\tag{25}$$

The results of applying SPAMs are summarized in Fig. 1, using the plug-in bandwidths

$$h_j = 0.6 \text{sd}(X_j)/n^{1/5}.$$

Fig. 1(a) shows regularization paths as the parameter λ varies; each curve is a plot of $\|\hat{f}_j(\lambda)\|$ versus

$$\sum_{k=1}^p \|\hat{f}_k(\lambda)\| / \max_{\lambda} \left\{ \sum_{k=1}^p \|\hat{f}_k(\lambda)\| \right\}\tag{26}$$

for a particular variable X_j . The estimates are generated efficiently over a sequence of λ -values by ‘warm starting’ $\hat{f}_j(\lambda_t)$ at the previous value $\hat{f}_j(\lambda_{t-1})$. Fig. 1(b) shows the C_p -statistic as a function of regularization level.

5.2. Functional sparse coding

Olshausen and Field (1996) proposed a method of obtaining sparse representations of data such as natural images; the motivation comes from trying to understand principles of neural coding. In this example we suggest a non-parametric form of sparse coding.

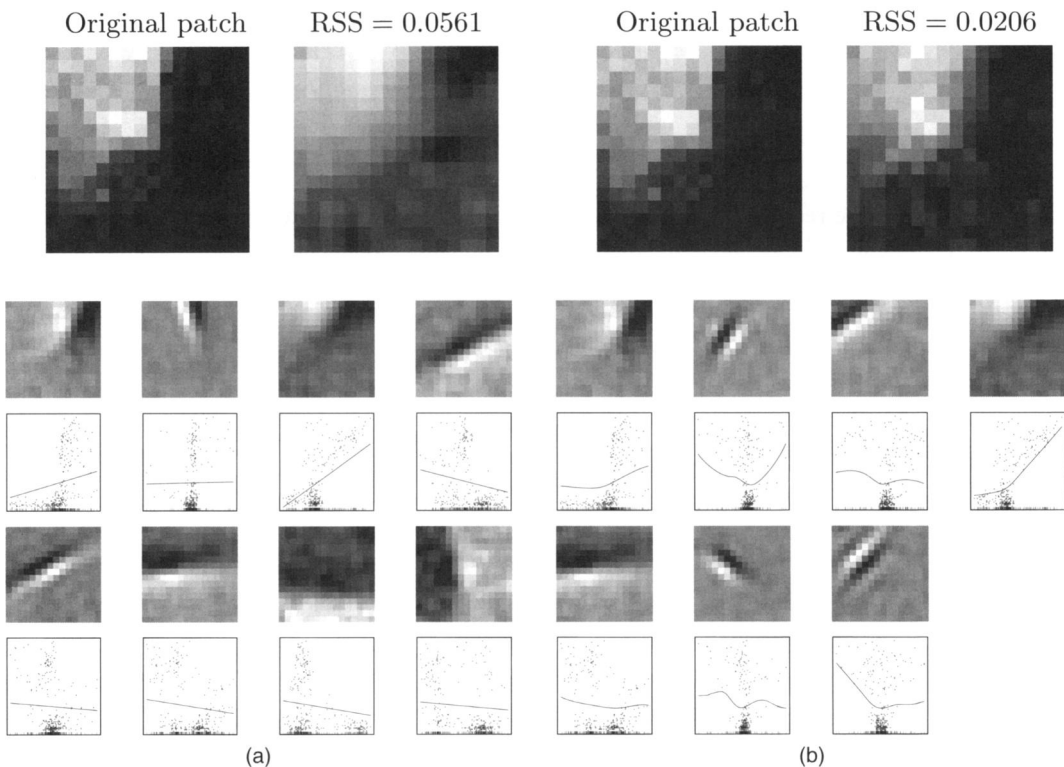


Fig. 2. Comparison of sparse reconstructions by using (a) the lasso and (b) SPAMs

Let $\{y^i\}_{i=1,\dots,N}$ be the data to be represented with respect to some learned basis, where each instance $y^i \in \mathbb{R}^n$ is an n -dimensional vector. The linear sparse coding optimization problem is

$$\min_{\beta, X} \left\{ \sum_{i=1}^N \left(\frac{1}{2n} \|y^i - X\beta^i\|^2 + \lambda \|\beta^i\|_1 \right) \right\} \quad (27)$$

such that

$$\|X_j\| \leq 1. \quad (28)$$

Here X is an $n \times p$ matrix with columns X_j , representing the ‘dictionary’ entries or basis vectors to be learned. It is not required that the basis vectors are orthogonal. The l_1 -penalty on the coefficients β^i encourages sparsity, so each data vector y^i is represented by only a small number of dictionary elements. Sparsity allows the features to specialize, and to capture salient properties of the data.

This optimization problem is not jointly convex in β^i and X . However, for fixed X , each weight vector β^i is computed by running the lasso. For fixed β^i , the optimization is similar to ridge regression and can be solved efficiently. Thus, an iterative procedure for (approximately) solving this optimization problem is easy to derive.

In the case of sparse coding of natural images, as in Olshausen and Field (1996), the basis vectors X_j encode basic edge features at different scales and spatial orientations. In the functional version, we no longer assume a linear parametric fit between the dictionary X and the data y . Instead, we model the relationship by using an additive model. This leads to the following optimization problem for functional sparse coding:

$$\min_{f, X} \left[\sum_{i=1}^N \left\{ \frac{1}{2n} \left\| y^i - \sum_{j=1}^p f_j^i(X_j) \right\|^2 + \lambda \sum_{j=1}^p \|f_j^i\| \right\} \right] \quad (29)$$

such that

$$\|X_j\| \leq 1, \quad j = 1, \dots, p. \quad (30)$$

Fig. 2 illustrates the reconstruction of various image patches by using the sparse linear model compared with the SPAM. Local linear smoothing was used with a Gaussian kernel having fixed bandwidth $h = 0.05$ for all patches and all codewords. The codewords X_j are those obtained by using the Olshausen-Field procedure; these become the design points in the regression estimators. Thus, a codeword for a 16×16 patch corresponds to a vector X_j of dimension 256, with each X_{ij} the grey level for a particular pixel.

6. Theoretical properties

6.1. Sparsistency

In the case of linear regression, with $f_j(X_j) = \beta_j^{*T} X_j$, several researchers have shown that, under certain conditions on n and p , the number of relevant variables $s = |\text{supp}(\beta^*)|$, and the design matrix X , the lasso recovers the sparsity pattern asymptotically, i.e. the lasso estimator $\hat{\beta}_n$ is *sparsistent*:

$$\mathbb{P}\{\text{supp}(\beta^*) = \text{supp}(\hat{\beta}_n)\} \rightarrow 1. \quad (31)$$

Here, $\text{supp}(\beta) = \{j : \beta_j \neq 0\}$. References include Wainwright (2006), Meinshausen and Bühlmann (2006), Zou (2005), Fan and Li (2001) and Zhao and Yu (2007). We show a similar result for SPAMs under orthogonal function regression.

