

CART AND BEST-ORTHO-BASIS: A CONNECTION¹

BY DAVID L. DONOHO

Stanford University and University of California, Berkeley

We study what we call “dyadic CART”—a method of nonparametric regression which constructs a recursive partition by optimizing a complexity penalized sum of squares, where the optimization is over all recursive partitions arising from midpoint splits. We show that the method is adaptive to unknown degrees of anisotropic smoothness. Specifically, consider the anisotropic smoothness classes of Nikol’skii, consisting of bivariate functions $f(x_1, x_2)$ whose finite difference of distance h in direction i is bounded in L^p norm by Ch^{δ_i} , $i = 1, 2$. We show that dyadic CART, with an appropriate complexity penalty parameter $\lambda \sim \sigma^2 \cdot \text{Const} \cdot \log(n)$, is within logarithmic terms of minimax over every anisotropic smoothness class $0 < C < \infty$, $0 < \delta_1, \delta_2 \leq 1$.

The proof shows that dyadic CART is identical to a certain adaptive best-ortho-basis algorithm based on the library of all anisotropic Haar bases. Then it applies empirical basis selection ideas of Donoho and Johnstone. The basis empirically selected by dyadic CART is shown to be nearly as good as a basis ideally adapted to the underlying f . The risk of estimation in an ideally adapted anisotropic Haar basis is shown to be comparable to the minimax risk over anisotropic smoothness classes.

Underlying the success of this argument is harmonic analysis of anisotropic smoothness classes. We show that, for each anisotropic smoothness class, there is an anisotropic Haar basis which is a best orthogonal basis for representing that smoothness class; the basis is optimal not just within the library of anisotropic Haar bases, but among all orthogonal bases of $L^2[0, 1]^2$.

1. Introduction. The CART methodology of tree-structured adaptive nonparametric regression [Breiman, Friedman, Olshen and Stone (1983)] has been widely used in statistical data analysis since its inception more than a decade ago. Built around ideas of recursive partitioning, it develops, based on an analysis of noisy data, a piecewise constant reconstruction, where the pieces are terminal nodes of a data-driven recursive partition.

The best-ortho-basis methodology of adaptive time-frequency analysis [Coifman, Meyer, Quake and Wickerhauser] has, more recently, caught the interest of a wide community of applied mathematicians and signal processing engineers. Based on ideas of recursive partitioning of the time-frequency

Received September 1995; revised November 1996.

¹ Research partially supported by NSF DMS-92-09130 and DMS-95-05151.

AMS 1991 subject classifications. Primary 62G07, 41A30; secondary 62G20, 41A25.

Key words and phrases. Wavelets, anisotropic smoothness, anisotropic Haar basis, best orthogonal basis, minimax estimation, spatial adaptation, oracle inequalities.

plane, it develops, from an analysis of a given signal, a segmented basis, where the segments are terminal nodes in a data-driven recursive segmentation of the time axis.

Both methods are concerned with recursive dyadic segmentation; therefore trees and tree pruning are key data structures and underlying algorithms in both areas. In addition, there is a mathematical connection between the areas.

Sudeshna Adak, while a graduate student at Stanford University, pointed out that central algorithms in the two subjects are really the same: namely, the optimal pruning algorithm in Theorem 10.7, page 285, in the CART book [Breiman, Friedman, Olshen and Stone (1983)] and in the proposition on page 717, in the best-basis paper [Coifman and Wickerhauser (1992)]. Both theorems assert that, given a function $\mathcal{E}(T)$, which assigns numerical values to a binary tree and its subtrees, and supposing that the function obeys a certain additivity property, the optimal subtree is obtained by breadth-first, bottom-up pruning of the complete tree.

On the other hand, the subjects are different, since in the CART case, one is searching for an optimal *function* on a multidimensional Cartesian product domain, and in the BOB case, one is searching for an optimal orthogonal *basis* for the vector space of $1 - d$ signals of length n .

This paper will exhibit a precise connection between CART and BOB in a specific setting—where one is seeking an optimal function–basis built from rectangular blocks on a product domain. In this setting we show that certain specific variants of the two apparently different methodologies lead to identical fast algorithms and identical solutions.

1.1. *An implication.* The connection between CART and best basis affords new insights about recursive partitioning methods. Donoho and Johnstone (1994b) have investigated the use of adaptively chosen bases for noise removal. They have developed so-called oracle inequalities which show that certain schemes for basis selection in the presence of noisy data will work well. By adapting such ideas from the best-basis setting to the CART setting, we are able to establish new results on the performance of optimal dyadic recursive partitioning. In particular, we are able to show that such methods can be nearly minimax simultaneously over a wide range of anisotropic smoothness spaces.

