

SPADES AND MIXTURE MODELS

FLORENTINA BUNEA, ALEXANDRE B. TSYBAKOV,
MARTEN H. WEGKAMP, ADRIAN BARBU

ABSTRACT. This paper studies sparse density estimation via ℓ_1 penalization (SPADES). We focus on estimation in high-dimensional mixture models and nonparametric adaptive density estimation. We show, respectively, that SPADES can recover, with high probability, the unknown components of a mixture of probability densities and that it yields minimax adaptive density estimates. These results are based on a general sparsity oracle inequality that the SPADES estimates satisfy. We offer a data driven method for the choice of the tuning parameter used in the construction of SPADES. The method uses the generalized bisection method first introduced in [10]. The suggested procedure bypasses the need for a grid search and offers substantial computational savings. We complement our theoretical results with a simulation study that employs this method for approximations of one and two dimensional densities with mixtures. The numerical results strongly support our theoretical findings.

MSC2000 Subject classification: Primary 62G08, Secondary 62C20, 62G05, 62G20

Keywords and phrases: Adaptive estimation, aggregation, lasso, minimax risk, mixture models, consistent model selection, nonparametric density estimation, oracle inequalities, penalized least squares, sparsity, statistical learning.

Acknowledgements: The research of Bunea and Wegkamp is partially supported by the National Science Foundation grant DMS 0706829. The research of Tsybakov is partially supported by the grant ANR-06-BLAN-0194 and by the PASCAL Network of Excellence. Part of the research was done while the authors were visiting the Isaac Newton Institute for Mathematical Sciences (Statistical Theory and Methods for Complex, High-Dimensional Data Programme) at Cambridge University during Spring 2008.

1. INTRODUCTION

Let X_1, \dots, X_n be independent random variables with common unknown density f in \mathbb{R}^d . Let $\{f_1, \dots, f_M\}$ be a finite set of functions with $f_j \in L_2(\mathbb{R}^d), j = 1, \dots, M$, called a dictionary. We consider estimators of f that belong to the linear span of $\{f_1, \dots, f_M\}$. We will be particularly interested in the case where $M \gg n$. Denote by f_λ the linear combinations

$$f_\lambda(x) = \sum_{j=1}^M \lambda_j f_j(x), \quad \lambda = (\lambda_1, \dots, \lambda_M) \in \mathbb{R}^M.$$

Let us mention some examples where such estimates are of importance.

Date: December 29, 2009.

- *Estimation in sparse mixture models.* Assume that the density f can be represented as a finite mixture $f = f_{\lambda^*}$ where f_j are known probability densities and λ^* is a vector of mixture probabilities. The number M can be very large, much larger than the sample size n , but we believe that the representation is sparse, i.e., that very few coordinates of λ^* are non-zero, with indices corresponding to a set $I^* \subseteq \{1, \dots, M\}$. Our goal is to estimate the weight vector λ^* by a vector $\hat{\lambda}$ that adapts to this unknown sparsity and to identify I^* , with high probability.
- *Adaptive nonparametric density estimation.* Assume that the density f is a smooth function, and $\{f_1, \dots, f_M\}$ are the first M functions from a basis in $L_2(\mathbb{R}^d)$. If the basis is orthonormal, a natural idea is to estimate f by an orthogonal series estimator which has the form $f_{\tilde{\lambda}}$ with $\tilde{\lambda}$ having the coordinates $\tilde{\lambda}_j = n^{-1} \sum_{i=1}^n f_j(X_i)$. However, it is well known that such estimators are very sensitive to the choice of M , and a data-driven selection of M or thresholding is needed to achieve adaptivity (cf., e.g., [37, 27, 6]); moreover these methods have been applied with $M \leq n$. We would like to cover more general problems where the system $\{f_j\}$ is not necessarily orthonormal, even not necessarily a basis, M is not necessarily smaller than n , but an estimate of the form $f_{\tilde{\lambda}}$ still achieves, adaptively, the optimal rates of convergence.
- *Aggregation of density estimators.* Assume now that f_1, \dots, f_M are some preliminary estimators of f constructed from a training sample independent of (X_1, \dots, X_n) , and we would like to aggregate f_1, \dots, f_M . This means that we would like to construct a new estimator, the aggregate, which is approximately as good as the best among f_1, \dots, f_M or approximately as good as the best linear or convex combination of f_1, \dots, f_M . General notions of aggregation and optimal rates are introduced in [33, 40]. Aggregation of density estimators is discussed in [38, 36, 35] and more recently in [5] where one can find further references. The aggregates that we have in mind here are of the form $f_{\hat{\lambda}}$ with suitably chosen weights $\hat{\lambda} = \hat{\lambda}(X_1, \dots, X_n) \in \mathbb{R}^M$.

In this paper, we suggest a data-driven choice of $\hat{\lambda}$ that can be used in all the examples mentioned above and also more generally. We define $\hat{\lambda}$ as a minimizer of an ℓ_1 -penalized criterion, that we call SPADES (SPARse Density ESTimation). This method was introduced in [14]. The idea of ℓ_1 penalized estimation is widely used in the statistical literature, mainly in linear regression where it is usually referred to as the Lasso criterion [39, 16, 19, 24, 32]. For Gaussian sequence models or for regression with orthogonal design matrix the Lasso

is equivalent to soft thresholding [18, 30]. Model selection consistency of the Lasso type linear regression estimators is treated in many papers including [32, 47, 46, 48, 31]. Recently, ℓ_1 penalized methods have been extended to nonparametric regression with general fixed or random design [11, 12, 13, 4], as well as to some classification and other more general prediction type models [28, 29, 42, 8].

In this paper we show that ℓ_1 penalized techniques can also be successfully used in density estimation. In Section 2 we give the construction of the SPADES estimates and we show that they satisfy general oracle inequalities in Section 3. In the remainder of the paper we discuss the implications of these results for two particular problems, identification of mixture components and adaptive nonparametric density estimation. For the application of SPADES in aggregation problems we refer to [14].

Section 4 is devoted to mixture models. A vast amount of literature exists on estimation in mixture models, especially when the number of components is known; see e.g. [43] for examples involving the EM algorithm. The literature on determining the number of mixture components is still developing, and we will focus on this aspect here. Recent works on the selection of the number of components (mixture complexity) are [26, 2]. A consistent selection procedure specialized to Gaussian mixtures is suggested in [26]. The method of [26] relies on comparing a nonparametric kernel density estimator with the best parametric fit of various given mixture complexities. Nonparametric estimators based on the combinatorial density method (see [17]) are studied in [2, 3]. These can be applied to estimating consistently the number of mixture components, when the components have known functional form. Both [26, 2] can become computationally infeasible when M , the number of candidate components, is large. The method proposed here bridges this gap and guarantees correct identification of the mixture components with probability close to 1.

In Section 4 we begin by giving conditions under which the mixture weights can be estimated accurately, with probability close to 1. This is an intermediate result that allows us to obtain the main result of Section 4, correct identification of the mixture components. We show that in identifiable mixture models, if the mixture weights are above the noise level, then the components of the mixture can be recovered with probability larger than $1 - \varepsilon$, for any given small ε . Our results are non-asymptotic, they hold for any M and n . Since the emphasis here is on correct component selection, rather than optimal density estimation, the tuning sequence that accompanies the ℓ_1 penalty needs to be slightly larger than the one

used for good prediction. The same phenomenon has been noted for ℓ_1 penalized estimation in linear and generalized regression models, see, e.g., [8].

Section 5 uses the oracle inequalities of Section 3 to show that SPADES estimates adaptively achieve optimal rates of convergence (up to a logarithmic factor) simultaneously on a large scale of functional classes, such as Hölder, Sobolev or Besov classes, as well as on the classes of sparse densities, i.e., densities having only a finite, but unknown, number of non-zero wavelet coefficients.

Section 6.1 offers an algorithm for computing the SPADES. Our procedure is based on coordinate descent optimization, recently suggested by [20]. In Section 6.2 we use this algorithm together with a tuning parameter chosen in a data adaptive manner. This choice employs the generalized bisection method first introduced in [10], a computationally efficient method for constructing candidate tuning parameters without performing a grid search. The final tuning parameter is chosen from the list of computed candidates by using a 10-fold cross-validated dimension-regularized criterion. The combined procedure works very well in practice, and we present a simulation study in Section 6.3.

2. DEFINITION OF SPADES

Consider the $L_2(\mathbb{R}^d)$ norm

$$\|g\| = \left(\int_{\mathbb{R}^d} g^2(x) dx \right)^{1/2}$$

associated with the inner product

$$\langle g, h \rangle = \int_{\mathbb{R}^d} g(x)h(x) dx$$

for $g, h \in L_2(\mathbb{R}^d)$. Note that if the density f belongs to $L_2(\mathbb{R}^d)$ and X has the same distribution as X_i , we have, for any $g \in L_2$,

$$\langle g, f \rangle = \mathbb{E}g(X),$$

where the expectation is taken under f . Moreover

$$(2.1) \quad \|f - g\|^2 = \|f\|^2 + \|g\|^2 - 2\langle g, f \rangle = \|f\|^2 + \|g\|^2 - 2\mathbb{E}g(X).$$

In view of identity (2.1), minimizing $\|f_\lambda - f\|^2$ in λ is the same as minimizing

$$\gamma(\lambda) = -2\mathbb{E}f_\lambda(X) + \|f_\lambda\|^2.$$

The function $\gamma(\lambda)$ depends on f but can be approximated by its empirical counterpart

$$(2.2) \quad \widehat{\gamma}(\lambda) = -\frac{2}{n} \sum_{i=1}^n f_{\lambda}(X_i) + \|f_{\lambda}\|^2.$$

This motivates the use of $\widehat{\gamma} = \widehat{\gamma}(\lambda)$ as the empirical criterion, see, for instance, [6, 37, 44].

We define the penalty

$$(2.3) \quad \text{pen}(\lambda) = 2 \sum_{j=1}^M \omega_j |\lambda_j|$$

with weights ω_j to be specified later, and we propose the following data-driven choice of λ :

$$(2.4) \quad \begin{aligned} \widehat{\lambda} &= \arg \min_{\lambda \in \mathbb{R}^M} \{ \widehat{\gamma}(\lambda) + \text{pen}(\lambda) \} \\ &= \arg \min_{\lambda \in \mathbb{R}^M} \left\{ -\frac{2}{n} \sum_{i=1}^n f_{\lambda}(X_i) + \|f_{\lambda}\|^2 + 2 \sum_{j=1}^M \omega_j |\lambda_j| \right\}. \end{aligned}$$

Our estimator of density f that we will further call the *SPADES estimator* is defined by

$$f^{\spadesuit}(x) = f_{\widehat{\lambda}}(x), \quad \forall x \in \mathbb{R}^d.$$

It is easy to see that, for an orthonormal system $\{f_j\}$, the SPADES estimator coincides with the soft thresholding estimator whose components are of the form $\widehat{\lambda}_j = (1 - \omega_j/|\tilde{\lambda}_j|)_+ \tilde{\lambda}_j$ where $\tilde{\lambda}_j = n^{-1} \sum_{i=1}^n f_j(X_i)$ and $x_+ = \max(0, x)$. We see that in this case ω_j is the threshold for the j th component of a preliminary estimator $\tilde{\lambda} = (\tilde{\lambda}_1, \dots, \tilde{\lambda}_M)$.

The SPADES estimate can be easily computed by convex programming even if $M \gg n$. We present an algorithm in Section 6 below. SPADES retains the desirable theoretical properties of other density estimators, the computation of which may become problematic for $M \gg n$. We refer to [17] for a thorough overview on combinatorial methods in density estimation, to [41] for density estimation using support vector machines and to [6] for density estimates using penalties proportional to the dimension.

3. ORACLE INEQUALITIES FOR SPADES

3.1. Preliminaries. For any $\lambda \in \mathbb{R}^M$, let

$$J(\lambda) = \{j \in \{1, \dots, M\} : \lambda_j \neq 0\}$$

be the set of indices corresponding to non-zero components of λ and

$$M(\lambda) = |J(\lambda)| = \sum_{j=1}^M I\{\lambda_j \neq 0\}$$

its cardinality. Here $I\{\cdot\}$ denotes the indicator function. Furthermore, set

$$\sigma_j^2 = \text{Var}(f_j(X_1)), \quad L_j = \|f_j\|_\infty$$

for $1 \leq j \leq M$, where $\text{Var}(\zeta)$ denotes the variance of random variable ζ and $\|\cdot\|_\infty$ is the $L_\infty(\mathbb{R}^d)$ norm.