In terms of an orthogonal basis ψ , we can write

$$Y_i = \sum_{j=1}^p \sum_{k=1}^{\infty} \beta_{jk}^* \psi_{jk}(X_{ij}) + \varepsilon_i. \quad (32)$$

To simplify the notation, let β_j be the d_n -dimensional vector $\{\beta_{jk}, k=1, \dots, d_n\}$ and let Ψ_j be the $n \times d_n$ matrix $\Psi_j(i, k) = \psi_{jk}(X_{ij})$. If $A \subset \{1, \dots, p\}$, we denote by Ψ_A the $n \times d|A|$ matrix where, for each $j \in A$, Ψ_j appears as a submatrix in the natural way.

We now analyse the sparse backfitting algorithm of Table 1 by assuming that an orthogonal series smoother is used to estimate the conditional expectation in its step 2. As noted earlier, an orthogonal series smoother for a predictor X_j is the least squares projection onto a truncated set of basis functions $\{\psi_{j1}, \dots, \psi_{jd}\}$. Our optimization problem in this setting is

$$\min_{\beta} \left\{ \frac{1}{2n} \left\| Y - \sum_{j=1}^p \Psi_j \beta_j \right\|_2^2 + \lambda \sum_{j=1}^p \sqrt{\left(\frac{1}{n} \beta_j^T \Psi_j^T \Psi_j \beta_j \right)} \right\}. \quad (33)$$

Combined with the soft thresholding step, the update for f_j in the algorithm in Table 1 can thus be seen to solve the problem

$$\min_{\beta} \left\{ \frac{1}{2n} \|R_j - \Psi_j \beta_j\|_2^2 + \lambda_n \sqrt{\left(\frac{1}{n} \beta_j^T \Psi_j^T \Psi_j \beta_j \right)} \right\}$$

where $\|v\|_2^2$ denotes $\sum_{i=1}^n v_i^2$ and $R_j = Y - \sum_{l \neq j} \Psi_l \beta_l$ is the residual for f_j . The sparse backfitting algorithm thus solves

$$\min_{\beta} \{R_n(\beta) + \lambda_n \Omega(\beta)\} = \min_{\beta} \left(\frac{1}{2n} \left\| Y - \sum_{j=1}^p \Psi_j \beta_j \right\|_2^2 + \lambda_n \sum_{j=1}^p \left\| \frac{1}{\sqrt{n}} \Psi_j \beta_j \right\|_2 \right) \quad (34)$$

where R_n denotes the squared error term and Ω denotes the regularization term, and each β_j is a d_n -dimensional vector. Let S denote the true set of variables $\{j: f_j \neq 0\}$, with $s = |S|$, and let S^c denote its complement. Let $\hat{S}_n = \{j: \hat{\beta}_j \neq 0\}$ denote the estimated set of variables from the minimizer $\hat{\beta}_n$, with corresponding function estimates $\hat{f}_j(x_j) = \sum_{k=1}^{d_n} \hat{\beta}_{jk} \psi_{jk}(x_j)$. For the results in this section, we shall treat the covariates as fixed. A preliminary version of the following result is stated, without proof, in Ravikumar *et al.* (2008).

Theorem 2. Suppose that the following conditions hold on the design matrix X in the orthogonal basis ψ :

$$\Lambda_{\max} \left(\frac{1}{n} \Psi_S^T \Psi_S \right) \leq C_{\max} < \infty, \quad (35)$$

$$\Lambda_{\min} \left(\frac{1}{n} \Psi_S^T \Psi_S \right) \geq C_{\min} > 0, \quad (36)$$

$$\max_{j \in S^c} \left\| \left(\frac{1}{n} \Psi_j^T \Psi_j \right) \left(\frac{1}{n} \Psi_S^T \Psi_S \right)^{-1} \right\| \leq \sqrt{\left(\frac{C_{\min}}{C_{\max}} \right) \frac{1-\delta}{\sqrt{s}}}, \quad \text{for some } 0 < \delta \leq 1. \quad (37)$$

Assume that the truncation dimension d_n satisfies $d_n \rightarrow \infty$ and $d_n = o(n)$. Furthermore, suppose the following conditions, which relate the regularization parameter λ_n to the design parameters n and p , the number of relevant variables s and the truncation size d_n :

$$\frac{s}{d_n \lambda_n} \rightarrow 0, \quad (38)$$

$$\frac{d_n \log\{d_n(p-s)\}}{n\lambda_n^2} \rightarrow 0, \quad (39)$$

$$\frac{1}{\rho_n^*} \left[\sqrt{\left\{ \frac{\log(sd_n)}{n} \right\}} + \frac{s^{3/2}}{d_n} + \lambda_n \sqrt{(sd_n)} \right] \rightarrow 0 \quad (40)$$

where $\rho_n^* = \min_{j \in S} \|\beta_j^*\|_\infty$. Then the solution $\hat{\beta}_n$ to problem (33) is unique and satisfies $\hat{S}_n = S$ with probability approaching 1.

This result parallels the theorem of Wainwright (2006) on model selection consistency of the lasso; however, technical subtleties arise because of the truncation dimension d_n which is increasing with sample size, and the matrix $\Psi_j^T \Psi$ which appears in the regularization of β_j . As a result, the operator norm rather than the ∞ -norm appears in the incoherence condition (37). Note, however, that condition (37) implies that

$$\|\Psi_{S^c}^T \Psi_S (\Psi_S^T \Psi_S)^{-1}\|_\infty = \max_{j \in S^c} \|\Psi_j^T \Psi_S (\Psi_S^T \Psi_S)^{-1}\|_\infty \quad (41)$$

$$\leq \sqrt{\left(\frac{C_{\min} d_n}{C_{\max}} \right)} (1 - \delta) \quad (42)$$

since $(1/\sqrt{n})\|A\|_\infty \leq \|A\| \leq \sqrt{m}\|A\|_\infty$ for an $m \times n$ matrix A . This relates it to the more standard incoherence conditions that have been used for sparsistency in the case of the lasso.

The following corollary, which imposes the additional condition that the number of relevant variables is bounded, follows directly. It makes explicit how to choose the design parameters d_n and λ_n , and implies a condition on the fastest rate at which the minimum norm ρ_n^* can approach 0.

Corollary 1. Suppose that $s = O(1)$, and assume that the design conditions (35)–(37) hold. If the truncation dimension d_n , regularization parameter λ_n and minimum norm ρ_n^* satisfy

$$d_n \asymp n^{1/3}, \quad (43)$$

$$\lambda_n \asymp \frac{\log(np)}{n^{1/3}}, \quad (44)$$

$$\frac{1}{\rho_n^*} = o\left\{ \frac{n^{1/6}}{\log(np)} \right\} \quad (45)$$

then $\mathbb{P}(\hat{S}_n = S) \rightarrow 1$.