We assume observations of the form

$$(1.1) \quad y(i_1, i_2) = \tilde{f}(i_1, i_2) + \sigma z(i_1, i_2), \quad 0 \leq i_1, i_2 < n,$$

where n is dyadic (an integral power of 2), $z(i_1, i_2)$ is a white Gaussian noise, and $\sigma > 0$ is a noise level. We assume the observations are related to the underlying f by cell averaging;

$$(1.2) \quad \tilde{f}(i_1, i_2) = \text{ave}\{f \mid [i_1/n, (i_1 + 1)/n) \times [i_2/n, (i_2 + 1)/n)\}.$$

Our goal is to recover the denoised cell averages with small mean-squared error $\text{MSE}(\hat{f}, \tilde{f}) = E \sum_{i_1, i_2} (\hat{f}(i_1, i_2) - \tilde{f}(i_1, i_2))^2 / n^2$. About f we will assume

that it belongs to a certain class \mathcal{F} , and we will compare performance of estimates with the best mean-squared error available uniformly over the class \mathcal{F} , that is, the minimax risk

$$(1.3) \quad M^*(\sigma, n; \mathcal{F}) = \inf_{\hat{f}(\cdot)} \sup_{f \in \mathcal{F}} \text{MSE}(\hat{f}(y), \tilde{f}).$$

For our \mathcal{F} we consider anisotropic smoothness classes $\mathcal{F}_p^{\delta_1, \delta_2}(C)$ consisting of functions on $[0, 1]^2$ obeying $\|D_h^1 f\|_p \leq Ch^{\delta_1}$, $\|D_h^2 f\|_p \leq Ch^{\delta_2}$, for all $h \in (0, 1)$, where D_h^i denotes the finite difference of distance h in direction i . Such spaces were introduced and systematically studied by Nikol'skii (1969) for structure and imbedding theorems; see Temlyakov (1993) for approximation theorems. We let \mathcal{AS} denote the scale of all such classes, where $0 < \delta_1, \delta_2 \leq 1$ and $0 < C < \infty$ and where p obeys the constraint $\frac{1}{p} < \rho + \frac{1}{2}$, with $\rho = \delta_1 \delta_2 / (\delta_1 + \delta_2)$.

Our main result is the following theorem.

THEOREM 1.1. *Dyadic CART (defined in Section 2), with the specific complexity penalty $\lambda = \lambda(\sigma, \log_e(n))$ defined in Section 7 ($\lambda \asymp \sigma^2 \log(n)$), comes within logarithmic factors of minimax over each functional class $\mathcal{F}_p^{\delta_1, \delta_2}(C)$, where $0 < \delta_1, \delta_2 \leq 1$, $C > 0$ and $1/p < \rho + 1/2$. If $\hat{f}^{*, \lambda}$ denotes the dyadic CART estimator, then*

$$(1.4) \quad \sup_{\mathcal{F}} \text{MSE}(\hat{f}^{*, \lambda}, \tilde{f}) \leq \text{Const}(\delta_1, \delta_2, p) \log(n) M^*(\sigma, n; \mathcal{F}) \quad \text{as } n \rightarrow \infty$$

for each $\mathcal{F} \in \mathcal{AS}$.

In short, the estimator behaves nearly as well over any class in the scale \mathcal{AS} as one could achieve *knowing precisely* which smoothness class were true. However, the construction of the optimal recursive partitioning estimator requires no knowledge of which smoothness class might actually be the case. (Indeed, we are unaware of any previous literature suggesting a connection between such smoothness classes and CART).

This type of near minimaxity is not possible by isotropic approaches, such as thresholding in standard isotropic wavelets bases or isotropic Fourier series. Those orthogonal bases do not provide sufficiently sparse decompositions of anisotropic smoothness classes; dyadic CART, in contrast, is associated with harmonic analysis in a specially adapted basis which provides optimal sparsity decompositions; see Section 8.4.

This type of near minimaxity is also not possible by using linear methods, such as anisotropic kernel methods. Even allowing a choice of global bandwidth which is different in each of the two directions and choosing those two directional bandwidths optimally for the class in question will not lead to near-minimax estimates for those classes with $p < 2$; see Section 11.3.

It is possible to get results of comparable near minimaxity by using anisotropic wavelet schemes; our aim here is not really to delineate all near-minimax approaches, but instead to demonstrate the near minimaxity

of a form of adaptive recursive partitioning. In particular we show that there is a theoretical motivation for using recursive partitioning in a setting where objects may possess different degrees of smoothness in different directions.

1.2. *Plan of the paper.* In Sections 2 through 6 we develop the connection between CART methods and best-basis methods. Section 2 defines dyadic CART and describes its fast algorithm. Section 3 defines a library of anisotropic Haar bases and describes a fast algorithm for finding a best anisotropic Haar basis from given data, where “best” is defined in the Coifman–Wickerhauser sense. In Sections 4 and 5, building on an insight of Engel (1994), we point out that, with traditional choices of entropy, best-ortho-basis is different from CART, but that, with a special *hereditary entropy*, the two methods are the same.

In Sections 7 and 8 we discuss ideas first developed in the best-basis setting. Section 7 develops oracle inequalities, which show how to select a basis empirically from noisy data to yield a basis that is nearly as good as the ideal basis which could be designed based on noiseless data. Section 8 describes the best-basis problem for anisotropic smoothness classes and shows that a certain kind of anisotropic Haar basis is, in one sense, a best basis.

In Section 9, building on Sections 7 and 8, we show that a certain best-basis denoising technique [introduced by Donoho and Johnstone (1994b)]—which is different from CART—is nearly minimax over the scale of anisotropic smoothness classes. Section 10 establishes our main result for CART by comparing the CART estimator with this best-basis denoising method and showing that the two estimates have comparable performance over anisotropic smoothness spaces. Section 11 mentions comparisons and generalizations.