We will prove sparsity oracle inequalities for the estimator $\hat{\lambda} = \hat{\lambda}(\omega_1, \dots, \omega_M)$, provided the weights ω_j are chosen large enough. We first consider a simple choice:

$$(3.1) \quad \omega_j = 4L_j r(\delta/2)$$

where $0 < \delta < 1$ is a user-specified parameter and

$$(3.2) \quad r(\delta) = r(M, n, \delta) = \sqrt{\frac{\log(M/\delta)}{n}}.$$

The oracle inequalities that we prove below hold with a probability of at least $1 - \delta$ and are non-asymptotic: they are valid for all integers M and n . The first of these inequalities is established under a coherence condition on the ‘‘correlations’’

$$\rho_M(i, j) = \frac{\langle f_i, f_j \rangle}{\|f_i\| \|f_j\|}, \quad i, j = 1, \dots, M.$$

For $\lambda \in \mathbb{R}^M$, we define a local coherence number (called *maximal local coherence*) by

$$\rho(\lambda) = \max_{i \in J(\lambda)} \max_{j \neq i} |\rho_M(i, j)|,$$

and we also define

$$F(\lambda) = \max_{j \in J(\lambda)} \frac{\omega_j}{r(\delta/2) \|f_j\|} = \max_{j \in J(\lambda)} \frac{4L_j}{\|f_j\|}$$

and

$$G = \max_{1 \leq j \leq M} \frac{r(\delta/2) \|f_j\|}{\omega_j} = \max_{1 \leq j \leq M} \frac{\|f_j\|}{4L_j}.$$

3.2. Main results.

Theorem 1. *Assume that $L_j < \infty$ for $1 \leq j \leq M$. Then with probability at least $1 - \delta$ for all $\lambda \in \mathbb{R}^M$ that satisfy*

$$(3.3) \quad 16GF(\lambda)\rho(\lambda)M(\lambda) \leq 1$$

and all $\alpha > 1$ and we have the following oracle inequality:

$$\|f^\spadesuit - f\|^2 + \frac{\alpha}{2(\alpha-1)} \sum_{j=1}^M \omega_j |\hat{\lambda}_j - \lambda_j| \leq \frac{\alpha+1}{\alpha-1} \|f_\lambda - f\|^2 + \frac{8\alpha^2}{\alpha-1} F^2(\lambda) r^2(\delta/2) M(\lambda).$$

Note that only a condition on the local coherence (3.3) is required to obtain the result of Theorem 1. However, even this condition can be too strong, because the bound on “correlations” should be *uniform* over $j \in J(\lambda), i \neq j$, cf. the definition of $\rho(\lambda)$. For example, this excludes the cases where the “correlations” can be relatively large for a small number of pairs (i, j) and almost zero for otherwise. To account for this situation, we suggest below another version of Theorem 1. Instead of maximal local coherence, we introduce *cumulative local coherence* defined by

$$\rho_*(\lambda) = \sum_{i \in J(\lambda)} \sum_{j > i} |\rho_M(i, j)|.$$

Theorem 2. *Assume that $L_j < \infty$ for $1 \leq j \leq M$. Then with probability at least $1 - \delta$ for all $\lambda \in \mathbb{R}^M$ that satisfy*

$$(3.4) \quad 16F(\lambda)G\rho_*(\lambda)\sqrt{M(\lambda)} \leq 1$$

and all $\alpha > 1$ we have the following oracle inequality:

$$\|f^\spadesuit - f\|^2 + \frac{\alpha}{2(\alpha - 1)} \sum_{j=1}^M \omega_j |\hat{\lambda}_j - \lambda_j| \leq \frac{\alpha + 1}{\alpha - 1} \|f_\lambda - f\|^2 + \frac{8\alpha^2}{\alpha - 1} F^2(\lambda) r^2 (\delta/2) M(\lambda).$$

Theorem 2 is useful when we deal with sparse Gram matrices $\Psi_M = (\langle f_i, f_j \rangle)_{1 \leq i, j \leq M}$ that have only a small number N of non-zero off-diagonal entries. This number will be called a *sparsity index* of matrix Ψ_M , and is defined as

$$N = |\{(i, j) : i, j \in \{1, \dots, M\}, i > j \text{ and } \psi_M(i, j) \neq 0\}|,$$

where $\psi_M(i, j)$ is the (i, j) th entry of Ψ_M and $|A|$ denotes the cardinality of a set A . Clearly, $N < M(M + 1)/2$. We therefore obtain the following immediate corollary of Theorem 2.

Corollary 1. *Let Ψ_M be a Gram matrix with sparsity index N . Then the assertion of Theorem 2 holds if we replace there (3.4) by the condition*

$$(3.5) \quad 16F(\lambda)N\sqrt{M(\lambda)} \leq 1.$$

We finally give an oracle inequality, which is valid under the assumption that the Gram matrix Ψ_M is positive definite. It is simpler to use than the above results when the dictionary is orthonormal or forms a frame. Note that the coherence assumptions considered above do not necessarily imply the positive definiteness of Ψ_M . Vice versa, the positive definiteness of Ψ_M does not imply these assumptions.

Theorem 3. *Assume that $L_j < \infty$ for $1 \leq j \leq M$ and that the Gram matrix Ψ_M is positive definite with minimal eigenvalue larger than or equal to $\kappa_M > 0$. Then, with probability at least $1 - \delta$, for all $\alpha > 1$ and all $\lambda \in \mathbb{R}^M$ we have*

$$(3.6) \quad \|f^\spadesuit - f\|^2 + \frac{\alpha}{\alpha - 1} \sum_{j=1}^M \omega_j |\hat{\lambda}_j - \lambda_j| \leq \frac{\alpha + 1}{\alpha - 1} \|f_\lambda - f\|^2 + \left(\frac{8\alpha^2}{\alpha - 1} \right) \frac{G(\lambda)}{n \kappa_M},$$

where

$$G(\lambda) \triangleq \sum_{j \in J(\lambda)} \omega_j^2 = \frac{16 \log(2M/\delta)}{n} \sum_{j \in J(\lambda)} L_j^2.$$

We can consider some other choices for ω_j without affecting the previous results. For instance,

$$(3.7) \quad \omega_j = 2\sqrt{2}\sigma_j r(\delta/2) + \frac{8}{3}L_j r^2(\delta/2)$$

or

$$(3.8) \quad \omega_j = 2\sqrt{2}T_j r(\delta/2) + \frac{8}{3}L_j r^2(\delta/2)$$

with

$$T_j^2 = \frac{2}{n} \sum_{i=1}^n f_j^2(X_i) + 2L_j^2 r^2(\delta/2).$$

yield the same conclusions. These modifications of (3.1) prove useful, for example, for situations where f_j are wavelet basis functions, cf. Section 5. The choice (3.8) of ω_j has an advantage of being completely data-driven.

Theorem 4. *Theorems 1–3 and Corollary 1 hold with the choices (3.7) or (3.8) for the weights ω_j without changing the assertions. They also remain valid if we replace these ω_j by any ω'_j such that $\omega'_j > \omega_j$.*

If ω_j is chosen as in (3.8), our bounds on the risk of SPADES estimator involve the random variables $(1/n) \sum_{i=1}^n f_j^2(X_i)$. These can be replaced in the bounds by deterministic values using the following lemma.

Lemma 1. *Assume that $L_j < \infty$ for $j = 1, \dots, M$. Then*

$$(3.9) \quad \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n f_j^2(X_i) \leq 2\mathbb{E}f_j^2(X_1) + \frac{4}{3}L_j^2 r^2(\delta/2), \forall j = 1, \dots, M \right) \geq 1 - \delta/2.$$

From Theorem 4 and Lemma 1 we find that, for the choice of ω_j as in (3.8), the oracle inequalities of Theorems 1–3 and Corollary 1 remain valid with probability at least $1 - 3\delta/2$ if we replace the ω_j in these inequalities by the expressions $2\sqrt{2}\tilde{T}_j r(\delta/2) + (8/3)L_j r^2(\delta/2)$ where $\tilde{T}_j = \left(2\mathbb{E}f_j^2(X_1) + (4/3)L_j^2 r^2(\delta/2)\right)^{1/2}$.

3.3. Proofs. We first prove the following preliminary lemma. Define the random variables

$$V_j = \frac{1}{n} \sum_{i=1}^n \{f_j(X_i) - \mathbb{E}f_j(X_i)\}$$

and the event

$$(3.10) \quad A = \bigcap_{j=1}^M \{2|V_j| \leq \omega_j\}.$$

Lemma 2. *Assume that $L_j < \infty$ for $j = 1, \dots, M$. Then for all $\lambda \in \mathbb{R}^M$ we have that, on the event A ,*

$$(3.11) \quad \|f^\spadesuit - f\|^2 + \sum_{j=1}^M \omega_j |\hat{\lambda}_j - \lambda_j| \leq \|f_\lambda - f\|^2 + 4 \sum_{j \in J(\lambda)} \omega_j |\hat{\lambda}_j - \lambda_j|.$$

Proof. By the definition of $\hat{\lambda}$,

$$-\frac{2}{n} \sum_{i=1}^n f_{\hat{\lambda}}(X_i) + \|f_{\hat{\lambda}}\|^2 + 2 \sum_{j=1}^M \omega_j |\hat{\lambda}_j| \leq -\frac{2}{n} \sum_{i=1}^n f_\lambda(X_i) + \|f_\lambda\|^2 + 2 \sum_{j=1}^M \omega_j |\lambda_j|$$

for all $\lambda \in \mathbb{R}^M$. We rewrite this inequality as

$$\begin{aligned} \|f^\spadesuit - f\|^2 &\leq \|f_\lambda - f\|^2 - 2 \langle f, f^\spadesuit - f_\lambda \rangle + \frac{2}{n} \sum_{i=1}^n (f^\spadesuit - f_\lambda)(X_i) + 2 \sum_{j=1}^M \omega_j |\lambda_j| - 2 \sum_{j=1}^M \omega_j |\hat{\lambda}_j| \\ &= \|f_\lambda - f\|^2 + 2 \sum_{j=1}^M \left(\frac{1}{n} \sum_{i=1}^n f_j(X_i) - \mathbb{E}f_j(X_i) \right) (\hat{\lambda}_j - \lambda_j) \\ &\quad + 2 \sum_{j=1}^M \omega_j |\lambda_j| - 2 \sum_{j=1}^M \omega_j |\hat{\lambda}_j|. \end{aligned}$$

Then, on the event A ,

$$\|f^\spadesuit - f\|^2 \leq \|f_\lambda - f\|^2 + \sum_{j=1}^M \omega_j |\hat{\lambda}_j - \lambda_j| + 2 \sum_{j=1}^M \omega_j |\lambda_j| - 2 \sum_{j=1}^M \omega_j |\hat{\lambda}_j|.$$

Add $\sum_j \omega_j |\widehat{\lambda}_j - \lambda_j|$ to both sides of the inequality to obtain

$$\begin{aligned}
& \|f^\spadesuit - f\|^2 + \sum_{j=1}^M \omega_j |\widehat{\lambda}_j - \lambda_j| \\
& \leq \|f_\lambda - f\|^2 + 2 \sum_{j=1}^M \omega_j |\widehat{\lambda}_j - \lambda_j| + 2 \sum_{j=1}^M \omega_j |\lambda_j| - 2 \sum_{j=1}^M \omega_j |\widehat{\lambda}_j| \\
& \leq \|f_\lambda - f\|^2 + 2 \sum_{j \in J(\lambda)} \omega_j |\widehat{\lambda}_j - \lambda_j| + 2 \sum_{j=1}^M \omega_j |\lambda_j| - 2 \sum_{j \in J(\lambda)} \omega_j |\widehat{\lambda}_j| \\
& \leq \|f_\lambda - f\|^2 + 4 \sum_{j \in J(\lambda)} \omega_j |\widehat{\lambda}_j - \lambda_j|
\end{aligned}$$

where we used that $\lambda_j = 0$ for $j \notin J(\lambda)$ and the triangle inequality. \square