The following proposition clarifies the implications of condition (45), by relating the sup-norm $\|\beta_j\|_\infty$ to the function norm $\|f_j\|_2$.

Proposition 1. Suppose that $f(x) = \sum_k \beta_k \psi_k(x)$ is in the Sobolev space of order $\nu > \frac{1}{2}$, so that $\sum_{i=1}^\infty \beta_i^2 i^{2\nu} \leq C^2$ for some constant C . Then

$$\|f\|_2 = \|\beta\|_2 \leq c \|\beta\|_\infty^{2\nu/(2\nu+1)} \quad (46)$$

for some constant c .

For instance, the result of corollary 1 allows the norms of the coefficients β_j to decrease as $\|\beta_j\|_\infty = \log^2(np)/n^{1/6}$. In the case $\nu = 2$, this would allow the norms $\|f_j\|_2$ of the relevant functions to approach 0 at the rate $\log^{8/5}(np)/n^{2/15}$.

6.2. Persistence

The previous assumptions are very strong. They can be weakened at the expense of obtaining weaker results. In particular, in this section we do not assume that the true regression function is additive. We use arguments like those in Juditsky and Nemirovski (2000) and Greenshtein and Ritov (2004) in the context of linear models. In this section we treat X as random and we use triangular array asymptotics, i.e. the joint distribution for the data can change with n . Let (X, Y) denote a new pair (independent of the observed data) and define the predictive risk when predicting Y with $v(X)$ by

$$R(v) = \mathbb{E}\{Y - v(X)\}^2. \quad (47)$$

When $v(x) = \sum_j \beta_j g_j(x_j)$ we also write the risk as $R(\beta, g)$ where $\beta = (\beta_1, \dots, \beta_p)$ and $g = (g_1, \dots, g_p)$. Following Greenshtein and Ritov (2004) we say that an estimator \hat{m}_n is persistent (risk consistent) relative to a class of functions \mathcal{M}_n , if

$$R(\hat{m}_n) - R(m_n^*) \xrightarrow{P} 0 \quad (48)$$

where

$$m_n^* = \arg \min_{v \in \mathcal{M}_n} \{R(v)\} \quad (49)$$

is the predictive oracle. Greenshtein and Ritov (2004) showed that the lasso is persistent for $\mathcal{M}_n = \{l(x) = x^T \beta : \|\beta\|_1 \leq L_n\}$ and $L_n = o\{n/\log(n)^{1/4}\}$. Note that m_n^* is the best linear approximation (in prediction risk) in \mathcal{M}_n but the true regression function is not assumed to be linear. Here we show a similar result for SPAMs.

In this section, we assume that the SPAM estimator \hat{m}_n is chosen to minimize

$$\frac{1}{n} \sum_{i=1}^n \left\{ Y_i - \sum_j \beta_j g_j(X_{ij}) \right\}^2 \quad (50)$$

subject to $\|\beta\|_1 \leq L_n$ and $g_j \in \mathcal{T}_j$. We make no assumptions about the design matrix. Let $\mathcal{M}_n \equiv \mathcal{M}_n(L_n)$ be defined by

$$\mathcal{M}_n = \left\{ m : m(x) = \sum_{j=1}^{p_n} \beta_j g_j(x_j) : \mathbb{E}(g_j) = 0, \mathbb{E}(g_j^2) = 1, \sum_j |\beta_j| \leq L_n \right\} \quad (51)$$

and let $m_n^* = \arg \min_{v \in \mathcal{M}_n} \{R(v)\}$.

Theorem 3. Suppose that $p_n \leq \exp(n^\xi)$ for some $\xi < 1$. Then,

$$R(\hat{m}_n) - R(m_n^*) = O_P\left(\frac{L_n^2}{n^{(1-\xi)/2}}\right) \quad (52)$$

and hence, if $L_n = o(n^{(1-\xi)/4})$, then the SPAM is persistent.

7. Discussion

The results that are presented here show how many of the recently established theoretical properties of l_1 -regularization for linear models extend to SPAMs. The sparse backfitting algorithm that we have derived is attractive because it decouples smoothing and sparsity, and can be used with any non-parametric smoother. It thus inherits the nice properties of the original backfitting procedure. However, our theoretical analyses have made use of a particular form of smoothing, using a truncated orthogonal basis. An important problem is thus to extend the theory to cover more general classes of smoothing operators. Convergence properties of the SPAM backfitting

algorithm should also be investigated; convergence of special cases of standard backfitting was studied by Buja *et al.* (1989).

An additional direction for future work is to develop procedures for automatic bandwidth selection in each dimension. We have used plug-in bandwidths and truncation dimensions d_n in our experiments and theory. It is of particular interest to develop procedures that are adaptive to different levels of smoothness in different dimensions. It would also be of interest to consider more general penalties of the form $p_\lambda(\|f_j\|)$, as in Fan and Li (2001).

Finally, we note that, although we have considered basic additive models that allow functions of individual variables, it is natural to consider interactions, as in the functional analysis-of-variance model. One challenge is to formulate suitable incoherence conditions on the functions that enable regularization-based procedures or greedy algorithms to recover the correct interaction graph. In the parametric setting, one result in this direction is Wainwright *et al.* (2007).

Acknowledgements

This research was supported in part by National Science Foundation grant CCF-0625879 and a Siebel scholarship to PR.

Appendix A: Proofs

A.1. Proof of theorem 1

Consider the minimization of the Lagrangian

$$\min_{\{f_j \in \mathcal{H}_j\}} \{\mathcal{L}(f, \lambda)\} \equiv \frac{1}{2} \mathbb{E} \left\{ Y - \sum_{j=1}^p f_j(X_j) \right\}^2 + \lambda \sum_{j=1}^p \sqrt{\mathbb{E}\{f_j(X_j)^2\}} \quad (53)$$

with respect to $f_j \in \mathcal{H}_j$, holding the other components $\{f_k, k \neq j\}$ fixed. The stationary condition is obtained by setting the Fréchet derivative to 0. Denote by $\partial_j \mathcal{L}(f, \lambda; \eta_j)$ the directional derivative with respect to f_j in the direction $\eta_j(X_j) \in \mathcal{H}_j \{\mathbb{E}(\eta_j) = 0, \mathbb{E}(\eta_j^2) < \infty\}$. Then the stationary condition can be formulated as

$$\partial_j \mathcal{L}(f, \lambda; \eta_j) = \frac{1}{2} \mathbb{E}\{(f_j - R_j + \lambda v_j) \eta_j\} = 0 \quad (54)$$

where $R_j = Y - \sum_{k \neq j} f_k$ is the residual for f_j , and $v_j \in \mathcal{H}_j$ is an element of the subgradient $\partial \sqrt{\mathbb{E}(f_j^2)}$, satisfying $v_j = f_j / \sqrt{\mathbb{E}(f_j^2)}$ if $\mathbb{E}(f_j^2) \neq 0$ and $v_j \in \{u_j \in \mathcal{H}_j | \mathbb{E}(u_j^2) \leq 1\}$ otherwise.