2. Dyadic CART. We change the notation slightly from (1.1). We observe noisy two-dimensional data on a regular square $n \times n$ array of “pixels”

$$(2.1) \quad y(i_1, i_2) = f(i_1, i_2) + \sigma z(i_1, i_2), \quad 0 \leq i_1, i_2 < n,$$

where (in a change from the last section) f is the object of interest—an $n \times n$ array—and z is a standard Gaussian white noise [i.i.d. $N(0, 1)$]. We also introduce a fruitful abuse of notation: we write $[0, n)$ for the *discrete interval* $\{0, \dots, n - 1\}$. Thus $[0, n)^2$ is a discrete square, and so on. Here and below we also write $\mathbf{i} = (i_1, i_2)$, so $y(\mathbf{i}) = f(\mathbf{i}) + \sigma z(\mathbf{i})$, for $\mathbf{i} \in [0, n)^2$ is an equivalent form of (2.1). Finally, we use the variable $N = n^2$ to stand for the cardinality of the $n \times n$ array y .

In this setting, the CART methodology constructs a piecewise constant estimator \hat{f} of f ; data adaptively, it builds a partition \mathcal{P} of $[0, n)^2$ and finds \hat{f} by the rule

$$(2.2) \quad \hat{f}(\mathbf{i} | \mathcal{P}) = \text{ave}\{y | R(\mathbf{i}; \mathcal{P})\},$$

where $R(\mathbf{i}; \mathcal{P})$ denotes the rectangle of the partition \mathcal{P} containing \mathbf{i} .

2.1. *Optimal dyadic CART.* There are several variants of CART, depending on the procedure used to construct the partition \mathcal{P} . In this paper, we are only interested in *optimal* (nongreedy) dyadic recursive partitioning. With an acknowledged risk of misunderstanding, we call this *dyadic CART*. We define terms.

Dyadic partitioning. Starting from the trivial partition $\mathcal{P}_0 = \{[0, n]^2\}$ we may generate new partitions by splitting $[0, n]^2$ into two pieces either vertically or horizontally, yielding either the partition $\{[0, n/2] \times [0, n], [n/2, n] \times [0, n]\}$ or $\{[0, n] \times [0, n/2], [0, n] \times [n/2, n]\}$. We can apply this splitting recursively, generating other partitions. Thus, let $P = \{R_1, \dots, R_k\}$ be a partition and let R stand for one of the rectangles in the partition. We can create a new partition by splitting R in half horizontally or vertically. If $R = [a, b] \times [c, d)$ then let $R^{1,0}$ and $R^{1,1}$ denote the results of horizontal splitting, that is,

$$R^{1,0} = [a, (a+b)/2] \times [c, d),$$

$$R^{1,1} = [(a+b)/2, b] \times [c, d);$$

while we let $R^{2,0}$ and $R^{2,1}$ denote the results of vertical splitting,

$$R^{2,0} = [a, b] \times [c, (c+d)/2),$$

$$R^{2,1} = [a, b] \times [(c+d)/2, d).$$

Note that if $b = a + 1$ or $d = c + 1$ then horizontal-vertical splitting is not possible; only nonempty rectangles are allowed.

As an example, if we split vertically the rectangle $R = R_l$, say, we produce the $k + 1$ -element partition $\{R_1, \dots, R_{l-1}, R_l^{2,0}, R_l^{2,1}, R_{l+1}, \dots, R_k\}$.

A *recursive dyadic partition* is any partition reachable by successive application of these rules.

Optimal partitions. CART is often used to refer to “greedy growing” followed by “optimal pruning,” where the partition \mathcal{P} is constructed in a heuristic, myopic fashion. For the purposes of this paper, we consider instead the use of optimizing partitions, where the dyadic partition \mathcal{P} is constructed as the optimum of the *complexity penalized residual sum of squares*. Thus, with

$$(2.3) \quad \text{CPRSS}(\mathcal{P}, \lambda) = \|y - \hat{f}(\cdot | \mathcal{P})\|_{l_k}^2 + \lambda \#(\mathcal{P}),$$

what we will call (again in perhaps a slight abuse of nomenclature) dyadic CART seeks the partition

$$(2.4) \quad \hat{\mathcal{P}}_\lambda = \underset{\mathcal{P}}{\operatorname{argmin}} \text{CPRSS}(\mathcal{P}, \lambda).$$

The idea of using globally optimal partitions is covered in passing in Breiman, Friedman, Olshen and Stone (1983), Chapter 10. For the moment we let λ be a free parameter; in Section 7 we will propose a specific choice.

Dyadic CART differs from what is usually called CART, in that dyadic CART can split rectangles only in half, while general CART can split rectan-

gles in all proportions. While the extra flexibility of general CART may be useful, this flexibility is sufficient to make the finding of an exactly optimal partition unwieldy. Dyadic CART allows a more limited range of possible partitions, which makes it possible to find an optimal partition in $O(N)$ time.

2.2. *Fast optimal partitioning.* To describe the algorithm, we introduce some notation.

Rectangles. We use I generically to denote dyadic intervals, that is, intervals $I = [a, b)$ with $a = nk/2^j$ and $b = n(k + 1)/2^j$ with $n \geq 2^j$ and $0 \leq k < 2^j$. We use R to denote dyadic rectangles, that is, rectangles $I_1 \times I_2$.