For the choice (3.1) for ω_j , we find by Hoeffding's inequality for sums of independent random variables $\zeta_{ij} = f_j(X_i) - \mathbb{E}f_j(X_i)$ with $|\zeta_{ij}| \leq 2L_j$ that

$$\mathbb{P}(A) \leq \sum_{j=1}^M \mathbb{P}\{2|V_j| > \omega_j\} \leq 2 \sum_{j=1}^M \exp\left(-\frac{2n\omega_j^2/4}{8L_j^2}\right) = \delta.$$

Proof of Theorem 1. In view of Lemma 2, we need to bound $\sum_{j \in J(\lambda)} \omega_j |\widehat{\lambda}_j - \lambda_j|$. Set

$$u_j = \widehat{\lambda}_j - \lambda_j, \quad U(\lambda) = \sum_{j \in J(\lambda)} |u_j| \|f_j\|, \quad U = \sum_{j=1}^M |u_j| \|f_j\| \quad r = r(\delta/2).$$

Then, by the definition of $F(\lambda)$,

$$\sum_{j \in J(\lambda)} \omega_j |\widehat{\lambda}_j - \lambda_j| \leq rF(\lambda)U(\lambda).$$

Since

$$\sum_{i,j \notin J(\lambda)} \langle f_i, f_j \rangle u_i u_j \geq 0,$$

we obtain

$$\begin{aligned}
 (3.12) \quad \sum_{j \in J(\lambda)} u_j^2 \|f_j\|^2 &= \|f^\spadesuit - f_\lambda\|^2 - \sum_{i, j \notin J(\lambda)} u_i u_j \langle f_i, f_j \rangle \\
 &\quad - 2 \sum_{i \notin J(\lambda)} \sum_{j \in J(\lambda)} u_i u_j \langle f_i, f_j \rangle - \sum_{i, j \in J(\lambda), i \neq j} u_i u_j \langle f_i, f_j \rangle \\
 &\leq \|f^\spadesuit - f_\lambda\|^2 + 2\rho(\lambda) \sum_{i \notin J(\lambda)} |u_i| \|f_i\| \sum_{j \in J(\lambda)} |u_j| \|f_j\| \\
 &\quad + \rho(\lambda) \sum_{i, j \in J(\lambda)} |u_i| |u_j| \|f_i\| \|f_j\| \\
 &= \|f^\spadesuit - f_\lambda\|^2 + 2\rho(\lambda) U(\lambda) U - \rho(\lambda) U^2(\lambda).
 \end{aligned}$$

The left-hand side can be bounded by $\sum_{j \in J(\lambda)} u_j^2 \|f_j\|^2 \geq U^2(\lambda)/M(\lambda)$ using the Cauchy-Schwarz inequality, and we obtain that

$$U^2(\lambda) \leq \|f^\spadesuit - f_\lambda\|^2 M(\lambda) + 2\rho(\lambda) M(\lambda) U(\lambda) U,$$

which immediately implies

$$(3.13) \quad U(\lambda) \leq 2\rho(\lambda) M(\lambda) U + \sqrt{M(\lambda)} \|f^\spadesuit - f_\lambda\|.$$

Hence, by Lemma 2, we have with probability at least $1 - \delta$,

$$\begin{aligned}
 &\|f^\spadesuit - f\|^2 + \sum_{j=1}^M \omega_j |\hat{\lambda}_j - \lambda_j| \\
 &\leq \|f_\lambda - f\|^2 + 4 \sum_{j \in J(\lambda)} \omega_j |\hat{\lambda}_j - \lambda_j| \\
 &\leq \|f_\lambda - f\|^2 + 4rF(\lambda) U(\lambda) \\
 &\leq \|f_\lambda - f\|^2 + 4rF(\lambda) \left\{ 2\rho(\lambda) M(\lambda) U + \sqrt{M(\lambda)} \|f^\spadesuit - f_\lambda\| \right\} \\
 &\leq \|f_\lambda - f\|^2 + 8F(\lambda) \rho(\lambda) M(\lambda) G \sum_{j=1}^M \omega_j |\hat{\lambda}_j - \lambda_j| + 4rF(\lambda) \sqrt{M(\lambda)} \|f^\spadesuit - f_\lambda\|.
 \end{aligned}$$

For all $\lambda \in \mathbb{R}^M$ that satisfy relation (3.3), we find that with probability exceeding $1 - \delta$,

$$\begin{aligned}
 \|f^\spadesuit - f\|^2 + \frac{1}{2} \sum_{j=1}^M \omega_j |\hat{\lambda}_j - \lambda_j| &\leq \|f_\lambda - f\|^2 + 4rF(\lambda) \sqrt{M(\lambda)} \|f^\spadesuit - f_\lambda\| \\
 &\leq \|f_\lambda - f\|^2 + 2 \left\{ 2rF(\lambda) \sqrt{M(\lambda)} \right\} \|f^\spadesuit - f\| \\
 &\quad + 2 \left\{ 2rF(\lambda) \sqrt{M(\lambda)} \right\} \|f_\lambda - f\|.
 \end{aligned}$$

After applying the inequality $2xy \leq x^2/\alpha + \alpha y^2$ ($x, y \in \mathbb{R}$, $\alpha > 0$) for each of the last two summands, we easily find the claim. \square

Proof of Theorem 2. The proof is similar to that of Theorem 1. With

$$U_*(\lambda) = \sqrt{\sum_{j \in J(\lambda)} u_j^2 \|f_j\|^2}$$

we obtain now the following analogue of (3.12):

$$\begin{aligned} U_*^2(\lambda) &\leq \|f^\spadesuit - f_\lambda\|^2 + 2\rho_*(\lambda) \max_{i \in J(\lambda), j > i} |u_i| \|f_i\| |u_j| \|f_j\| \\ &\leq \|f^\spadesuit - f_\lambda\|^2 + 2\rho_*(\lambda) U_*(\lambda) \sum_{j=1}^M |u_j| \|f_j\| \\ &= \|f^\spadesuit - f_\lambda\|^2 + 2\rho_*(\lambda) U_*(\lambda) U. \end{aligned}$$

Hence, as in the proof of Theorem 1, we have

$$U_*(\lambda) \leq 2\rho_*(\lambda)U + \|f^\spadesuit - f_\lambda\|,$$

and using the inequality $U_*(\lambda) \geq U(\lambda)/\sqrt{M(\lambda)}$ we find

$$(3.14) \quad U(\lambda) \leq 2\rho_*(\lambda)\sqrt{M(\lambda)}U + \sqrt{M(\lambda)}\|f^\spadesuit - f_\lambda\|.$$

Note that (3.14) differs from (3.13) only in the fact that the factor $2\rho(\lambda)M(\lambda)$ on the right hand side is now replaced by $2\rho_*(\lambda)\sqrt{M(\lambda)}$. Up to this modification, the rest of the proof is identical to that of Theorem 1. \square

Proof of Theorem 3. By the assumption on Ψ_M we have

$$\|f_\lambda\|^2 = \sum_{1 \leq i, j \leq M} \lambda_i \lambda_j \int_{\mathbb{R}^d} f_i(x) f_j(x) dx \geq \kappa_M \sum_{j \in J(\lambda)} \lambda_j^2.$$

By the Cauchy-Schwarz inequality, we find

$$\begin{aligned} 4 \sum_{j \in J(\lambda)} \omega_j |\widehat{\lambda}_j - \lambda_j| &\leq 4 \sqrt{\sum_{j \in J(\lambda)} \omega_j^2} \sqrt{\sum_{j \in J(\lambda)} |\widehat{\lambda}_j - \lambda_j|^2} \\ &\leq 4 \left(\frac{\sum_{j \in J(\lambda)} \omega_j^2}{n\kappa_M} \right)^{1/2} \|f^\spadesuit - f_\lambda\|. \end{aligned}$$

Combination with Lemma 2 yields that, with probability at least $1 - \delta$,

$$\begin{aligned} (3.15) \quad \|f^\spadesuit - f\|^2 + \sum_{j=1}^M \omega_j |\widehat{\lambda}_j - \lambda_j| &\leq \|f_\lambda - f\|^2 + 4 \left(\frac{\sum_{j \in J(\lambda)} \omega_j^2}{n\kappa_M} \right)^{1/2} \|f^\spadesuit - f_\lambda\| \\ &\leq \|f_\lambda - f\|^2 + b \left(\|f^\spadesuit - f\| + \|f_\lambda - f\| \right) \end{aligned}$$

where $b = 4\sqrt{\sum_{j \in J(\lambda)} \omega_j^2 / \sqrt{n\kappa_M}}$. Applying the inequality $2xy \leq x^2/\alpha + \alpha y^2$ ($x, y \in \mathbb{R}$, $\alpha > 0$) for each of the last two summands in (3.15) we get the result. \square

Proof of Theorem 4. Write $\bar{\omega}_j = 2\sqrt{2}\sigma_j r(\delta/2) + (8/3)L_j r^2(\delta/2)$ for the choice of ω_j in (3.7). Using Bernstein's exponential inequality for sums of independent random variables $\zeta_{ij} = f_j(X_i) - \mathbb{E}f_j(X_i)$ with $|\zeta_{ij}| \leq 2L_j$, we obtain that

$$\begin{aligned}
 (3.16) \quad \mathbb{P}(A^c) &= \mathbb{P}\left(\bigcup_{j=1}^M \{2|V_j| > \bar{\omega}_j\}\right) \leq \sum_{j=1}^M \mathbb{P}\{2|V_j| > \bar{\omega}_j\} \\
 &\leq \sum_{j=1}^M \exp\left(-\frac{n\bar{\omega}_j^2/4}{2\text{Var}(f_j(X_1)) + 2L_j\bar{\omega}_j/3}\right) \\
 &\leq M \exp(-nr^2(\delta/2)) = \delta/2.
 \end{aligned}$$

Let now ω_j be defined by (3.8). Then, using (3.16), we can write

$$\begin{aligned}
 (3.17) \quad \mathbb{P}(A^c) &= \mathbb{P}\left(\bigcup_{j=1}^M \{2|V_j| > \omega_j\}\right) \\
 &\leq \sum_{j=1}^M \mathbb{P}\{2|V_j| > \bar{\omega}_j\} + \sum_{j=1}^M \mathbb{P}\{\bar{\omega}_j > \omega_j\} \\
 &\leq \delta/2 + \sum_{j=1}^M \mathbb{P}\{\bar{\omega}_j > \omega_j\}.
 \end{aligned}$$

Define

$$t_j = 2 \frac{\mathbb{E}f_j^4(X_1) \log(2M/\delta)}{\mathbb{E}f_j^2(X_1) n}$$

and note that

$$\frac{2}{n} \sum_{i=1}^n f_j^2(X_i) + t_j \leq T_j^2.$$

Then

$$\begin{aligned}
\mathbb{P}\{\bar{\omega}_j > \omega_j\} &= \mathbb{P}\{\text{Var}(f_j(X_1)) > T_j^2\} \\
&\leq \mathbb{P}\{\mathbb{E}f_j^2(X_1) > \frac{2}{n} \sum_{i=1}^n f_j^2(X_i) + t_j\} \\
&\leq \exp\left(-\frac{n\{\mathbb{E}f_j^2(X_1) + t_j\}^2}{8\mathbb{E}f_j^4(X_1)}\right) \\
&\quad \text{using Proposition 2.6 in [45]} \\
&\leq \exp\left(-\frac{nt_j\mathbb{E}f_j^2(X_1)}{2\mathbb{E}f_j^4(X_1)}\right) \\
&\quad \text{since } (x+y)^2 \geq 4xy
\end{aligned}$$

which is less than $\delta/(2M)$. Plugging this in (3.17) concludes the proof. \square

Proof of Lemma 1. Using Bernstein's exponential inequality for sums of independent random variables $f_j^2(X_i) - \mathbb{E}f_j^2(X_i)$ and the fact that $\mathbb{E}f_j^4(X_1) \leq L_j^2\mathbb{E}f_j^2(X_1)$ we find

$$\begin{aligned}
&\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n f_j^2(X_i) \geq 2\mathbb{E}f_j^2(X_1) + \frac{4}{3}L_j^2r^2(\delta/2)\right) \\
&= \mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n f_j^2(X_i) - \mathbb{E}f_j^2(X_1) \geq \mathbb{E}f_j^2(X_1) + \frac{4}{3}L_j^2r^2(\delta/2)\right) \\
&\leq \exp\left(-\frac{n(\mathbb{E}f_j^2(X_1) + \frac{4}{3}L_j^2r^2(\delta/2))^2}{2\mathbb{E}f_j^4(X_1) + \frac{4}{3}L_j^2\{\mathbb{E}f_j^2(X_1) + \frac{4}{3}L_j^2r^2(\delta/2)\}}\right) \\
&\leq \exp(-nr^2(\delta/2)) = \frac{\delta}{2M},
\end{aligned}$$

which implies the lemma. \square

4. SPARSE ESTIMATION IN MIXTURE MODELS

In this section we assume that the true density f can be represented as a finite mixture

$$f(x) = \sum_{j \in I^*} \bar{\lambda}_j p_j(x),$$

where $I^* \subseteq \{1, \dots, M\}$ is unknown, p_j are known probability densities and $\bar{\lambda}_j > 0$ for all $j \in I^*$. We focus in this section on model selection, i.e., on the correct identification of the set I^* . It will be convenient for us to normalize the densities p_j by their L_2 norms and to

write the model in the form

$$f(x) = \sum_{j \in I^*} \lambda_j^* f_j(x),$$

where $I^* \subseteq \{1, \dots, M\}$ is unknown, $f_j = p_j/\|p_j\|$ are known functions and $\lambda_j^* > 0$ for all $j \in I^*$. We set $\lambda^* = (\lambda_1^*, \dots, \lambda_M^*)$ where $\lambda_j^* = 0, j \notin I^*$.