Using iterated expectations, the above condition can be rewritten as

$$\mathbb{E}\{(f_j + \lambda v_j - \mathbb{E}(R_j | X_j)) \eta_j\} = 0. \quad (55)$$

But, since $f_j - \mathbb{E}(R_j | X_j) + \lambda v_j \in \mathcal{H}_j$, we can compute the derivative in the direction $\eta_j = f_j - \mathbb{E}(R_j | X_j) + \lambda v_j \in \mathcal{H}_j$, implying that

$$\mathbb{E}\{(f_j(x_j) - \mathbb{E}(R_j | X_j = x_j) + \lambda v_j(x_j))^2\} = 0, \quad (56)$$

i.e.

$$f_j + \lambda v_j = \mathbb{E}(R_j | X_j) \quad \text{almost everywhere.} \quad (57)$$

Denote the conditional expectation $\mathbb{E}(R_j | X_j)$ —also the projection of the residual R_j onto \mathcal{H}_j —by P_j . Now, if $\mathbb{E}(f_j^2) \neq 0$, then $v_j = f_j / \sqrt{\mathbb{E}(f_j^2)}$, which from condition (57) implies

$$\sqrt{\mathbb{E}(P_j^2)} = \sqrt{\mathbb{E}\{f_j + \lambda f_j / \sqrt{\mathbb{E}(f_j^2)}\}^2} \quad (58)$$

$$= \left\{ 1 + \frac{\lambda}{\sqrt{\mathbb{E}(f_j^2)}} \right\} \sqrt{\mathbb{E}(f_j^2)} \quad (59)$$

$$= \sqrt{\mathbb{E}(f_j^2)} + \lambda \quad (60)$$

$$\geq \lambda. \quad (61)$$

If $\mathbb{E}(f_j^2) = 0$, then $f_j = 0$ almost everywhere, and $\sqrt{\mathbb{E}(v_j^2)} \leq 1$. Equation (57) then implies that

$$\sqrt{\mathbb{E}(P_j^2)} \leq \lambda. \quad (62)$$

We thus obtain the equivalence

$$\sqrt{\mathbb{E}(P_j^2)} \leq \lambda \Leftrightarrow f_j = 0 \quad \text{almost everywhere.} \quad (63)$$

Rewriting equation (57) in light of result (63), we obtain

$$\begin{cases} 1 + \frac{\lambda}{\sqrt{\mathbb{E}(f_j^2)}} \} f_j = P_j & \text{if } \sqrt{\mathbb{E}(P_j^2)} > \lambda, \\ f_j = 0 & \text{otherwise.} \end{cases}$$

Using equation (60), we thus arrive at the soft thresholding update for f_j :

$$f_j = \left[1 - \frac{\lambda}{\sqrt{\mathbb{E}(P_j^2)}} \right]_+ P_j \quad (64)$$

where $[\cdot]_+$ denotes the positive part and $P_j = \mathbb{E}[R_j | X_j]$.

A.2. Proof of theorem 2

A vector $\hat{\beta} \in \mathbb{R}^{d_n p}$ is an optimum of the objective function in expression (34) if and only if there is a subgradient $\hat{g} \in \partial\Omega(\hat{\beta})$, such that

$$\frac{1}{n} \Psi^T \left(\sum_j \Psi_j \hat{\beta}_j - Y \right) + \lambda_n \hat{g} = 0. \quad (65)$$

The subdifferential $\partial\Omega(\beta)$ is the set of vectors $g \in \mathbb{R}^{p d_n}$ satisfying

$$\begin{aligned} g_j &= \frac{(1/n) \Psi_j^T \Psi_j \beta_j}{\sqrt{\{(1/n) \beta_j^T \Psi_j^T \Psi_j \beta_j\}}} & \text{if } \beta_j \neq 0, \\ g_j^T \left(\frac{1}{n} \Psi_j^T \Psi_j \right)^{-1} g_j &\leq 1 & \text{if } \beta_j = 0. \end{aligned}$$

Our argument is based on the technique of a *primal dual witness*, which has been used previously in the analysis of the lasso (Wainwright, 2006). In particular, we construct a coefficient subgradient pair $(\hat{\beta}, \hat{g})$ which satisfies $\text{supp}(\hat{\beta}) = \text{supp}(\beta^*)$ and in addition satisfies the optimality conditions for the objective (34) with high probability. Thus, when the procedure succeeds, the constructed coefficient vector $\hat{\beta}$ is equal to the solution of the convex objective (34), and \hat{g} is an optimal solution to its dual. From its construction, the support of $\hat{\beta}$ is equal to the true support $\text{supp}(\beta^*)$, from which we can conclude that the solution of the objective (34) is sparsistent. The construction of the primal dual witness proceeds as follows.

- (a) Set $\hat{\beta}_{S^c} = 0$.
- (b) Set $\hat{g}_S = \partial\Omega(\beta^*)_S$.
- (c) With these settings of $\hat{\beta}_{S^c}$ and \hat{g}_S , obtain $\hat{\beta}_S$ and \hat{g}_{S^c} from the stationary conditions in equation (65).

For the witness procedure to succeed, we must show that $(\hat{\beta}, \hat{g})$ is optimal for the objective (34), meaning that

$$\hat{\beta}_j \neq 0 \quad \text{for } j \in S, \quad (66a)$$

$$g_j^T \left(\frac{1}{n} \Psi_j^T \Psi_j \right)^{-1} g_j < 1 \quad \text{for } j \in S^c. \quad (66b)$$

For uniqueness of the solution, we require strict dual feasibility, meaning strict inequality in condition (66b). In what follows, we show that these two conditions hold with high probability.