Parents and siblings. Two dyadic rectangles are *siblings* if their union is a dyadic rectangle. This is equivalent to saying that we can write either

$$(2.5) \quad R_i = I_i \times I_0, \quad i = 1, 2,$$

or

$$(2.6) \quad R_i = I_0 \times I_i, \quad i = 1, 2,$$

where I_0, I_1, I_2 are dyadic intervals and

$$(2.7) \quad \begin{aligned} I_1 &= [n \times 2k/2^j, n \times (2k + 1)/2^j), \\ I_2 &= [n \times (2k + 1)/2^j, n \times (2k + 2)/2^j), \end{aligned}$$

with $0 \leq k < 2^{j-1}$, $0 \leq j < \log_2(n) - 1$. A pair satisfying (2.5) is a pair of *horiz-sibs*; a pair satisfying (2.6) is a pair of *vert-sibs*.

The union of two siblings is the parent rectangle. Each rectangle generally has two siblings—a *vert-sib* and a *horiz-sib*—and two parents—a *vert-parent* and a *horiz-parent*. Parents generally have two sets of children: a pair of *horiz-kids* and a pair of *vert-kids*. In extreme cases a rectangle may have only a *vert-sib* [if it is very wide, such as $[0, n) \times [0, n/2)$], or only a *horiz-sib* [if it is very tall, such as $[0, n/2) \times [0, n)$]. In some cases a rectangle may have only *vert-kids* [if it is very narrow, such as $[0, 1) \times [0, n/2)$] or only *horiz-kids* [if it is very short, such as $[0, n/2) \times [0, 1)$].

Inheritance. CPRSS has an “inheritance property” which we see more easily by taking a general point of view. Let $\text{CART}(R)$ denote the problem of finding the optimal partition for *just the data falling in the dyadic rectangle* R :

$$[\text{CART}(R)] \quad \hat{\mathcal{P}}(R) = \operatorname{argmin} \|y - \hat{f}(\cdot | \mathcal{P}(R))\|_{l^2_{(R)}}^2 + \lambda \#(\mathcal{P}(R)).$$

Here $\mathcal{P}(R)$ denotes a recursive dyadic partition of R , and $\|\cdot\|_{l^2_{(R)}}$ refers to the sum-of-squares only of data falling in the rectangle R .

Here is the inheritance property of optimal partitions. Let R be a dyadic rectangle and suppose it has both *vert-children* and *horiz-children*. Then the optimal partition of R is either (1) the trivial partition $\{R\}$, or (2) the union of optimal partitions of the *horiz-kids* $\hat{\mathcal{P}}(R^{1,0}) \cup \hat{\mathcal{P}}(R^{1,1})$, or (3) the union of

optimal partitions of the vert-kids $\hat{\mathcal{P}}(R^{2,0}) \cup \hat{\mathcal{P}}(R^{2,1})$. Which of these three cases holds can be determined by holding a “tournament,” selecting the winner as the smallest of the three numbers

$$\|y - \text{ave}\{y \mid R\}\|_{l^2(R)}^2, \text{CART}(R^{1,0}) + \text{CART}(R^{1,1}), \\ \text{CART}(R^{2,0}) + \text{CART}(R^{2,1}).$$

The exception to this rule is of course at the finest scale: a 1×1 rectangle has no children, and so the optimal partition of such an R is just the trivial partition $\{R\}$.

By starting from the next-to-finest scale and applying the inheritance property, we can get the optimal partitions of all 2×1 rectangles, and of all 1×2 rectangles. By going to the next coarser level and applying inheritance, we can get the optimal partitions of all 4×1 , of all 2×2 and of all 1×4 rectangles and so on. Continuing in a fine-to-coarse or bottom-up fashion, we eventually get to the coarsest level and obtain an optimal partition for $[0, n]^2$.

There are approximately $2n$ dyadic intervals and hence approximately $4n^2 = 4N$ dyadic rectangles. Each dyadic rectangle is visited once in the main loop of the algorithm and there are at most a certain constant number C of additions and multiplications per visit. The total work is $\leq C4N$ flops and $\leq 16N$ storage locations. See the appendix in Donoho (1995) for a formal description of the algorithm.

3. Best-ortho-basis. We now turn attention away from CART. We recall the standard notation for Haar functions in dimension 1. Let I be a dyadic subinterval of $[0, n]$ and let $\chi_I(i) = |I|^{-1/2} \mathbf{1}_I(i)$. If I contains at least two points, set $h_I(i) = (\mathbf{1}_{I^{(0)}}(i) - \mathbf{1}_{I^{(1)}}(i))|I|^{-1/2}$, where $I^{(1)}$ is the right half of I and $I^{(0)}$ is the left half of I .

Using these, we can build anisotropic Haar functions in two dimensions. Let R be a dyadic rectangle $I_1 \times I_2$; we can form three atoms

$$\phi_R^0(i_1, i_2) = \chi_{I_1}(i_1) \chi_{I_2}(i_2), \\ \phi_R^1(i_1, i_2) = h_{I_1}(i_1) \chi_{I_2}(i_2), \\ \phi_R^2(i_1, i_2) = \chi_{I_1}(i_1) h_{I_2}(i_2).$$

These are naturally associated with the rectangle R ; ϕ_R^0 is, up to scaling, the indicator of R , while ϕ_R^1 and ϕ_R^2 are associated with horizontal and vertical midpoint splits of R .

Adapting terminology proposed by Mallat and Zhang (1993) in a different setting, we call the ϕ_R^s atoms, and the collection of all such atoms ϕ_R^s indexed by (R, s) makes up a dictionary of atoms. This dictionary is overcomplete; it contains less than or equal to $3n^2 \approx 3N$ atoms, while the span of these elements is of dimension only N .