For clarity of exposition we consider a simplified version of the general set-up introduced above. We compute the estimates of λ^* via (2.4), with weights defined by (cf. (3.1)):

$$\omega_j = 4Lr, \text{ for all } j,$$

where $r > 0$ is a constant that we specify below, and for clarity of exposition we replaced all $L_j = \|f_j\|_\infty$ by an upper bound L on $\max_{1 \leq j \leq M} L_j$. Recall that, by construction, $\|f_j\| = 1$ for all j . Under these assumptions condition (3.3) takes the form

$$(4.1) \quad \rho(\lambda) \leq \frac{1}{16M(\lambda)}.$$

We state (4.1) for the true vector λ^* in the following form.

Condition (A).

$$\rho^* \leq \frac{1}{16k^*}$$

where $k^* = |I^*| = M(\lambda^*)$ and $\rho^* = \rho(\lambda^*)$.

Similar conditions are quite standard in the literature on sparse regression estimation and compressed sensing, cf., e.g., [19, 47, 11, 13, 4, 8]. The difference is that those papers use the empirical version of the correlation ρ^* and the numerical constant in the inequality is, in general, different from 1/16. Note that Condition (A) is quite intuitive. Indeed, the sparsity index k^* can be viewed as the effective dimension of the problem. When k^* increases the problem becomes harder, so that we need stronger conditions (smaller correlations ρ^*) in order to obtain our results. The interesting case that we have in mind is when the effective dimension k^* is small, i.e., the model is sparse.

The results of Section 3 are valid for any r larger or equal to $r(\delta/2) = \{\log(2M/\delta)/n\}^{1/2}$. They give bounds on the predictive performance of SPADES. As noted in, e.g., [8], for ℓ_1 -penalized model selection in regression, the tuning sequence ω_j required for correct selection is typically larger than the one that yields good prediction. We show below that the same is true for selecting the components of a mixture of densities. Specifically, in this section we will take the value

$$(4.2) \quad r = r(M, n, \delta/(2M)) = \sqrt{\frac{\log(2M^2/\delta)}{n}}.$$

We will use the following corollary of Theorem 1, obtained for $\alpha = \sqrt{2}$.

Corollary 2. *Assume that Condition (A) holds. Then with probability at least $1 - \delta/M$ we have*

$$(4.3) \quad \sum_{j=1}^M |\hat{\lambda}_j - \lambda_j^*| \leq \frac{4\sqrt{2}}{L} k^* \sqrt{\frac{\log(2M^2/\delta)}{n}}.$$

Inequality (4.3) guarantees that the estimate $\hat{\lambda}$ is close to the true λ^* in ℓ_1 norm, if the number of mixture components k^* is substantially smaller than \sqrt{n} . We regard this as an intermediate step for the next result that deals with the identification of I^* .

4.1. Correct identification of the mixture components. We now show that I^* can be identified with probability close to 1 by our procedure. Let $\hat{I} = J(\hat{\lambda})$ be the set of indices of the non-zero components of $\hat{\lambda}$ given by (2.4). In what follows we investigate when $P(\hat{I} = I^*) \geq 1 - \varepsilon$ for a given $0 < \varepsilon < 1$. Our results are non-asymptotic, they hold for any fixed M and n .

We need two conditions to ensure that correct recovery of I^* is possible. The first one is the identifiability of the model, as quantified by *Condition (A)* above. The second condition requires that the weights of the mixture are above the noise level, quantified by r . We state it as follows.

Condition (B).

$$\min_{j \in I^*} |\lambda_j^*| > 4(\sqrt{2} + 1)rL$$

where $L = \max(1/\sqrt{3}, \max_{1 \leq j \leq M} L_j)$ and r is given in (4.2).

Theorem 5. *Let $0 < \delta < 1/2$ be a given number. Assume that Conditions (A) and (B) hold. Then $\mathbb{P}(\hat{I} = I^*) \geq 1 - 2\delta(1 + 1/M)$.*

Remark. Since all λ_j^* are nonnegative, it seems reasonable to restrict the minimization in (2.4) to λ with nonnegative components. Inspection of the proofs shows that all the results of this section remain valid for such a modified estimator. However, in practice the non-negativity issue is not so important. Indeed, the estimators of the weights are quite close to the true values and turn out to be positive for positive λ_j^* . For example, this was the case in our simulations discussed in Section 6 below. On the other hand, adding the non-negativity constraint in (2.4) introduces some extra burden on the numerical algorithm. More generally, it is trivial to note that the results of this and previous sections extend verbatim

to the setting where $\lambda \in \Lambda$ with Λ being any subset of \mathbb{R}^M . Then the minimization in (2.4) should be performed on Λ , in the theorems of Section 3 we should replace $\lambda \in \mathbb{R}^M$ by $\lambda \in \Lambda$ and in this section λ^* should be supposed to belong to Λ .

Proof of Theorem 5. We begin by noticing that

$$\mathbb{P}(\hat{I} \neq I^*) \leq \mathbb{P}(I^* \not\subseteq \hat{I}) + \mathbb{P}(\hat{I} \not\subseteq I^*),$$

and we control each of the probabilities on the right hand side separately.

Control of $\mathbb{P}(I^* \not\subseteq \hat{I})$. By the definitions of the sets \hat{I} and I^* we have

$$\begin{aligned} \mathbb{P}(I^* \not\subseteq \hat{I}) &\leq \mathbb{P}(\hat{\lambda}_k = 0 \text{ for some } k \in I^*) \\ &\leq k^* \max_{k \in I^*} \mathbb{P}(\hat{\lambda}_k = 0). \end{aligned}$$

We control the last probability by using the characterization (6.1) of $\hat{\lambda}$ given in Lemma 3 of the Appendix. We also recall that $\mathbb{E}f_k(X_1) = \sum_{j \in I^*} \lambda_j^* \langle f_k, f_j \rangle = \sum_{j=1}^M \lambda_j^* \langle f_k, f_j \rangle$, since we assumed that the density of X_1 is the mixture $f^* = \sum_{j \in I^*} \lambda_j^* f_j$. We therefore obtain, for $k \in I^*$,

$$\begin{aligned} \mathbb{P}(\hat{\lambda}_k = 0) &= \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n f_k(X_i) - \sum_{j=1}^M \hat{\lambda}_j \langle f_j, f_k \rangle\right| \leq 4rL; \hat{\lambda}_k = 0\right) \\ &= \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n f_k(X_i) - \mathbb{E}f_k(X_1) - \sum_{j=1}^M (\hat{\lambda}_j - \lambda_j^*) \langle f_j, f_k \rangle\right| \leq 4rL; \hat{\lambda}_k = 0\right) \\ &\leq \mathbb{P}\left(\left|\lambda_k^* \|f_k\|^2 + \frac{1}{n} \sum_{i=1}^n f_k(X_i) - \mathbb{E}f_k(X_1) - \sum_{j \neq k} (\hat{\lambda}_j - \lambda_j^*) \langle f_j, f_k \rangle\right| \leq 4rL\right) \\ (4.4) \quad &\leq \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n f_k(X_i) - \mathbb{E}f_k(X_1)\right| \geq \frac{|\lambda_k^*| \|f_k\|^2}{2} - 2rL\right) \\ (4.5) \quad &+ \mathbb{P}\left(\left|\sum_{j \neq k} (\hat{\lambda}_j - \lambda_j^*) \langle f_j, f_k \rangle\right| \geq \frac{|\lambda_k^*| \|f_k\|^2}{2} - 2rL\right). \end{aligned}$$

To bound (4.4) we use Hoeffding's inequality, as in the course of Lemma 2. We first recall that $\|f_k\| = 1$ for all k and that, by *Condition (B)*, $\min_{k \in I^*} |\lambda_k^*| \geq 4(\sqrt{2} + 1)Lr$, with

$r = r(\delta/(2M)) = \{\log(2M^2/\delta)/n\}^{1/2}$. Therefore

$$(4.6) \quad \begin{aligned} & \mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n f_k(X_i) - \mathbb{E} f_k(X_1) \right| \geq \frac{|\lambda_k^*|}{2} - 2rL \right) \\ & \leq \mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n f_k(X_i) - \mathbb{E} f_k(X_1) \right| \geq 2\sqrt{2}rL \right) \leq \frac{\delta}{M^2}. \end{aligned}$$

To bound (4.5) notice that, by Conditions (A) and (B),

$$\begin{aligned} & \mathbb{P} \left(\left| \sum_{j \neq k} (\hat{\lambda}_j - \lambda_j^*) \langle f_j, f_k \rangle \right| \geq \frac{|\lambda_k^*|}{2} - 2rL \right) \\ & \leq \mathbb{P} \left(\sum_{j=1}^M |\hat{\lambda}_j - \lambda_j^*| \geq 32\sqrt{2}rLk^* \right) \\ & \leq \mathbb{P} \left(\sum_{j=1}^M |\hat{\lambda}_j - \lambda_j^*| \geq \frac{4\sqrt{2}rk^*}{L} \right) \leq \frac{\delta}{M}, \end{aligned}$$

where the penultimate inequality holds since, by definition, $L^2 \geq 1/3$ and the last inequality holds by Corollary 4.

Combining the above results we obtain

$$\mathbb{P}(I^* \not\subseteq \hat{I}) \leq k^* \frac{\delta}{M^2} + k^* \frac{\delta}{M} \leq \frac{\delta}{M} + \delta.$$

Control of $\mathbb{P}(\hat{I} \not\subseteq I^*)$. Let

$$(4.7) \quad h(\mu) = -\frac{2}{n} \sum_{i=1}^n \sum_{j \in I^*} \mu_j f_j(X_i) + \left\| \sum_{j \in I^*} \mu_j f_j \right\|^2 + 8rL \sum_{j \in I^*} |\mu_j|.$$

Let

$$(4.8) \quad \tilde{\mu} = \arg \min_{\mu \in \mathbb{R}^{k^*}} h(\mu).$$

Consider the random event

$$(4.9) \quad \mathcal{B} = \bigcap_{k \notin I^*} \left\{ \left| -\frac{1}{n} \sum_{i=1}^n f_k(X_i) + \sum_{j \in I^*} \tilde{\mu}_j \langle f_j, f_k \rangle \right| \leq 4Lr \right\}.$$

Let $\bar{\mu} \in \mathbb{R}^M$ be the vector that has the components of $\tilde{\mu}$ given by (4.8) in positions corresponding to the index set I^* and zero components elsewhere. By the first part of Lemma 3 in the Appendix we have that $\bar{\mu} \in \mathbb{R}^M$ is a solution of (2.4) on the event \mathcal{B} . Recall that $\hat{\lambda}$ is

a also solution of (2.4). By the definition of the set \hat{I} we have that $\hat{\lambda}_k \neq 0$ for $k \in \hat{I}$. By construction, $\tilde{\mu}_k \neq 0$ for some subset $S \subseteq I^*$. By the second part of Lemma 3 in the Appendix, any two solutions have non-zero elements in the same positions. Therefore $\hat{I} = S \subseteq I^*$ on \mathcal{B} . Thus,

$$\begin{aligned}
 (4.10) \quad & \mathbb{P}(\hat{I} \not\subseteq I^*) \leq \mathbb{P}(\mathcal{B}^c) \\
 & \leq \sum_{k \notin I^*} \mathbb{P} \left\{ \left| -\frac{1}{n} \sum_{i=1}^n f_k(X_i) + \sum_{j \in I^*} \tilde{\mu}_j \langle f_j, f_k \rangle \right| \geq 4rL \right\} \\
 & \leq \sum_{k \notin I^*} \mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n f_k(X_i) - E f_k(X_1) \right| \geq 2\sqrt{2}rL \right) \\
 & \quad + \sum_{k \notin I^*} \mathbb{P} \left(\sum_{j \in I^*} |\tilde{\mu}_j - \lambda_j^*| |\langle f_j, f_k \rangle| \geq (4 - 2\sqrt{2})rL \right).
 \end{aligned}$$