A.2.1. Condition (66a)

Setting $\hat{\beta}_{sc} = 0$ and

$$\hat{g}_j = \frac{(1/n)\Psi_j^T \Psi_j \beta_j^*}{\sqrt{\{(1/n)\beta_j^{*T} \Psi_j^T \Psi_j \beta_j^*\}}} \quad \text{for } j \in S,$$

the stationarity condition for $\hat{\beta}_S$ is given by

$$\frac{1}{n} \Psi_S^T (\Psi_S \hat{\beta}_S - Y) + \lambda_n \hat{g}_S = 0. \quad (67)$$

Let $V = Y - \Psi_S \beta_S^* - W$ denote the error due to finite truncation of the orthogonal basis, where $W = (\varepsilon_1, \dots, \varepsilon_n)^T$. Then the stationarity condition (67) can be simplified as

$$\frac{1}{n} \Psi_S^T \Psi_S (\hat{\beta}_S - \beta_S^*) - \frac{1}{n} \Psi_S^T W - \frac{1}{n} \Psi_S^T V + \lambda_n \hat{g}_S = 0,$$

so that

$$\hat{\beta}_S - \beta_S^* = \left(\frac{1}{n} \Psi_S^T \Psi_S \right)^{-1} \left(\frac{1}{n} \Psi_S^T W + \frac{1}{n} \Psi_S^T V - \lambda_n \hat{g}_S \right), \quad (68)$$

where we have used the assumption that $(1/n)\Psi_S^T \Psi_S$ is non-singular. Recalling our definition of the minimum function norm $\rho_n^* = \min_{j \in S} \|\beta_j^*\|_\infty > 0$, it suffices to show that $\|\hat{\beta}_S - \beta_S^*\|_\infty < \rho_n^*/2$, to ensure that

$$\text{supp}(\beta_S^*) = \text{supp}(\hat{\beta}_S) = \{j : \|\hat{\beta}_j\|_\infty \neq 0\},$$

so that condition (66a) would be satisfied. Using $\Sigma_{SS} = (1/n)(\Psi_S^T \Psi_S)$ to simplify the notation, we have the l_∞ -bound,

$$\|\hat{\beta}_S - \beta_S^*\|_\infty \leq \underbrace{\left\| \Sigma_{SS}^{-1} \left(\frac{1}{n} \Psi_S^T W \right) \right\|_\infty}_{T_1} + \underbrace{\left\| \Sigma_{SS}^{-1} \left(\frac{1}{n} \Psi_S^T V \right) \right\|_\infty}_{T_2} + \lambda_n \underbrace{\|\Sigma_{SS}^{-1} \hat{g}_S\|_\infty}_{T_3}. \quad (69)$$

We now proceed to bound the quantities T_1 , T_2 and T_3 .

A.2.2. Bounding T_3

Note that, for $j \in S$,

$$1 = g_j^T \left(\frac{1}{n} \Psi_j^T \Psi_j \right)^{-1} g_j \geq \frac{1}{C_{\max}} \|g_j\|^2,$$

and thus $\|g_j\| \leq \sqrt{C_{\max}}$. Noting further that

$$\|g_S\|_\infty = \max_{j \in S} (\|g_j\|_\infty) \leq \max_{j \in S} (\|g_j\|_2) \leq \sqrt{C_{\max}}, \quad (70)$$

it follows that

$$T_3 := \|\Sigma_{SS}^{-1} \hat{g}_S\|_\infty \leq \sqrt{C_{\max}} \|\Sigma_{SS}^{-1}\|_\infty. \quad (71)$$

A.2.3. Bounding T_2

We proceed in two steps; we first bound $\|V\|_\infty$ and use this to bound $\|(1/n)\Psi_S^T V\|_\infty$. Note that, as we are working over the Sobolev spaces \mathcal{S}_j of order 2,

$$\begin{aligned} |V_i| &= \left| \sum_{j \in S} \sum_{k=d_n+1}^{\infty} \beta_{jk}^* \Psi_{jk}(X_{ij}) \right| \leq B \sum_{j \in S} \sum_{k=d_n+1}^{\infty} |\beta_{jk}^*| \\ &= B \sum_{j \in S} \sum_{k=d_n+1}^{\infty} \frac{|\beta_{jk}^*| k^2}{k^2} \leq B \sum_{j \in S} \sqrt{\left(\sum_{k=d_n+1}^{\infty} \beta_{jk}^{*2} k^4 \right)} \sqrt{\left(\sum_{k=d_n+1}^{\infty} \frac{1}{k^4} \right)} \\ &\leq sBC \sqrt{\left(\sum_{k=d_n+1}^{\infty} \frac{1}{k^4} \right)} \leq \frac{sB'}{d_n^{3/2}}, \end{aligned}$$

for some constant $B' > 0$. It follows that

$$\left| \frac{1}{n} \Psi_{jk}^T V \right| \leq \left| \frac{1}{n} \sum_i \Psi_{jk}(X_{ij}) \right| \|V\|_\infty \leq \frac{Ds}{d_n^{3/2}}, \quad (72)$$

where D denotes a generic constant. Thus,

$$T_2 := \left\| \Sigma_{SS}^{-1} \left(\frac{1}{n} \Psi_S^T V \right) \right\|_\infty \leq \|\Sigma_{SS}^{-1}\|_\infty \frac{Ds}{d_n^{3/2}}. \quad (73)$$

A.2.4. Bounding T_1

Let $Z = T_1 = \Sigma_{SS}^{-1} (1/n) \Psi_S^T W$. Note that $W \sim N(0, \sigma^2 I)$, so that Z is Gaussian as well, with mean 0. Consider its l th component, $Z_l = e_l^T Z$. Then $\mathbb{E}(Z_l) = 0$, and

$$\text{var}(Z_l) = \frac{\sigma^2}{n} e_l^T \Sigma_{SS}^{-1} e_l \leq \frac{\sigma^2}{C_{\min} n}.$$

By Gaussian comparison results (Ledoux and Talagrand, 1991), we have then that

$$\mathbb{E}(\|Z\|_\infty) \leq 3\sqrt{\{\log(sd_n)\|\text{var}(Z)\|_\infty\}} \leq 3\sigma\sqrt{\left\{\frac{\log(sd_n)}{nC_{\min}}\right\}}. \quad (74)$$

Substituting the bounds for T_2 and T_3 from equations (73) and (71) respectively into equation (69), and using the bound for the expected value of T_1 from inequality (74), it follows from an application of Markov's inequality that

$$\begin{aligned} \mathbb{P}\left(\|\hat{\beta}_S - \beta_S^*\|_\infty > \frac{\rho_n^*}{2}\right) &\leq \mathbb{P}\left\{\|Z\|_\infty + \|\Sigma_{SS}^{-1}\|_\infty (Ds d_n^{-3/2} + \lambda_n \sqrt{C_{\max}}) > \frac{\rho_n^*}{2}\right\} \\ &\leq \frac{2}{\rho_n^*} \{\mathbb{E}(\|Z\|_\infty) + \|\Sigma_{SS}^{-1}\|_\infty (Ds d_n^{-3/2} + \lambda_n \sqrt{C_{\max}})\} \\ &\leq \frac{2}{\rho_n^*} \left[3\sigma\sqrt{\left\{\frac{\log(sd_n)}{nC_{\min}}\right\}} + \|\Sigma_{SS}^{-1}\|_\infty \left(\frac{Ds}{d_n^{3/2}} + \lambda_n \sqrt{C_{\max}} \right) \right], \end{aligned}$$

which converges to 0 under the condition that

$$\frac{1}{\rho_n^*} \left[\sqrt{\left\{\frac{\log(sd_n)}{n}\right\}} + \left\| \left(\frac{1}{n} \Psi_S^T \Psi_S \right)^{-1} \right\|_\infty \left(\frac{s}{d_n^{3/2}} + \lambda_n \right) \right] \rightarrow 0. \quad (75)$$