3.1. *Anisotropic Haar bases.* Certain structured subcollections of the elements of \mathcal{D} make up orthogonal bases. These subcollections are in correspon-

dence with *complete recursive partitions*, that is to say, recursive dyadic partitions in which all terminal nodes are 1×1 rectangles $[i_1, i_1 + 1) \times [i_2, i_2 + 1)$ containing a single point $\mathbf{i} = (i_1, i_2)$.

Given a complete recursive partition \mathcal{P}^* , the corresponding orthobasis \mathcal{B} is constructed as follows. Let $NT(\mathcal{P}^*)$ be the collection of all rectangles encountered at nonterminal stages of the recursive partitioning leading to \mathcal{P}^* . Let $R \in NT(\mathcal{P}^*)$. As R is nonterminal it will be further subdivided in forming \mathcal{P}^* ; that is, it will be split either horizontally or vertically; let $s(R) = 1$ or 2 according to the splitting variable chosen. Then define \mathcal{B} as the collection of all such $\phi_R^{s(R)}$ and $\chi_{[0, n]^2}$:

$$(3.1) \quad \mathcal{B}(\mathcal{P}^*) = \{ \chi_{[0, n]^2} \} \cup \{ \phi_R^{s(R)} \}_{R \in NT(\mathcal{P}^*)}.$$

THEOREM 3.1. *Let \mathcal{P}^* be a complete recursive dyadic partition of $[0, n]^2$ and let $\mathcal{B}(\mathcal{P}^*)$ be constructed as in (3.1). This is an orthobasis for the N -dimensional vector space of $n \times n$ arrays.*

PROOF. Indeed, \mathcal{B} has cardinality N , and the elements of \mathcal{B} are normalized and pairwise orthogonal. The pairwise orthogonality comes from two simple facts. Take any two distinct elements in \mathcal{B} ; then either they have disjoint support, or the support of one is included in the support of the other. In the first instance, orthogonality is immediate; in the second instance, orthogonality follows from two observations: (i) one element of the pair, call it ϕ , is supported in a rectangle on which the other element, φ say, is constant; and (ii) the element ϕ has zero mean, and so is orthogonal to any function which is constant on its support, that is, to φ . \square

Each such basis \mathcal{B} has a fast transform, produced in a fashion similar to the Haar transform in dimension 1. Indeed, the coefficients in such a basis can be computed in terms of block averages and differences of block averages. If $S(R) = \sum_{\mathbf{i} \in R} y(\mathbf{i})$ denotes the sum of values in a rectangle R , then of course

$$(3.2) \quad \langle y, \chi_R \rangle = S(R)|R|^{-1/2},$$

while, if $(R^{1,0}, R^{1,1})$ are horizontal kids of R ,

$$(3.3) \quad \langle y, \phi_R^1 \rangle = (S(R^{1,1}) - S(R^{1,0}))|R|^{-1/2},$$

and, if $(R^{2,0}, R^{2,1})$ are vertical kids of R ,

$$(3.4) \quad \langle y, \phi_R^2 \rangle = (S(R^{2,1}) - S(R^{2,0}))|R|^{-1/2}.$$

These relations are useful because there is a simple “pyramid-of-adders” for calculating all $(S(R): R \in NT(\mathcal{P}^*))$ in order N time. See the appendix in Donoho (1995b) for a formal description of the algorithm.

3.2. Best basis algorithm. The collection of all anisotropic Haar bases and fast transforms makes for a potentially very useful library. It contains bases

associated with partitions which subdivide much more finely in i_1 than in i_2 in some subsets of $[0, n]^2$ and more finely in i_2 than in i_1 in other subsets. There is therefore the possibility of finding bases very well adapted to certain anisotropic problems.

How to choose a “best-adapted” basis? In the general framework set up in the context of cosine packet–wavelet packet bases by Coifman and Wickerhauser (1992), one specifies an additive “entropy” functional of the vector $\theta \in R^N$,

$$(3.5) \quad \mathcal{E}(\theta) = \sum_{i=1}^N e(\theta_i),$$

where $e(t)$ is a scalar function. Coifman and Wickerhauser’s original proposal was $e_{CW}(t) = -t^2 \log(t^2)$, but $e_p(t) = |t|^p$, where $0 < p < 2$ also makes sense, as well as other choices—see below. One uses such a functional to evaluate the quality of a basis; if $\theta(f, \mathcal{B})$ denotes the vector of coefficients of the object f in basis \mathcal{B} , then $\mathcal{E}(\theta(f, \mathcal{B}))$ is a measure of the usefulness of a basis for representing f , and the best basis \mathcal{B} in a library \mathcal{L} of ortho bases solves the problem

$$(3.6) \quad \min_{\mathcal{B} \in \mathcal{L}} \mathcal{E}(\theta(f, \mathcal{B})).$$

In the specific case of interest, there are as many bases in the library as there are complete recursive partitions. Elementary arguments show that the number of bases is exponential in N .

While such exponential behavior makes brute force calculation of the optimum in (3.6) practically impossible, judicious application of dynamic programming gives a practical algorithm.