Reasoning as in (4.6) above we find

$$\sum_{k \notin I^*} \mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n f_k(X_i) - E f_k(X_1) \right| \geq 2\sqrt{2}rL \right) \leq \frac{\delta}{M}.$$

To bound the last sum in (4.10) we first notice that Theorem 1 (if we replace there $r(\delta/2)$ by the larger value $r(\delta/(2M))$), cf. Theorem 4) applies to $\tilde{\mu}$ given by (4.8). In particular

$$\mathbb{P} \left(\sum_{j \in I^*} |\tilde{\mu}_j - \lambda_j^*| \geq \frac{4\sqrt{2}}{L} k^* r \right) \leq \frac{\delta}{M}.$$

Therefore, by *Condition (A)*, we have

$$\begin{aligned}
 & \sum_{k \notin I^*} \mathbb{P} \left(\sum_{j \in I^*} |\tilde{\mu}_j - \lambda_j^*| |\langle f_j, f_k \rangle| \geq (4 - 2\sqrt{2})rL \right) \\
 & \leq \sum_{k \notin I^*} \mathbb{P} \left(\sum_{j \in I^*} |\tilde{\mu}_j - \lambda_j^*| \geq 32(4 - 2\sqrt{2})k^* rL \right) \\
 & \leq \sum_{k \notin I^*} \mathbb{P} \left(\sum_{j \in I^*} |\tilde{\mu}_j - \lambda_j^*| \geq \frac{4\sqrt{2}}{L} k^* r \right) \leq \delta,
 \end{aligned}$$

which holds since $L^2 \geq 1/3$. Collecting all the bounds above we obtain

$$P(\hat{I} \neq I^*) \leq 2\delta + \frac{2\delta}{M},$$

which concludes the proof. \blacksquare

4.2. Example: Identifying true components in mixtures of Gaussian densities.

Consider an ensemble of M Gaussian densities p_j 's in \mathbb{R}^d with means μ_j and covariance matrices $\tau_j \mathbb{I}_d$, where \mathbb{I}_d is the unit $d \times d$ matrix. In what follows we show that *Condition (A)* holds if the means of the Gaussian densities are well separated and we make this precise below. Therefore, in this case, Theorem 5 guarantees that if the weights of the mixture are above the threshold given in Condition B, we can recover the true mixture components with high probability via our procedure. The densities are

$$p_j(x) = \frac{1}{(2\pi\tau_j^2)^{d/2}} \exp\left(-\frac{\|x - \mu_j\|_2^2}{2\tau_j^2}\right),$$

where $\|\cdot\|_2$ denotes the Euclidean norm. Consequently, $f_j = p_j/\|p_j\|$ with $\|p_j\| = (4\pi\tau_j^2)^{-d/4}$. Recall that *Condition (A)* requires $16\rho^* = 16 \max_{i \in I^*, j \neq i} |\langle f_i, f_j \rangle| \leq 1/k^*$. Let $\tau_{\max} = \max_{1 \leq j \leq M} \tau_j$ and $D_{\min}^2 = \min_{k \neq j} \|\mu_k - \mu_j\|_2^2$. Via simple algebra we obtain

$$\rho^* \leq \exp\left(-\frac{D_{\min}^2}{4\tau_{\max}^2}\right).$$

Therefore, *Condition (A)* holds if

$$(4.11) \quad D_{\min}^2 \geq 4\tau_{\max}^2 \log(16k^*).$$

Using this and Theorem 5 we see that SPADES identifies the true components in a mixture of Gaussian densities if the square Euclidean distance between any two means is large enough as compared to the largest variance of the components in the mixture.

Note that *Condition (B)* on the size of the mixture weights involves the constant L , which in this example can be taken as

$$L = \max\left(\frac{\sqrt{3}}{3}, \max_{1 \leq j \leq M} \|f_j\|_{\infty}\right) = \max\left(\frac{\sqrt{3}}{3}, (\pi\tau_{\min}^2)^{-d/4}\right),$$

where $\tau_{\min} = \min_{1 \leq j \leq M} \tau_j$.

Remark. Often both the location and scale parameters are unknown. In this situation, as suggested by the Associate Editor, the SPADES procedure can be applied to a family of densities with both scale and location parameters chosen from an appropriate grid. By Theorem 1 the resulting estimate will be a good approximation of the unknown target density. An immediate modification of Theorem 5, as in [9], further guarantees that SPADES identifies correctly the important components of this approximation.

5. SPADES FOR ADAPTIVE NONPARAMETRIC DENSITY ESTIMATION

We assume in this section that the density f is defined on a bounded interval of \mathbb{R} that we take without loss of generality to be the interval $[0, 1]$. Consider a countable system of functions $\{\psi_{lk}, l \geq -1, k \in V(l)\}$ in L_2 , where the set of indices $V(l)$ satisfies $|V(-1)| \leq C$, $2^l \leq |V(l)| \leq C2^l$, $l \geq 0$, for some constant C , and where the functions ψ_{lk} satisfy

$$(5.1) \quad \|\psi_{lk}\| \leq C_1, \quad \|\psi_{lk}\|_\infty \leq C_1 2^{l/2}, \quad \left\| \sum_{k \in V(l)} \psi_{lk}^2 \right\|_\infty \leq C_1 2^l,$$

for all $l \geq -1$ and for some $C_1 < \infty$. Examples of such systems $\{\psi_{lk}\}$ are given, for instance, by compactly supported wavelet bases, see, e.g., [25]. In this case $\psi_{lk}(x) = 2^{l/2} \psi(2^l x - k)$ for some compactly supported function ψ . We assume that $\{\psi_{lk}\}$ is a frame, i.e., there exist positive constants c_1 and c_2 depending only on $\{\psi_{lk}\}$ such that, for any two sequences of coefficients β_{lk}, β'_{lk} ,

$$(5.2) \quad c_1 \sum_{l=-1}^{\infty} \sum_{k \in V(l)} (\beta_{lk} - \beta'_{lk})^2 \leq \left\| \sum_{l=-1}^{\infty} \sum_{k \in V(l)} (\beta_{lk} - \beta'_{lk}) \psi_{lk} \right\|^2 \leq c_2 \sum_{l=-1}^{\infty} \sum_{k \in V(l)} (\beta_{lk} - \beta'_{lk})^2.$$

If $\{\psi_{lk}\}$ is an orthonormal wavelet basis, this condition is satisfied with $c_1 = c_2 = 1$.

Now, choose $\{f_1, \dots, f_M\} = \{\psi_{lk}, -1 \leq l \leq l_{\max}, k \in V(l)\}$ where l_{\max} is such that $2^{l_{\max}} \asymp n/(\log n)$. Then also $M \asymp n/(\log n)$. The coefficients λ_j are now indexed by $j = (l, k)$, and we set by definition $\lambda_{(l,k)} = 0$ for $(l, k) \notin \{-1 \leq l \leq l_{\max}, k \in V(l)\}$. Assume that there exist coefficients β_{lk}^* such that

$$f = \sum_{l=-1}^{\infty} \sum_{k \in V(l)} \beta_{lk}^* \psi_{lk}$$

where the series converges in L_2 . Then Theorem 3 easily implies the following result.

Theorem 6. *Let f_1, \dots, f_M be as defined above with $M \asymp n/(\log n)$, and let ω_j be given by (3.8) for $\delta = n^{-2}$. Then for all $n \geq 1$, $\lambda \in \mathbb{R}^M$ we have with probability at least $1 - n^{-2}$,*

$$(5.3) \quad \|f^\spadesuit - f\|^2 \leq K \left(\sum_{l=-1}^{\infty} \sum_{k \in V(l)} (\lambda_{(l,k)} - \beta_{lk}^*)^2 + \sum_{(l,k) \in J(\lambda)} \left[\frac{1}{n} \sum_{i=1}^n \psi_{lk}^2(X_i) \frac{\log n}{n} + 2^l \left(\frac{\log n}{n} \right)^2 \right] \right)$$

where K is a constant independent of f .

This is a general oracle inequality that allows one to show that the estimator f^\spadesuit attains minimax rates of convergence, up to a logarithmic factor simultaneously on various functional

classes. We will explain this in detail for the case where f belongs to a class of functions \mathcal{F} satisfying the following assumption for some $s > 0$.

Condition (C). For any $f \in \mathcal{F}$ and any $l' \geq 0$ there exists a sequence of coefficients $\lambda = \{\lambda_{(l,k)}, -1 \leq l \leq l', k \in V(l)\}$ such that

$$(5.4) \quad \sum_{l=-1}^{\infty} \sum_{k \in V(l)} (\lambda_{(l,k)} - \beta_{lk}^*)^2 \leq C_2 2^{-2l's}$$

for a constant C_2 independent of f .

It is well known that Condition (C) holds for various functional classes \mathcal{F} , such as Hölder, Sobolev, Besov classes, if $\{\psi_{lk}\}$ is an appropriately chosen wavelet basis, see, e.g., [25] and the references cited therein. In this case s is the smoothness parameter of the class. Moreover, the basis $\{\psi_{lk}\}$ can be chosen so that Condition (C) is satisfied with C_2 independent of s for all $s \leq s_{\max}$, where s_{\max} is a given positive number. This allows for adaptation in s .

Under Condition (C) we obtain from (5.3) that, with probability at least $1 - n^{-2}$,

$$(5.5) \quad \|f^\spadesuit - f\|^2 \leq \min_{l' \leq l_{\max}} K \left(C_2 2^{-2l's} + \sum_{(l,k): l \leq l'} \left[\frac{1}{n} \sum_{i=1}^n \psi_{lk}^2(X_i) \frac{\log n}{n} + 2^l \left(\frac{\log n}{n} \right)^2 \right] \right)$$

From (5.5) and the last inequality in (5.1) we find for some constant K' , with probability at least $1 - n^{-2}$,

$$(5.6) \quad \begin{aligned} \|f^\spadesuit - f\|^2 &\leq \min_{l' \leq l_{\max}} K' \left(2^{-2l's} + 2^{l'} \left(\frac{\log n}{n} \right) + 2^{2l'} \left(\frac{\log n}{n} \right)^2 \right) \\ &= O \left(\left(\frac{\log n}{n} \right)^{-2s/(2s+1)} \right) \end{aligned}$$

where the last expression is obtained by choosing l' such that $2^{l'} \asymp (n/\log n)^{1/(2s+1)}$. It follows from (5.6) that f^\spadesuit converges with the optimal rate (up to a logarithmic factor) simultaneously on all the functional classes satisfying Condition (C). Note that the definition of the functional class is not used in the construction of the estimator f^\spadesuit , so this estimator is optimal adaptive in the rate of convergence (up to a logarithmic factor) on this scale of functional classes for $s \leq s_{\max}$. Results of such type, and even more pointed (without extra logarithmic factors in the rate and sometimes with exact asymptotic minimax constants) are known for various other adaptive density estimators, see, for instance, [22, 6, 25, 27, 35, 36] and the references cited therein. These papers consider classes of densities that are uniformly bounded by a fixed constant, see the recent discussion in [5]. This prohibits, for example, free scale transformations of densities within a class. Inequality (5.6) does not have this

drawback. It allows to get the rates of convergence for classes of unbounded densities f as well.