Noting that

$$\left\| \left(\frac{1}{n} \Psi_S^T \Psi_S \right)^{-1} \right\|_\infty \leq \frac{\sqrt{(sd_n)}}{C_{\min}}, \quad (76)$$

it follows that condition (75) holds when

$$\frac{1}{\rho_n^*} \left[\sqrt{\left\{\frac{\log(sd_n)}{n}\right\}} + \frac{s^{3/2}}{d_n} + \lambda_n \sqrt{(sd_n)} \right] \rightarrow 0. \quad (77)$$

But this is satisfied by assumption (40) in the theorem. We have thus shown that condition (66a) is satisfied with probability converging to 1.

A.2.5. Condition (66b)

We now must consider the dual variables \hat{g}_{S^c} . Recall that we have set $\hat{\beta}_{S^c} = \beta_{S^c}^* = 0$. The stationarity condition for $j \in S^c$ is thus given by

$$\frac{1}{n} \Psi_j^T (\Psi_S \hat{\beta}_S - \Psi_S \beta_S^* - W - V) + \lambda_n \hat{g}_j = 0.$$

It then follows from equation (68) that

$$\begin{aligned}\hat{g}_{S^c} &= \frac{1}{\lambda_n} \left\{ \frac{1}{n} \Psi_{S^c}^T \Psi_S (\beta_S^* - \hat{\beta}_S) + \frac{1}{n} \Psi_{S^c}^T (W + V) \right\} \\ &= \frac{1}{\lambda_n} \left\{ \frac{1}{n} \Psi_{S^c}^T \Psi_S \left(\frac{1}{n} \Psi_S^T \Psi_S \right)^{-1} \left(\lambda_n \hat{g}_S - \frac{1}{n} \Psi_S^T W - \frac{1}{n} \Psi_S^T V \right) + \frac{1}{n} \Psi_{S^c}^T (W + V) \right\},\end{aligned}$$

so

$$\hat{g}_{S^c} = \frac{1}{\lambda_n} \left\{ \Sigma_{S^c S} \Sigma_{SS}^{-1} \left(\lambda_n \hat{g}_S - \frac{1}{n} \Psi_S^T W - \frac{1}{n} \Psi_S^T V \right) + \frac{1}{n} \Psi_{S^c}^T (W + V) \right\}. \quad (78)$$

Condition (66b) requires that

$$g_j^T \left(\frac{1}{n} \Psi_j^T \Psi_j \right)^{-1} g_j < 1, \quad (79)$$

for all $j \in S^c$. Since

$$g_j^T \left(\frac{1}{n} \Psi_j^T \Psi_j \right)^{-1} g_j \leq \frac{1}{C_{\min}} \|g_j\|^2 \quad (80)$$

it suffices to show that $\max_{j \in S^c} \|g_j\| < \sqrt{C_{\min}}$. From equation (78), we see that \hat{g}_j is Gaussian, with mean μ_j as

$$\mu_j = \mathbb{E}(\hat{g}_j) = \Sigma_{jS} \Sigma_{SS}^{-1} \left(\hat{g}_S - \frac{1}{\lambda_n} \frac{1}{n} \Psi_S^T V \right) - \frac{1}{\lambda_n} \frac{1}{n} \Psi_j^T V.$$

This can be bounded as

$$\begin{aligned}\|\mu_j\| &\leq \|\Sigma_{jS} \Sigma_{SS}^{-1}\| \left(\|\hat{g}_S\| + \frac{1}{\lambda_n} \left\| \frac{1}{n} \Psi_S^T V \right\| \right) + \frac{1}{\lambda_n} \left\| \frac{1}{n} \Psi_j^T V \right\| \\ &= \|\Sigma_{jS} \Sigma_{SS}^{-1}\| \left\{ \sqrt{(sC_{\max})} + \frac{1}{\lambda_n} \left\| \frac{1}{n} \Psi_S^T V \right\| \right\} + \frac{1}{\lambda_n} \left\| \frac{1}{n} \Psi_j^T V \right\|.\end{aligned} \quad (81)$$

Using the bound $\|\Psi_j^T V\|_\infty \leq Ds/d_n^{3/2}$ from equation (72), we have

$$\left\| \frac{1}{n} \Psi_j^T V \right\| \leq \sqrt{d_n} \left\| \frac{1}{n} \Psi_j^T V \right\|_\infty \leq \frac{Ds}{d_n},$$

and hence

$$\left\| \frac{1}{n} \Psi_S^T V \right\| \leq \sqrt{s} \left\| \frac{1}{n} \Psi_S^T V \right\|_\infty \leq \frac{Ds^{3/2}}{d_n}.$$

Substituting in the bound (81) on the mean μ_j ,

$$\|\mu_j\| \leq \|\Sigma_{jS} \Sigma_{SS}^{-1}\| \left\{ \sqrt{(sC_{\max})} + \frac{Ds^{3/2}}{\lambda_n d_n} \right\} + \frac{Ds}{\lambda_n d_n}. \quad (82)$$

Assumptions (37) and (38) of the theorem can be rewritten as

$$\|\Sigma_{jS} \Sigma_{SS}^{-1}\| \leq \sqrt{\left(\frac{C_{\min}}{C_{\max}} \right) \frac{1-\delta}{\sqrt{s}}} \quad \text{for some } \delta > 0, \quad (83)$$

$$\frac{s}{\lambda_n d_n} \rightarrow 0. \quad (84)$$

Thus the bound on the mean becomes

$$\|\mu_j\| \leq \sqrt{C_{\min}}(1-\delta) + \frac{2Ds}{\lambda_n d_n} < \sqrt{C_{\min}},$$

for sufficiently large n . It therefore suffices, for condition (66b) to be satisfied, to show that

$$\mathbb{P}\left(\max_{j \in S^c} \|\hat{g}_j - \mu_j\|_\infty > \frac{\delta}{2\sqrt{d_n}}\right) \rightarrow 0, \quad (85)$$

since this implies that

$$\begin{aligned} \|\hat{g}_j\| &\leq \|\mu_j\| + \|\hat{g}_j - \mu_j\| \\ &\leq \|\mu_j\| + \sqrt{d_n} \|\hat{g}_j - \mu_j\|_\infty \\ &\leq \sqrt{C_{\min}}(1-\delta) + \frac{\delta}{2} + o(1), \end{aligned}$$

with probability approaching 1. To show result (85), we again appeal to Gaussian comparison results. Define

$$Z_j = \Psi_j^T (I - \Psi_S (\Psi_S^T \Psi_S)^{-1} \Psi_S^T) \frac{W}{n}, \quad (86)$$

for $j \in S^c$. Then Z_j are zero-mean Gaussian random variables, and we need to show that

$$\mathbb{P}\left\{\max_{j \in S^c} \left(\frac{\|Z_j\|_\infty}{\lambda_n}\right) \geq \frac{\delta}{2\sqrt{d_n}}\right\} \rightarrow 0. \quad (87)$$