In order to express the key analytic feature of the objective functional, we take a more general point of view, and consider the problem of finding an optimal basis for *just the data falling in the dyadic rectangle* R . Each complete recursive dyadic partition of R , $\mathcal{P}^*(R)$ say, leads to an anisotropic Haar basis, $\mathcal{B}(R)$ say, for the collection of $n \times n$ arrays supported only in R . Hence we can define the optimization problem

$$[\text{BOB}(R)] \quad \hat{\mathcal{B}}(R) = \underset{\mathcal{B}(R)}{\operatorname{argmin}} \tilde{\mathcal{E}}(\theta(y, \mathcal{B}(R))).$$

Here $\theta(y, \mathcal{B}(R))$ refers to the coefficients in an anisotropic basis for $l^2(R)$, and $\tilde{\mathcal{E}}(\theta) = \sum_{i=2}^{\dim(\theta)} e(\theta_i)$ refers to a relative entropy, which ignores the first coordinate. We let $\hat{\mathcal{P}}^*(R)$ denote the corresponding optimal complete recursive dyadic partition of R .

Solutions to $\text{BOB}(R)$ have a key inheritance property. Let R be a dyadic rectangle and suppose it has both vert-children and horiz-children. Then the optimal basis of R is generated by a complete recursive dyadic partition $\hat{\mathcal{P}}^*(R)$ formed in one of two ways. This partition is either (1) the union of optimal partitions of the horiz-children $\hat{\mathcal{P}}^*(R^{1,0}) \cup \hat{\mathcal{P}}^*(R^{1,1})$, or (2) the union of optimal partitions of the vert-children $\hat{\mathcal{P}}^*(R^{2,0}) \cup \hat{\mathcal{P}}^*(R^{2,1})$. Which of these two cases holds can be determined by holding a “tournament,” selecting

the winner as the smallest of the numbers

$$\text{BOB}(R^{1,0}) + \text{BOB}(R^{1,1}) + e_1, \text{BOB}(R^{2,0}) + \text{BOB}(R^{2,1}) + e_2,$$

where $e_i = e(\alpha_R^i)$.

The exception to this rule is of course at the finest scale: a 2×1 or 1×2 rectangle has only one complete recursive partition, and no tournament is necessary to select a “best” one.

By starting from the next-to-finest scale and applying the inheritance property, we can get the optimal partitions of all 4×1 , of all 2×2 and of all 1×4 rectangles (omitting again the tournament for 4×1 and 1×4 rectangles) and so on. Continuing in a fine-to-coarse or ‘bottom-up’ fashion, we eventually get to the coarsest level and obtain an optimal partition for $[0, n]^2$.

Once again there are approximately $4n^2 = 4N$ dyadic rectangles. Each dyadic rectangle is visited once in the main part of the algorithm, and there are at most a certain constant number C of additions and multiplications per visit. The total work is less than or equal to $C4N$ flops and less than or equal to $4N$ storage locations. See the appendix in Donoho (1995b) for a formal description of the algorithm.

4. Best basis denoising. CART has to do with removing noise from the data y to produce a reconstruction \hat{f} approximating the noiseless data f . The philosophy of BOB is much less specific: it may be used for many purposes, for example, in data compression and for fast numerical analysis [Coifman, Meyer, Quake and Wickerhauser (1994)]. The application determines the choice of entropy, and the use of the expansion in the best basis.

To use best-basis ideas for noise removal, one could apply the proposals of Donoho and Johnstone (1994b). Define

$$(4.1) \quad \mathcal{E}_\lambda(\theta) = \sum_1^N \min(\theta_i^2, \lambda^2 \sigma^2)$$

and obtain an optimal basis

$$(4.2) \quad \hat{\mathcal{B}} = \min_{\mathcal{B} \in \mathcal{L}} \mathcal{E}_\lambda(\theta(y, \mathcal{B})).$$

Then apply hard thresholding in the selected basis, at threshold level $\lambda\sigma$:

$$(4.3) \quad \hat{\theta}_i = \theta(y, \hat{\mathcal{B}})_i \mathbf{1}_{\{|\theta(y, \hat{\mathcal{B}})_i| > \lambda\sigma\}}, \quad 1 \leq i \leq N.$$

Reconstruct object \hat{f} having coefficients $\hat{\theta}$ in basis $\hat{\mathcal{B}}$. This is the denoised object.

Donoho and Johnstone (1994b) developed results, to be discussed in Section 7, showing that with an appropriate choice of λ , the empirical basis chosen by this scheme was *near ideal*.

In the current setting, where \mathcal{L} is the library of anisotropic Haar bases, (4.2) is amenable to treatment by the fast best-basis algorithm of the last section, so it may be computed in order N time. This denoising estimate,

while possessing certain nice characteristics, lacks one of the attractive features of CART: an interpretation as a spatially adaptive averaging method. Such a spatially adaptive method would have the form

$$\hat{f}(\mathbf{i}) = \sum_{R \in \mathcal{P}} \langle y, \chi_R \rangle \chi_R(\mathbf{i}),$$

giving a piecewise constant reconstruction based on rectangular averages of the noisy data y over rectangles R . Here the partition $\mathcal{P} = \mathcal{P}(y)$ would be chosen data adaptively, and once the partition were chosen, the reconstruction would take a simple form of averaging. While we will mention this procedure further and use its properties, we mention it now only to show that threshold denoising in a best-ortho-basis is not identical to CART.