Another example is given by the classes of sparse densities defined as follows:

$$\mathcal{L}_0(m) = \left\{ f : [0, 1] \rightarrow \mathbb{R} : f \text{ is a probability density and } \left| \{j : \langle f, f_j \rangle \neq 0\} \right| \leq m \right\}$$

where $m \leq M$ is an unknown integer. If f_1, \dots, f_M is a wavelet system as defined above and $J^* = \{j = (l, k) : \langle f, f_j \rangle \neq 0\}$, then under the conditions of Theorem 6 for any $f \in \mathcal{L}_0(m)$ we have, with probability at least $1 - n^{-2}$,

$$(5.7) \quad \|f^\spadesuit - f\|^2 \leq K \left(\sum_{(l,k) \in J^*} \left[\frac{1}{n} \sum_{i=1}^n \psi_{lk}^2(X_i) \frac{\log n}{n} + 2^l \left(\frac{\log n}{n} \right)^2 \right] \right).$$

From (5.7), using Lemma 1 and the first two inequalities in (5.1) we obtain the following result.

Corollary 3. *Let the assumptions of Theorem 6 hold. Then, for every $L < \infty$ and $n \geq 1$,*

$$(5.8) \quad \sup_{f \in \mathcal{L}_0(m) \cap \{f : \|f\|_\infty \leq L\}} \mathbb{P} \left\{ \|f^\spadesuit - f\|^2 \geq b \left(\frac{m \log n}{n} \right) \right\} \leq (3/2)n^{-2}, \quad \forall m \leq M,$$

where $b > 0$ is a constant depending only on L .

Corollary 3 can be viewed as an analogue for density estimation of the adaptive minimax results for \mathcal{L}_0 classes obtained in the Gaussian sequence model [1, 23] and in the random design regression model [13].

6. NUMERICAL EXPERIMENTS

In this section we describe the algorithm used for the minimization problem (2.4) and we assess the performance of our procedure via a simulation study.

6.1. A coordinate descent algorithm. Since the criterion given in (2.4) is convex, but not differentiable, we adopt an optimization by coordinate descent instead of a gradient-based approach (gradient descent, conjugate gradient, etc) in the spirit of [20, 21]. Coordinate descent is an iterative greedy optimization technique that starts at an initial location $\lambda \in \mathbb{R}^M$ and at each step chooses one coordinate $\lambda_j \in \mathbb{R}$ of λ at random or in order and finds the optimum in that direction, keeping the other variables λ_{-j} fixed at their current values. For convex functions, it usually converges to the global optimum, see [20]. The method is based on the obvious observation that for functions of the type

$$H(\lambda) = g(\lambda) + \omega |\lambda|_1$$

where g is a generic convex and differentiable function, $\omega > 0$ is a given parameter, and $|\lambda|_1$ denotes the ℓ_1 norm, the optimum in a direction $\lambda_j \in \mathbb{R}$ is to the left, right or at $\lambda_j = 0$, depending on the signs of the left and right partial derivatives of H at zero. Specifically, let g_j denote the partial derivative of g with respect to λ_j , and denote by λ_{-j}^0 the vector λ with j -th coordinate set to 0. Then, the minimum in direction j of $H(\lambda)$ is at λ_{-j}^0 if and only if $|g_j(\lambda_{-j}^0)| < \omega$. This observation makes the coordinate descent become the iterative thresholding algorithm described below.

Coordinate Descent.

Given ω , initialize all λ_j , $1 \leq j \leq M$, e.g. with $1/M$.

1. Choose a direction $j \in \{1, \dots, M\}$ and set $\lambda^{old} = \lambda$.
2. If $|g_j(\lambda_{-j}^0)| < \omega$ then set $\lambda = \lambda_{-j}^0$, otherwise obtain λ by line minimization in direction j .
3. If $|\lambda^{old} - \lambda| > \epsilon$ go to 1, where $\epsilon > 0$ is a given precision level.

For line minimization, we used the procedure `linmin` from Numerical Recipes [34], page 508.

6.2. Estimation of mixture weights using the generalized bisection method and a penalized cross-validated loss function. We apply the coordinate descent algorithm described above to optimize the function $H(\lambda)$ given by (2.4), where the tuning parameters ω_j are all set to be equal to the same quantity ω . The theoretical choice of this quantity described in detail in the previous sections may be too conservative in practice. In this section we propose a data driven method for choosing the tuning parameter ω , following the procedure first introduced in [10], which we briefly describe here for completeness.

The procedure chooses adaptively the tuning parameter from a list of candidate values, and it has two distinctive features: the list of candidates is not given by a fine grid of values and the adaptive choice is not given by cross validation, but by a dimension stabilized cross validated criterion. We begin by describing the principle underlying our construction of the set of candidate values which, by avoiding a grid search, provides significant computational savings. We use a generalization of the bisection method to find, for each $0 \leq k \leq M$, a preliminary tuning parameter $\omega = w_k$ that gives a solution $\hat{\lambda}^k$ with exactly k nonzero elements. Formally, denote by $\hat{n}(\omega)$ the number of non-zero elements in the λ obtained by minimizing (2.4) with $\omega_j \equiv \omega$ for a given value of the tuning parameter ω . The Generalized Bisection Method will find a sequence of values of the tuning parameter, w_0, \dots, w_M , such

that $\hat{n}(w_k) = k$, for each $0 \leq k \leq M$. It proceeds as follows, using a queue consisting of pairs (w_i, w_j) such that $\hat{n}(w_i) < \hat{n}(w_j) - 1$.

The General Bisection Method for all k (GBM).

Initialize all w_i with -1 .

1. Choose w_0 very large, such that $\hat{n}(w_0) = 0$. Choose $w_n = 0$, hence $\hat{n}(w_n) = n$.
2. Initialize a queue q with the pair (w_0, w_n) .
3. Pop the first pair (a, b) from the queue.
4. Take $w = (a + b)/2$. Compute $k = \hat{n}(w)$.
5. If $w_k = -1$ make $w_k = w$.
6. If $|\hat{n}(a) - k| > 1$ and $|a - w| > \alpha$, add (a, w) to the back of the queue.
7. If $|\hat{n}(b) - k| > 1$ and $|b - w| > \alpha$, add (w, b) to the back of the queue.
8. If the queue is not empty, go to 3.

This algorithm generalizes the Basic Bisection Method (**BBM**), which is a well established computationally efficient method for finding a root $z \in \mathbb{R}$ of a function $h(z)$, see e.g. [15]. We experimentally observed (see also [10] for a detailed discussion) that using the **GBM** is about 50 times faster than a grid search with the same accuracy.

Our procedure finds the final tuning parameter ω by combining the **GBM** with the dimension stabilized p -fold cross-validation procedure summarized below. Let D denote the whole data set, and let $D = D_1 \cup \dots \cup D_p$ be a partition of D in p disjoint subsets. Let $D_{-j} = D \setminus D_j$. We will denote by w_k^j a candidate tuning parameter determined using the **GBM** on D_{-j} . We denote by I_k^j the set of indices corresponding to the non-zero coefficients of the estimator of λ given by (2.4), for tuning parameter w_k^j on D_{-j} . We denote by $\hat{\lambda}^{kj}$ the minimizers on D_{-j} of the unpenalized criterion $\hat{\gamma}(\lambda)$, with respect only to those λ_l with $l \in I_k^j$. Let $L_k^j =: \hat{\gamma}(\hat{\lambda}^{kj})$, computed on D_j . With this notation, the procedure becomes:

Weight Selection Procedure.

Given: a data set D partitioned into p disjoint subsets, $D = D_1 \cup \dots \cup D_p$. Let $D_{-j} = D \setminus D_j$ for all j .

1. For each $1 \leq k \leq M$ and each fold j of the partition, $1 \leq j \leq p$:
 Use the **GBM** to find w_k^j and I_k^j such that $\hat{n}(w_k^j) = |I_k^j| = k$ on D_{-j} .
 Compute $L_k^j =: \hat{\gamma}(\hat{\lambda}^{kj})$, as defined above, on D_j .

2. For each $1 \leq k \leq M$:

$$\text{Compute } L_k =: \frac{1}{p} \sum_{j=1}^p L_k^j.$$

3. Obtain

$$\hat{k} = \arg \min_k (L_k + 0.5k \frac{\log n}{n}).$$

4. With \hat{k} from Step 3, use the **BBM** on the whole data set D to find the tuning sequence $w_{\hat{k}}$ and then compute the final estimators using the coordinate descent algorithm and tuning parameter $\omega = w_{\hat{k}}$.

In all the numerical experiments described below we took the number of splits $p = 10$.

Remark. We recall that the theoretical results of Section 4.1 show that for correct identification of the mixture components one needs to work with a value of the tuning sequence that is slightly larger than the one needed for good approximations with mixtures of a given density. A good practical approximation of the latter tuning value is routinely obtained by cross-validation; this approximation is however not appropriate if the goal is correct selection, when the theoretical results indicate that a different value is needed. Our modification of the cross-validated loss function via a BIC-type penalty is motivated by the known properties of the BIC-type criteria to yield consistent model selection in a large array of models, see e.g. [7] for results on regression models. The numerical experiments presented below show that this is also the case for our criterion in the context of selecting mixture components. The theoretical investigation of this method is beyond the scope of this paper and will be undertaken in future research.

6.3. Numerical results. In this subsection we illustrate the performance of our procedure via a simulation study.

6.3.1. One-dimensional densities. We begin by investigating the ability of SPADES, with its tuning parameter chosen as above, to (i) approximate well, with respect to the L_2 norm, a true mixture; (ii) to identify the true mixture components. We conducted a simulation study where the true density is a mixture of Gaussian densities with $k^* = 2$ and, respectively, $k^* = 5$ true mixture components. The mixture components are chosen at random from a larger pool of M Gaussians $\mathcal{N}(aj, 1)$, $1 \leq j \leq M$, where for $k^* = 2$ we take $a = 4$, and for $k^* = 5$ we take $a = 5$. These choices for a ensure that the identifiability condition (4.11) is satisfied. The true components correspond to the first k^* Gaussian densities from our list, and their weights in the true mixture are all equal to $1/k^*$. The maximum size M of the

candidate list we considered is $M = 200$, for $k^* = 2$ and $M = 600$, for $k^* = 5$. All the results obtained below are relative to $S = 100$ simulations. Each time, a sample of size n is obtained from the true mixture and is the input of the procedure described in Section 6.2.

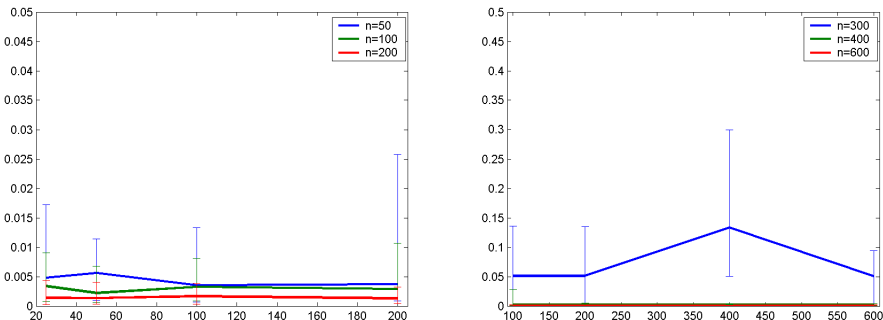


FIGURE 1. Median L_2 error $\|f^* - f^\spadesuit\|^2$ for $|I^*| = 2$ respectively $|I^*| = 5$. The error bars are the 25 and 75 percentiles.

We begin by evaluating the accuracy with respect to the L_2 norm of the estimates of f^* . We investigate the sensitivity of our estimates relative to an increase in the dictionary size and k^* . In Figure 1, we plot the median over 100 simulations of $\|f^* - f^\spadesuit\|^2$ versus the size M of the dictionary, when the true mixture cardinality is $k^* = 2$ (left panel) and $k^* = 5$ (right panel). For $k^* = 2$ we considered three instances of sample sizes $n = 50, 100, 200$ and we varied M up to 200. For $k^* = 5$ we considered three larger instances of sample sizes $n = 300, 400, 600$ and we varied M up to 600. These experiments provide strong support for our theoretical results: the increase in M does not significantly affect the quality of estimation, and an increase in k^* does. For larger values of k^* we need larger sample sizes to obtain good estimation accuracy.

We next investigated the ability of the SPADES to find the exact mixture components. Figure 2 shows a plot of the percentage of times the exact mixture components were found versus M . We considered the same combinations (n, M) as in Figure 1. Again, observe that the performance does not seriously degrade with the dictionary size M , and is almost unaffected by its increase once a threshold sample size is being used. However, notice that on the difference from the results presented in Figure 1, correct identification is poor below the threshold sample size, which is larger for larger k^* . This is in accordance with our theoretical results: recall *Condition (B)* of Section 4.1 on the minimum size of the mixture weights.