A calculation shows that $\mathbb{E}(Z_{jk}^2) \leq \sigma^2/n$. Therefore, we have by Markov's inequality and Gaussian comparison that

$$\begin{aligned} \mathbb{P}\left\{\max_{j \in S^c} \left(\frac{\|Z_j\|_\infty}{\lambda_n}\right) \geq \frac{\delta}{2\sqrt{d_n}}\right\} &\leq \frac{2\sqrt{d_n}}{\delta \lambda_n} \mathbb{E}(\max_{jk} |Z_{jk}|) \\ &\leq \frac{2\sqrt{d_n}}{\delta \lambda_n} [3\sqrt{\log\{(p-s)d_n\}} \max_{jk} \{\sqrt{\mathbb{E}(Z_{jk}^2)}\}] \\ &\leq \frac{6\sigma}{\delta \lambda_n} \sqrt{\left[\frac{d_n \log\{(p-s)d_n\}}{n}\right]}, \end{aligned}$$

which converges to 0 given the assumption (39) of the theorem that

$$\frac{\lambda_n^2 n}{d_n \log\{(p-s)d_n\}} \rightarrow \infty.$$

Thus condition (66b) is also satisfied with probability converging to 1, which completes the proof.

A.3. Proof of proposition 1

For any index k we have that

$$\|f\|_2^2 = \sum_{i=1}^{\infty} \beta_i^2 \quad (88)$$

$$\leq \|\beta\|_\infty \sum_{i=1}^{\infty} |\beta_i| \quad (89)$$

$$= \|\beta\|_\infty \sum_{i=1}^k |\beta_i| + \|\beta\|_\infty \sum_{i=k+1}^{\infty} |\beta_i| \quad (90)$$

$$\leq k \|\beta\|_\infty^2 + \|\beta\|_\infty \sum_{i=k+1}^{\infty} \frac{i^\nu |\beta_i|}{i^\nu} \quad (91)$$

$$\leq k\|\beta\|_\infty^2 + \|\beta\|_\infty \sqrt{\left(\sum_{i=1}^{\infty} \beta_i^2 i^{2\nu}\right) \left(\sum_{i=k+1}^{\infty} \frac{1}{i^{2\nu}}\right)} \quad (92)$$

$$\leq k\|\beta\|_\infty^2 + \|\beta\|_\infty C \sqrt{\left(\frac{k^{1-2\nu}}{2\nu-1}\right)}, \quad (93)$$

where the last inequality uses the bound

$$\sum_{i=k+1}^{\infty} i^{-2\nu} \leq \int_k^{\infty} x^{-2\nu} dx = \frac{k^{1-2\nu}}{2\nu-1}. \quad (94)$$

Let k^* be the index that minimizes expression (93). Some calculus shows that k^* satisfies

$$c_1 \|\beta\|_\infty^{-2/(2\nu+1)} \leq k^* \leq c_2 \|\beta\|_\infty^{-2/(2\nu+1)} \quad (95)$$

for some constants c_1 and c_2 . Using the above expression in expression (93) then yields

$$\|f\|_2^2 \leq \|\beta\|_\infty (c_2 \|\beta\|_\infty^{(2\nu-1)/(2\nu+1)} + c'_1 \|\beta\|_\infty^{(2\nu-1)/(2\nu+1)}) \quad (96)$$

$$= c \|\beta\|_\infty^{4\nu/(2\nu+1)} \quad (97)$$

for some constant c , and the result follows.

A.4. Proof of theorem 3

We begin with some notation. If \mathcal{M} is a class of functions then the L_∞ bracketing number $N_{[]}(\varepsilon, \mathcal{M})$ is defined as the smallest number of pairs $B = \{(l_1, u_1), \dots, (l_k, u_k)\}$ such that $\|u_j - l_j\|_\infty \leq \varepsilon$, $1 \leq j \leq k$, and such that for every $m \in \mathcal{M}$ there exists $(l, u) \in B$ such that $l \leq m \leq u$. For the Sobolev space \mathcal{T}_j ,

$$\log\{N_{[]}(\varepsilon, \mathcal{T}_j)\} \leq K \left(\frac{1}{\varepsilon}\right)^{1/2} \quad (98)$$

for some $K > 0$; see van der Vaart (1998). The bracketing integral is defined to be

$$J_{[]}(\delta, \mathcal{M}) = \int_0^\delta \sqrt{\log\{N_{[]}(\varepsilon, \mathcal{M})\}} d\varepsilon. \quad (99)$$

From corollary 19.35 of van der Vaart (1998),

$$\mathbb{E} \left\{ \sup_{g \in \mathcal{M}} |\hat{\mu}(g) - \mu(g)| \right\} \leq \frac{C J_{[]}(\|F\|_\infty, \mathcal{M})}{\sqrt{n}} \quad (100)$$

for some $C > 0$, where $F(x) = \sup_{g \in \mathcal{M}} |g(x)|$, $\mu(g) = \mathbb{E}\{g(X)\}$ and $\hat{\mu}(g) = n^{-1} \sum_{i=1}^n g(X_i)$.