5. Tree constraints in the one-dimensional Haar system. In the context of the ordinary one-dimensional Haar transform, Engel (1994) has shown that a special type of reconstruction in the Haar system can be related to recursive partitioning. Let, temporarily, $y = (y_i)_{i=0}^{n-1}$ and suppose $y_i = g(i) + v_i$, with v_i noise. Consider reconstructions \hat{g} of the form

$$(5.1) \quad \hat{g}(i) = \bar{y} + \sum_I w_I \langle y, h_I \rangle h_I(i),$$

where the sum is over dyadic subintervals of $[0, n)$ and the w_I are scalar “weights.” Now impose on the weights (w_I) two constraints.

1. [Tree-i]. *Keep-or-kill.* Each weight is 1 or 0.
2. [Tree-ii]. *Hereditiy.* w_I can be 1 only if also $w_{I'} = 1$ whenever $I \subset I'$. If $w_I = 0$, then $w_{I'} = 0$ for every $I' \subset I$.

Each set of weights satisfying these constraints selects the nodes of a dyadic tree T . Engel has called such constraints tree constraints and shown that reconstructions obeying these constraints may be put in the form of spatial averages.

THEOREM 5.1 [Engel (1994)]. *Suppose that \hat{g} defined by (5.1) obeys the tree constraints (Tree-i) and (Tree-ii). Say that I is terminal if $w_I = 1$ but every interval $I' \subset I$ has $w_{I'} = 0$. The collection of terminal intervals forms a partition \mathcal{P} , and*

$$(5.2) \quad \hat{g}(i) = \sum_{I \in \mathcal{P}} \langle y, \chi_I \rangle \chi_I(i).$$

6. Hereditary constraints and CART. Tree constraints make sense also in the setting of two-dimensional anisotropic Haar bases. We consider reconstructions

$$(6.1) \quad \hat{f}(\mathbf{i}) = \bar{y} + \sum_{R \in NT(\mathcal{P}^*)} w_R \langle y, \phi_R^{s(R)} \rangle \phi_R^{s(R)}(\mathbf{i}),$$

where \mathcal{P}^* is a complete recursive dyadic partition, $\{\chi_{[0, n)^2}, (\phi_R^{s(R)})\}$ the associated orthogonal basis, and the weights (w_R) obey two *hereditary constraints*

1. [Hered-i]. *Keep-or-kill*. Each weight w_R is 0 or 1.
2. [Hered-ii]. *Heredity*. $w_R = 1$ implies $w_{R'} = 1$ for all ancestors R' of R in \mathcal{P}^* ; $w_R = 0$ implies $w_{R'} = 0$ for all descendants of R in \mathcal{P}^* . We state without proof the analog of Engel’s theorem.

THEOREM 6.1. *The reconstruction \hat{f} obeying (6.1), [Hered-i], and [Hered-ii] has precisely the form*

$$\hat{f}(\mathbf{i}) = \sum_{R \in \mathcal{P}} \langle y, \chi_R \rangle \chi_R(\mathbf{i})$$

for some possibly incomplete recursive dyadic partition \mathcal{P} .

Three questions arise naturally about reconstructions obeying hereditary constraints.

- Q1. How can the best hereditary reconstruction in a given basis be found?
- Q2. How can the basis in which hereditary reconstruction works best be found?
- Q3. How can the hereditary best-basis be efficiently calculated?

All three questions have attractive answers.

6.1. *Best hereditary reconstruction in given basis.* Let T^* denote the complete binary tree of depth $\log_2(N)$. Identifying subtrees $T \subset T^*$ with sequences of weights (w_R) obeying [Hered-i]-[Hered-ii], we write $\hat{f}_{\mathcal{B}, T}$ for the reconstruction (6.1) in basis \mathcal{B} having weights (w_I) associated with the tree T .

We define the “best” hereditary reconstruction in terms of the hereditary CPRSS

$$(6.2) \quad \text{CPRSS}(T; \lambda, \mathcal{B}) = \|y - \hat{f}_{\mathcal{B}, T}\|_{l^2}^2 + \lambda \#(T).$$

The optimization problem is the one achieving the minimal CPRSS among all such reconstructions:

$$(6.3) \quad \min_{T \subset T^*} \text{CPRSS}(T; \lambda, \mathcal{B}).$$

By orthogonality of the basis \mathcal{B} , we can reformulate this in terms of coordinates. Let $\theta = \theta(y, \mathcal{B})$ denote the vector of coordinates and $(w_R \theta_R(y, \mathcal{B}))$ denote the same vector after applying weights w_R associated with the subtree T . Then we have the following equivalent form of (6.2):

$$\text{CPRSS}(T) = \sum_R ((w_R - 1)^2 \theta_R^2 + \lambda w_R).$$

This quantity has an inheritance property, which we express as follows. Let $T^*(R)$ denote the complete tree of depth $\log_2(\#R)$ and define the optimization

problem

$$[\text{Hered}(R)] \quad \min_{T \subset T^*(R)} \sum_{R'} ((w_{R'} - 1)^2 \theta_{R'}^2 + \lambda w_{R'}).$$

The optimization problem implicitly defines an optimal subtree $\hat{T}(R)$. The inheritance property: the optimal subtree $\hat{T}(R)$ is a function of the optimal subtrees of the children problems $\hat{T}(R^{s(R), b})$, $b = 0, 1$. The tree $\hat{T}(R)$ is either the empty subtree, or else it has $\hat{T}(R^{s(R), b})$ as subtrees joined at root($\hat{T}(R)$).