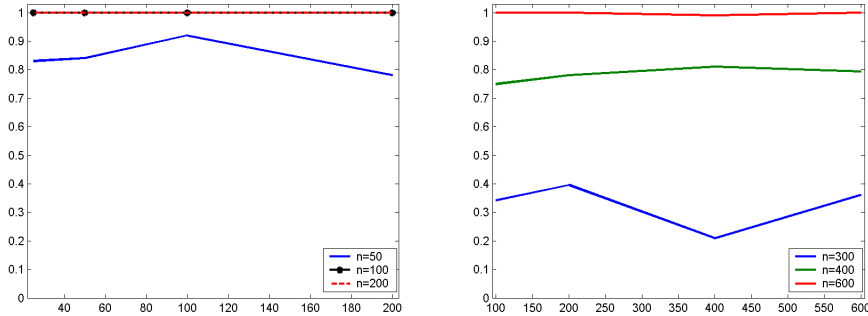


FIGURE 2. Percentage of times $I^* = \hat{I}$ obtained from 100 runs, for $|I^*| = 2$ respectively $|I^*| = 5$.

Indeed, we designed our simulations so that the weights are relatively small for $k^* = 5$, they are all equal to $1/k^* = 0.2$, and a larger sample size is needed for their correct identification.

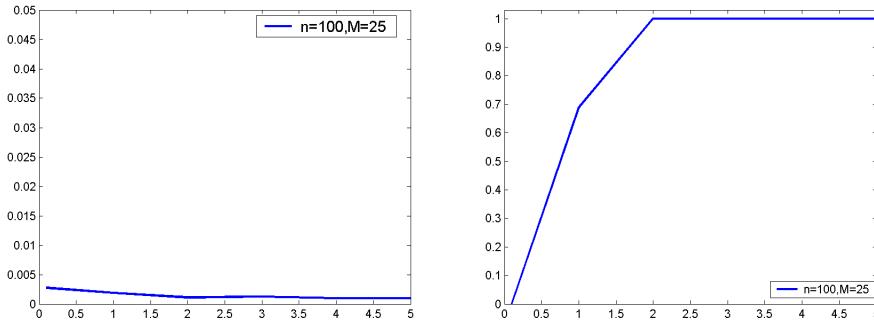


FIGURE 3. Dependence on the distance $D_{\min} = \min_{k \neq j} |\mu_k - \mu_j|$ of the L_2 error $\|f^* - \hat{f}\|^2$ and the percentage of times $I^* = \hat{I}$. In this example, $n = 100, M = 25, |I^*| = 2$.

Finally, we evaluated in Figure 3 the dependence of the error and hit rate (i.e. the percentage of times $I^* = \hat{I}$) on the smallest distance $D_{\min} = \min_{k \neq j} |\mu_k - \mu_j|$ between the means of the densities in the dictionary. The results presented in Figures 1 and 2 above were obtained for the value $D_{\min} = 4$, which satisfies the theoretical requirement for correct mixture identification. On the other hand, D_{\min} can be smaller for good L_2 mixture approximation. It is interesting to see what happens when D_{\min} decreases, so that the mixture elements become very close to one another. In Figure 3 we present the simulations for $k^* = 2, M = 25$ and $n = 100$, which is sufficient to illustrate this point. We see that, although the L_2 error increases slightly when D_{\min} decreases, the deterioration is not crucial. However, as

our theoretical results suggest, the percentage of times we can correctly identify the mixture decreases to zero when the dictionary functions are very close to each other.

6.3.2. Two-dimensional densities. In a second set of experiments our aim was to approximate a two-dimensional probability density on a thick circle (cf. the left panel of Figure 5) with a mixture of isotropic Gaussians. A sample of size 2000 from the circle density is shown in the middle panel of Figure 5. We use a set of isotropic Gaussian candidates with covariance $\Sigma = \mathbb{I}_2$ centered at some of the 2000 locations, such that the Euclidean distance between the means of any two such Gaussians is at least 1. We select from these candidate mixture densities in a greedy iterative manner, each time choosing one of the 2000 locations that is at distance at least 1 from each of those already chosen. As a result, we obtain a dictionary of $M = 248$ candidate densities.

The circle density cannot be exactly represented as a finite mixture of Gaussian components. This is a standard instance of many practical applications in Computer Vision, as the statistics of natural images are highly kurtotic and cannot be exactly approximated by isotropic Gaussians. However, in many practical applications a good approximation of an object that reflects its general shape is sufficient and constitutes a first crucial step in any analysis. We show below that SPADES offers such an approximation.

Depending on the application, different tradeoffs between the number of mixture components (which relates to the computational demand of the mixture model) and accuracy might be appropriate. For example, in real-time applications a small number of mixture elements would be required to fit into the computational constraints of the system, as long as there is no significant loss in accuracy.

For the example presented below we used the **GBM** to determine the mixture weights $\hat{\lambda}^k$, for mixtures with $k = 1, 2, \dots, 248$ components. Let $\gamma_0 = \min_k \hat{\gamma}(\hat{\lambda}^k)$, where we recall that the loss function $\hat{\gamma}$ is given by (2.2) above. We used the quantity $\hat{\gamma}(\hat{\lambda}^k) - \gamma_0$ to measure the accuracy of the mixture approximation. In Figure 4 we plotted $\hat{\gamma}(\hat{\lambda}^k) - \gamma_0$ as a function of k and used this plot to determine the desired trade-off between accuracy and mixture complexity. Based on this plot, we selected the number of mixture components to be 80; indeed, including more components does not yield any significant improvement. The obtained mixture is displayed in the right panel of Figure 5. We see that it successfully approximates the circle density with a relatively small number of components.

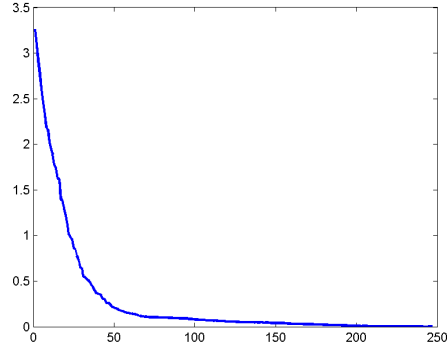


FIGURE 4. Plot of $\widehat{\gamma}(\widehat{\lambda}^k) - \gamma_0$ as a function of the mixture components k

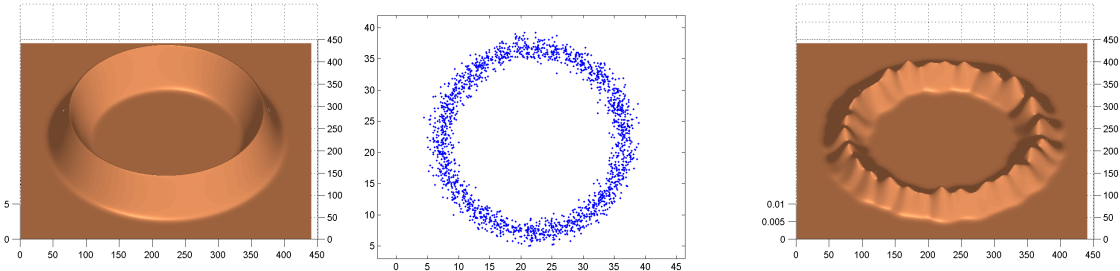


FIGURE 5. A thick circle density, a sample of size 2000 from this density and approximations using a mixture of 80 isotropic Gaussians.

APPENDIX

Lemma 3. (I) Let $\tilde{\mu}$ be given by (4.8). Then $\bar{\mu} = (\tilde{\mu}, 0) \in \mathbb{R}^M$ is a minimizer in $\lambda \in \mathbb{R}^M$ of

$$g(\lambda) = -\frac{2}{n} \sum_{i=1}^n f_{\lambda}(X_i) + \|f_{\lambda}\|^2 + 8Lr \sum_{k=1}^M |\lambda_k|.$$

on the random event \mathcal{B} defined in (4.9).

(II) Any two minimizers of $g(\lambda)$ have non-zero components in the same positions.

Proof. (I). Since g is convex, by standard results in convex analysis, $\bar{\lambda} \in \mathbb{R}^M$ is a minimizer of g if and only if $0 \in D_{\bar{\lambda}}$ where D_{λ} is the subdifferential of $g(\lambda)$:

$$D_{\lambda} = \left\{ w \in \mathbb{R}^M : w_k = -\frac{2}{n} \sum_{i=1}^n f_k(X_i) + 2 \sum_{j=1}^M \lambda_j \langle f_j, f_k \rangle + 8rv_k, v_k \in V_k(\lambda_k), 1 \leq k \leq M \right\}$$

where

$$V_k(\lambda_k) = \begin{cases} \{L\} & \text{if } \lambda_k > 0, \\ \{-L\} & \text{if } \lambda_k < 0, \\ [-L, L] & \text{if } \lambda_k = 0. \end{cases}$$

Therefore, $\bar{\lambda}$ minimizes $g(\cdot)$ if and only if, for all $1 \leq k \leq M$,

$$(6.1) \quad \frac{1}{n} \sum_{i=1}^n f_k(X_i) - \sum_{j=1}^M \bar{\lambda}_j \langle f_j, f_k \rangle = 4Lr \operatorname{sign}(\bar{\lambda}_k), \quad \text{if } \bar{\lambda}_k \neq 0,$$

$$(6.2) \quad \left| \frac{1}{n} \sum_{i=1}^n f_k(X_i) - \sum_{j=1}^M \bar{\lambda}_j \langle f_j, f_k \rangle \right| \leq 4Lr, \quad \text{if } \bar{\lambda}_k = 0.$$

We now show that $\bar{\mu} = (\tilde{\mu}, 0) \in \mathbb{R}^M$ with $\tilde{\mu}$ given in (4.8) satisfies (6.1)–(6.2) on the event \mathcal{B} and therefore is a minimizer of $g(\lambda)$ on this event. Indeed, since $\tilde{\mu}$ is a minimizer of the convex function $h(\mu)$ given in (4.7), the same convex analysis argument as above implies that

$$\frac{1}{n} \sum_{i=1}^n f_k(X_i) - \sum_{j \in I^*} \tilde{\mu}_j \langle f_j, f_k \rangle = 4Lr \operatorname{sign}(\tilde{\mu}_k), \quad \text{if } \tilde{\mu}_k \neq 0, \quad k \in I^*,$$

$$\left| \frac{1}{n} \sum_{i=1}^n f_k(X_i) - \sum_{j \in I^*} \tilde{\mu}_j \langle f_j, f_k \rangle \right| \leq 4Lr, \quad \text{if } \tilde{\mu}_k = 0, \quad k \in I^*.$$

Note that on the event \mathcal{B} we also have

$$\left| \frac{1}{n} \sum_{i=1}^n f_k(X_i) - \sum_{j \in I^*} \tilde{\mu}_j \langle f_j, f_k \rangle \right| \leq 4Lr, \quad \text{if } k \notin I^* \text{ (for which } \bar{\mu}_k = 0, \text{ by construction).}$$

Here $\bar{\mu}_k$ denotes the k th coordinate of $\bar{\mu}$. The above three displays and the fact that $\bar{\mu}_k = \tilde{\mu}_k, k \in I^*$, show that $\bar{\mu}$ satisfies conditions (6.1)–(6.2) and is therefore a minimizer of $g(\lambda)$ on the event \mathcal{B} .

(II). We now prove the second assertion of the lemma. In view of (6.1) the index set S of the non-zero components of any minimizer $\bar{\lambda}$ of $g(\lambda)$ satisfies

$$S = \left\{ k \in \{1, \dots, M\} : \left| \frac{1}{n} \sum_{i=1}^n f_k(X_i) - \sum_{j=1}^M \bar{\lambda}_j \langle f_j, f_k \rangle \right| = 4rL \right\}.$$

Therefore, if for any two minimizers $\bar{\lambda}^{(1)}$ and $\bar{\lambda}^{(2)}$ of $g(\lambda)$ we have

$$(6.3) \quad \sum_{j=1}^M (\bar{\lambda}_j^{(1)} - \bar{\lambda}_j^{(2)}) \langle f_j, f_k \rangle = 0, \quad \text{for all } k,$$

then S is the same for all minimizers of $g(\lambda)$.