Set $Z \equiv (Z_0, \dots, Z_p) = (Y, X_1, \dots, X_p)$ and note that

$$R(\beta, g) = \sum_{j=0}^p \sum_{k=0}^p \beta_j \beta_k \mathbb{E}\{g_j(Z_j) g_k(Z_k)\} \quad (101)$$

where we define $g_0(z_0) = z_0$ and $\beta_0 = -1$. Also define

$$\hat{R}(\beta, g) = \frac{1}{n} \sum_{i=1}^n \sum_{j=0}^p \sum_{k=0}^p \beta_j \beta_k g_j(Z_{ij}) g_k(Z_{ik}). \quad (102)$$

Hence \hat{m}_n is the minimizer of $\hat{R}(\beta, g)$ subject to the constraint $\sum_j \beta_j g_j(x_j) \in \mathcal{M}_n(L_n)$ and $g_j \in \mathcal{T}_j$. For all (β, g) ,

$$|\hat{R}(\beta, g) - R(\beta, g)| \leq \|\beta\|_1^2 \max_{jk} \sup_{g_j \in \mathcal{T}_j, g_k \in \mathcal{T}_k} |\hat{\mu}_{jk}(g) - \mu_{jk}(g)| \quad (103)$$

where

$$\hat{\mu}_{jk}(g) = n^{-1} \sum_{i=1}^n \sum_{k=0}^p g_j(Z_{ij}) g_k(Z_{ik})$$

and $\mu_{jk}(g) = \mathbb{E}\{g_j(Z_j)g_k(Z_k)\}$. From inequality (98) it follows that

$$\log\{N_{\square}(\varepsilon, \mathcal{M}_n)\} \leq 2 \log(p_n) + K \left(\frac{1}{\varepsilon}\right)^{1/2}. \quad (104)$$

Hence, $J_{\square}(C, \mathcal{M}_n) = O\{\sqrt{\log(p_n)}\}$ and it follows from inequality (100) and Markov's inequality that

$$\max_{jk} \sup_{g_j \in \mathcal{S}_j, g_k \in \mathcal{S}_k} |\hat{\mu}_{jk}(g) - \mu_{jk}(g)| = O_P \left[\sqrt{\left\{ \frac{\log(p_n)}{n} \right\}} \right] = O_P \left(\frac{1}{n^{(1-\xi)/2}} \right). \quad (105)$$

We conclude that

$$\sup_{g \in \mathcal{M}} |\hat{R}(g) - R(g)| = O_P \left(\frac{L_n^2}{n^{(1-\xi)/2}} \right). \quad (106)$$

Therefore,

$$\begin{aligned} R(m^*) &\leq R(\hat{m}_n) \leq \hat{R}(\hat{m}_n) + O_P \left(\frac{L_n^2}{n^{(1-\xi)/2}} \right) \\ &\leq \hat{R}(m^*) + O_P \left(\frac{L_n^2}{n^{(1-\xi)/2}} \right) \leq R(m^*) + O_P \left(\frac{L_n^2}{n^{(1-\xi)/2}} \right) \end{aligned}$$

and the conclusion follows.

References

- Antoniadis, A. and Fan, J. (2001) Regularized wavelet approximations (with discussion). *J. Am. Statist. Ass.*, **96**, 939–967.
- BuJa, A., Hastie, T. and Tibshirani, R. (1989) Linear smoothers and additive models. *Ann. Statist.*, **17**, 453–510.
- Bunea, F., Tsybakov, A. and Wegkamp, M. (2007) Sparsity oracle inequalities for the lasso. *Electron. J. Statist.*, **1**, 169–194.
- Daubechies, I., Defrise, M. and DeMol, C. (2004) An iterative thresholding algorithm for linear inverse problems. *Commun. Pure Appl. Math.*, **57**, 1413–1457.
- Daubechies, I., Fornasier, M. and Loris, I. (2007) Accelerated projected gradient method for linear inverse problems with sparsity constraints. *Technical Report*. Princeton University, Princeton. (Available from arXiv:0706.4297.)
- Fan, J. and Jiang, J. (2005) Nonparametric inference for additive models. *J. Am. Statist. Ass.*, **100**, 890–907.
- Fan, J. and Li, R. Z. (2001) Variable selection via penalized likelihood. *J. Am. Statist. Ass.*, **96**, 1348–1360.
- Greenshtein, E. and Ritov, Y. (2004) Persistency in high dimensional linear predictor-selection and the virtue of over-parametrization. *Bernoulli*, **10**, 971–988.
- Hastie, T. and Tibshirani, R. (1999) *Generalized Additive Models*. New York: Chapman and Hall.
- Juditsky, A. and Nemirovski, A. (2000) Functional aggregation for nonparametric regression. *Ann. Statist.*, **28**, 681–712.
- Koltchinskii, V. and Yuan, M. (2008) Sparse recovery in large ensembles kernel machines. In *Proc. 21st A. Conf. Learning Theory*, pp. 229–238. Eastbourne: Omnipress.
- Ledoux, M. and Talagrand, M. (1991) *Probability in Banach Spaces: Isoperimetry and Processes*. New York: Springer.
- Lin, Y. and Zhang, H. H. (2006) Component selection and smoothing in multivariate nonparametric regression. *Ann. Statist.*, **34**, 2272–2297.
- Meier, L., van de Geer, S. and Bühlmann, P. (2008) High-dimensional additive modelling. (Available from arXiv.)
- Meinshausen, N. and Bühlmann, P. (2006) High dimensional graphs and variable selection with the lasso. *Ann. Statist.*, **34**, 1436–1462.
- Meinshausen, N. and Yu, B. (2006) Lasso-type recovery of sparse representations for high-dimensional data. *Technical Report 720*. Department of Statistics, University of California, Berkeley.
- Olshausen, B. A. and Field, D. J. (1996) Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, **381**, 607–609.
- Ravikumar, P., Liu, H., Lafferty, J. and Wasserman, L. (2008) Spam: sparse additive models. In *Advances in Neural Information Processing Systems*, vol. 20 (eds J. Platt, D. Koller, Y. Singer and S. Roweis), pp. 1201–1208. Cambridge: MIT Press.
- Tibshirani, R. (1996) Regression shrinkage and selection via the lasso. *J. R. Statist. Soc. B*, **58**, 267–288.
- van der Vaart, A. W. (1998) *Asymptotic Statistics*. Cambridge: Cambridge University Press.

- Wainwright, M. (2006) Sharp thresholds for high-dimensional and noisy recovery of sparsity. *Technical Report 709*. Department of Statistics, University of California, Berkeley.
- Wainwright, M. J., Ravikumar, P. and Lafferty, J. D. (2007) High-dimensional graphical model selection using l_1 -regularized logistic regression. In *Advances in Neural Information Processing Systems*, vol. 19 (eds B. Schölkopf, J. Platt and T. Hoffman), pp. 1465–1472. Cambridge: MIT Press.
- Wasserman, L. and Roeder, K. (2007) Multi-stage variable selection: screen and clean. Carnegie Mellon University, Pittsburgh. (Available from arXiv:0704.1139.)
- Yuan, M. (2007) Nonnegative garrote component selection in functional ANOVA models. *Proc. Artif. Intell. Statist.* (Available from www.stat.umn.edu/~aistat/proceedings.)
- Yuan, M. and Lin, Y. (2006) Model selection and estimation in regression with grouped variables. *J. R. Statist. Soc. B*, **68**, 49–67.
- Zhao, P. and Yu, B. (2007) On model selection consistency of lasso. *J. Mach. Learn. Res.*, **7**, 2541–2567.
- Zou, H. (2005) The adaptive lasso and its oracle properties. *J. Am. Statist. Ass.*, **101**, 1418–1429.