Thus, it remains to show (6.3). We use simple properties of convex functions. First, we recall that the set of minima of a convex function is convex. Then, if $\bar{\lambda}^{(1)}$ and $\bar{\lambda}^{(2)}$ are two distinct points of minima, so is $\rho \bar{\lambda}^{(1)} + (1 - \rho) \bar{\lambda}^{(2)}$, for any $0 < \rho < 1$. Re-write this convex combination as $\bar{\lambda}^{(2)} + \rho \eta$, where $\eta = \bar{\lambda}^{(1)} - \bar{\lambda}^{(2)}$. Recall that the minimum value of any convex

function is unique. Therefore, for any $0 < \rho < 1$, the value of $g(\lambda)$ at $\lambda = \bar{\lambda}^2 + \rho\eta$ is equal to some constant C :

$$\begin{aligned} F(\rho) &\triangleq -\frac{2}{n} \sum_{i=1}^n \sum_{j=1}^M (\bar{\lambda}_j^{(2)} + \rho\eta_j) f_j(X_i) + \int \left(\sum_{j=1}^M (\bar{\lambda}_j^{(2)} + \rho\eta_j) f_j(x) \right)^2 dx \\ &\quad + 8rL \sum_{j=1}^M |\bar{\lambda}_j^{(2)} + \rho\eta_j| = C. \end{aligned}$$

By taking the derivative with respect to ρ of $F(\rho)$ we obtain that, for all $0 < \rho < 1$,

$$\begin{aligned} F'(\rho) &= -\frac{2}{n} \sum_{i=1}^n \sum_{j=1}^M \eta_j f_j(X_i) + 8rL \sum_{j=1}^M \eta_j \text{sign}(\bar{\lambda}_j^{(2)} + \rho\eta_j) \\ &\quad + 2 \int \left(\sum_{j=1}^M (\bar{\lambda}_j^{(2)} + \rho\eta_j) f_j(x) \right) \left(\sum_{j=1}^M \eta_j f_j(x) \right) dx = 0. \end{aligned}$$

By continuity of $\rho \mapsto \bar{\lambda}_j^{(2)} + \rho\eta_j$, there exists an open interval in $(0, 1)$ on which $\rho \mapsto \text{sign}(\bar{\lambda}_j^{(2)} + \rho\eta_j)$ is constant for all j . Therefore, on that interval,

$$F'(\rho) = 2\rho \int \left(\sum_{j=1}^M \eta_j f_j(x) \right)^2 dx + C'$$

where C' does not depend on ρ . This is compatible with $F'(\rho) = 0$, $\forall 0 < \rho < 1$, (cf. (6.4)) only if

$$\sum_{j=1}^M \eta_j f_j(x) = 0, \text{ for all } x,$$

and therefore

$$\sum_{j=1}^M \eta_j \langle f_j, f_k \rangle = 0, \text{ for all } k \in \{1, \dots, M\},$$

which is the desired result. This completes the proof of the lemma. \square

REFERENCES

- [1] ABRAMOVICH, F., BENJAMINI, Y., DONOHO, D.L. and JOHNSTONE, I.M.. (2006). Adapting to unknown sparsity by controlling the False Discovery Rate. *Annals of Statistics* **34** 584–653.
- [2] BIAU, G. and DEVROYE, L. (2005). Density estimation by the penalized combinatorial method. *Journal of Multivariate Analysis*, **94**, 196–208.
- [3] BIAU, G., CADRE, B., DEVROYE, L. and GYÖRFI, L. (2008). Strongly consistent model selection for densities. *Test*, in press.
- [4] BICKEL, P.J., RITOV, Y. and TSYBAKOV, A.B. (2007). Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statist.*, to appear.
- [5] BIRGÉ, L. (2008). Model selection for density estimation with L_2 loss. [arXiv:0808.1416](https://arxiv.org/abs/0808.1416)

- [6] BIRGÉ, L. and MASSART, P. (1997). From model selection to adaptive estimation. In *Festschrift for Lucien LeCam: Research Papers in Probability and Statistics*, D. Pollard, E. Torgersen, and G. Yang, Eds., pp. 55-87. Springer, New York.
- [7] BUNEA, F. (2004) Consistent covariate selection and post model selection inference in semiparametric regression, *Annals of Statistics*, **32(3)** 898-927.
- [8] BUNEA, F. (2008) Honest variable selection in linear and logistic models via ℓ_1 and $\ell_1 + \ell_2$ penalization, *The Electronic Journal of Statistics*, **2** 1153 - 1194.
- [9] BUNEA, F. (2008) Consistent selection via the Lasso for high dimensional approximating regression models, *IMS Collections*, **3** 122 - 138.
- [10] BUNEA F. and BARBU A.(2009) Dimension reduction and variable selection in case control studies via regularized likelihood optimization. *Electronic Journal of Statistics*, **3** 1257 - 1287
- [11] BUNEA, F., TSYBAKOV, A.B. and WEGKAMP, M.H. (2007). Aggregation for Gaussian regression. *The Annals of Statistics*, **35** 1674 - 1697.
- [12] BUNEA, F., TSYBAKOV, A.B. and WEGKAMP, M.H. (2006a). Aggregation and sparsity via ℓ_1 -penalized least squares. *Proceedings of 19th Annual Conference on Learning Theory, COLT 2006. Lecture Notes in Artificial Intelligence* **4005** 379–391. Springer-Verlag, Heidelberg.
- [13] BUNEA, F., TSYBAKOV, A.B. and WEGKAMP, M.H. (2007). Sparsity oracle inequalities for the Lasso. *The Electronic Journal of Statistics* **1** 169 - 194.
- [14] BUNEA, F., TSYBAKOV, A.B. and WEGKAMP, M.H. (2007). Sparse density estimation with l1 penalties. *Proceedings of 20th Annual Conference on Learning Theory, COLT 2007. Lecture Notes in Artificial Intelligence* 530 - 544. Springer-Verlag, Heidelberg.
- [15] BURDEN, R.L. and FAIRES, J.D. (2001), Numerical analysis, 7th ed., Pacific Grove, CA: Brooks/Cole
- [16] CHEN, S., DONOHO, D. and SAUNDERS, M. (2001) Atomic decomposition by basis pursuit. *SIAM Review* **43** 129 - 159.
- [17] DEVROYE, L. and LUGOSI, G. (2000) *Combinatorial Methods in density estimation*, Springer.
- [18] DONOHO, D.L. (1995) Denoising via soft-thresholding. *IEEE Trans. Info. Theory* **41** 613-627.
- [19] DONOHO, D.L., ELAD, M. and TEMLYAKOV, V. (2006). Stable Recovery of Sparse Overcomplete Representations in the Presence of Noise. *IEEE Trans. on Information Theory* **52** 6–18.
- [20] FRIEDMAN, J., HASTIE, T.,HOFLING, H. and TIBSHIRANI, R. (2007) Pathwise coordinate optimization, *Annals of Applied Statistics*. **1** 302-332.
- [21] FRIEDMAN, J., HASTIE, T. and TIBSHIRANI, R.(2008) Regularization paths for generalized linear models via coordinate descent, *Technical Report*. Available at <http://www-stat.stanford.edu/jhf/ftp/glmnet.pdf>.
- [22] GOLUBEV, G.K. (1992). Nonparametric estimation of smooth probability densities in L_2 . *Problems of Information Transmission* **28** 44-54.
- [23] GOLUBEV, G.K. (2002). Reconstruction of sparse vectors in white Gaussian noise. *Problems of Information Transmission* **38** 65–79.
- [24] GREENSHTEIN, E. and RITOV, Y. (2004). Persistency in high dimensional linear predictor-selection and the virtue of over-parametrization. *Bernoulli* **10** 971–988.
- [25] HÄRDLE, W., KERKYACHARIAN, G., PICARD, D., and TSYBAKOV, A. (1998). *Wavelets, Approximation and Statistical Applications*. Lecture Notes in Statistics, vol. 129. Springer, New York.
- [26] JAMES, L., PRIEBE, C. and MARCHETTE, D. (2001). Consistent estimation of mixture complexity. *Annals of Statistics*, **29**, 1281–1296.
- [27] KERKYACHARIAN, G., PICARD, D. and TRIBOULEY, K. (1996). L^p adaptive density estimation. *Bernoulli* **2** 229–247.
- [28] KOLTCHINSKII, V. (2005). Model selection and aggregation in sparse classification problems. *Oberwolfach Reports* **2** 2663–2667, Mathematisches Forschungsinstitut Oberwolfach.
- [29] KOLTCHINSKII, V. (2006). Sparsity in penalized empirical risk minimization. *Submitted*.
- [30] LOUBES, J. – M. and VAN DE GEER, S. A. (2002). Adaptive estimation in regression, using soft thresholding type penalties. *Statistica Neerlandica* **56** 453 – 478.
- [31] LOUNICI, K. (2008). Sup-norm convergence rate and sign concentration property of Lasso and Dantzig estimators. *Electronic Journal of Statistics* **2** 90–102.
- [32] MEINSHAUSEN, N. and BÜHLMANN, P. (2006). High-dimensional graphs and variable selection with the Lasso. *Annals of Statistics* **34** 1436–1462.

- [33] NEMIROVSKI, A. (2000). Topics in non-parametric statistics. In P. Bernard, editor, *Ecole d'Eté de Probabilités de Saint-Flour 1998*, volume XXVIII of *Lecture Notes in Mathematics*. Springer, New York.
- [34] PRESS, W.H., TEUKOLSKY, S.A., VETTERLING, W.T. and FLANNERY, B.P. (2007), Numerical recipes 3rd edition: The art of scientific computing, *Cambridge University Press New York, NY, USA*
- [35] RIGOLLET, PH. (2006). Inégalités d'oracle, agrégation et adaptation. PhD thesis, University of Paris 6.
- [36] RIGOLLET, PH., and TSYBAKOV, A. B. (2007). Linear and convex aggregation of density estimators. *Mathematical Methods of Statistics* **16** 260–280.
- [37] RUDEMO, M. (1982). Empirical choice of histograms and kernel density estimators. *Scandinavian Journal of Statistics* **9** 65–78.
- [38] SAMAROV, A. and TSYBAKOV, A. (2007). Aggregation of density estimators and dimension reduction. *Advances in Statistical Modeling and Inference. Essays in Honor of Kjell A. Doksum* (V.Nair, ed.), World Scientific, Singapore e.a., 233–251.
- [39] TIBSHIRANI, R. (1996). Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B.* **58** 267–288.
- [40] TSYBAKOV, A. B. (2003). Optimal rates of aggregation. In *Lecture Notes in Artificial Intelligence*, volume 2777 of *Proceedings of 16th Annual Conference on Learning Theory (COLT) and 7th Annual Workshop on Kernel Machines*. Springer-Verlag, Heidelberg.
- [41] VAPNIK, V. N. (1999). *The Nature of Statistical Learning Theory (Information Science and Statistics)*, Springer.
- [42] VAN DE GEER, S.A. (2008). High dimensional generalized linear models and the Lasso. *The Annals of Statistics* **26** 225–287.
- [43] WASSERMAN, L. A. (2004) *All of Statistics*, Springer.
- [44] WEGKAMP, M. H. (1999). Quasi-Universal Bandwidth Selection for Kernel Density Estimators. *Canadian Journal of Statistics* **27** 409 – 420.
- [45] WEGKAMP, M.H. (2003). Model selection in nonparametric regression. *Annals of Statistics* **31** 252–273.
- [46] ZHANG, C.H. AND HUANG, J. (2008). The sparsity and bias of the Lasso selection in high-dimensional linear regression. *Annals of Statistics* **36** 1567-1594.
- [47] ZHAO, P. AND YU, B. (2007). On model selection consistency of Lasso. *Journal of Machine Learning Research* **7** 2541-2567.
- [48] ZOU, H. (2006). The adaptive Lasso and its oracle properties. *Journal of the American Statistical Association* **101** 1418-1429.

FLORENTINA BUNEA, MARTEN WEGKAMP, ADRIAN BARBU

DEPARTMENT OF STATISTICS

FLORIDA STATE UNIVERSITY

TALLAHASSEE, FL 32306-4330

E-MAIL: {BUNEA,WEGKAMP,BARBU}@STAT.FSU.EDU

ALEXANDRE TSYBAKOV

LABORATOIRE DE STATISTIQUE, CREST

92240 MALAKOFF, FRANCE &

LPMA (UMR CNRS 7599), UNIVERSITÉ PARIS 6

75252 PARIS, CEDEX 05, FRANCE

E-MAIL: ALEXANDRE.TSYBAKOV@UPMC.FR