It follows by this inheritance property that the optimal subtree may be computed by a bottom-up pruning exactly as in the optimal pruning algorithm of CART, Algorithm 10.1, page 294 of the CART book. Hence, a minimizing subtree may be found in order N time. A formal statement of the algorithm is given in the appendix of Donoho (1995b).

6.2. *Best basis for hereditary reconstruction.* We can define the quality of a basis for hereditary reconstruction by considering the optimum value of the CPRSS functional over all hereditary reconstructions in that basis. Hence, define the hereditary entropy

$$(6.4) \quad \mathcal{H}_\lambda(\mathcal{B}) = \min_{T \subset T^*} \text{CPRSS}(T; \lambda, \mathcal{B}).$$

A best basis for hereditary reconstruction is then the solution of

$$(6.5) \quad \min_{\mathcal{B} \subset \mathcal{L}} \mathcal{H}_\lambda(\mathcal{B}),$$

where \mathcal{L} is a library of orthogonal bases. This may be motivated in two ways. First, the goal is intrinsically reasonable, as it seeks a best tradeoff, over all bases and all subtrees, of complexity $\#(T)$ against fidelity to the data $\|y - \hat{f}_{\mathcal{B}, T}\|_2^2$. Second, we will prove below that the reconstruction obtained in the optimum basis has a near-ideal mean-squared error.

6.3. *Fast algorithm via CART.* The entropy $\mathcal{H}_\lambda(\mathcal{B})$ is not an additive functional $\sum_{i=1}^N e(\theta_i(y, \mathcal{B}))$ of the coordinates of y in basis \mathcal{B} . Therefore the best-basis algorithm of Section 3, strictly speaking, does not apply. Luckily, we can use the fast CART algorithm. By now this is obvious; we summarize this fact formally, though without writing out the proof.

THEOREM 6.2. *When λ is the same in both, CART and BOB with hereditary constraints have the same answers. More precisely,*

$$(6.6) \quad \min_{\mathcal{B} \in \mathcal{L}} \mathcal{H}_\lambda(\mathcal{B}) = \min_{\mathcal{P}} \text{CPRSS}(\mathcal{P}, \lambda).$$

The solution of the best-basis problem (6.5) gives, explicitly, an anisotropic basis $\hat{\mathcal{B}}$ and, implicitly by (6.3), an optimal subtree \hat{T} ; the solution of the CART problem (2.4) gives an optimizing partition $\hat{\mathcal{P}}$, and we have

$$\hat{f}_{\hat{\mathcal{B}}, \hat{T}}(\cdot) = \hat{f}(\cdot | \hat{\mathcal{P}}).$$

REMARK 1. Although $\mathcal{H}_\lambda(\mathcal{B})$ is not additive, a fast algorithm for computing it is available—the dyadic CART algorithm of Section 1. This shows that fast best-basis algorithms may exist for certain nonadditive entropies.

REMARK 2. Although CART and best-ortho-basis are *not* the same in general, *in this case*, with a specific set of definitions of best-ortho-basis and a specific set of restrictions on the splits employed by CART, the two methods are the same.

7. Oracle inequalities. CART and BOB define objects which are the solutions of certain optimization problems and hence are in some sense “optimal.” However, we should stress that they are optimal only in the very artificial sense that they solve certain optimization problems we have defined.

We now turn to the question of performance according to externally defined standards, which will lead ultimately to a proof of our main result. This will entail a certain kind of near optimality with a more significant and useful meaning.

In accordance with the philosophy laid out in Donoho (1995a), we approach this in two stages. First, there is a statistical decision theory component of the problem which we deal with in Section 7; second, there is a harmonic analysis component of the problem, which we deal with in Section 8.

7.1. *Oracle inequalities.* Once more we are in the model (2.1), and we wish to recover f with small mean-squared error. We evaluate an estimator $\hat{f} = \hat{f}(y)$ by its risk

$$R(\hat{f}, f) = \text{MSE}(\hat{f}(y), f).$$

Suppose we have a collection of estimators $\hat{\Phi} = \{\hat{f}(\cdot)\}$; we wish to use the one best adapted to the problem at hand. The best performance we can hope for is what Donoho and Johnstone (1994c) call the *ideal* risk:

$$\mathcal{R}^*(\hat{\Phi}, f) = \inf\{R(\hat{f}, f) : \hat{f} \in \hat{\Phi}\}.$$

We call this ideal because it can be attained only with an oracle, who in full knowledge of the underlying f (but not revealing this to us) selects the best estimator for this f from the collection $\hat{\Phi}$.

We optimistically propose $\mathcal{R}^*(\hat{\Phi}, f)$ as a target, and seek true estimators which can approach this target. It turns out that in several examples, one can find estimators which achieve this to within logarithmic terms. The inequalities which establish this are of the form

$$R(\hat{f}^*, f) \leq \text{Const} \cdot \log(N) (\sigma^2 + \mathcal{R}^*(\hat{\Phi}, f)) \quad \forall f,$$

which Donoho and Johnstone (1994a,b,c) call *oracle inequalities*, because they compare the risk of valid procedures with the risk achievable by idealized procedures which depend on oracles.

7.2. *Example: keep-or-kill de-noising.* Suppose we are operating in a fixed orthogonal basis \mathcal{B} and consider the family $\hat{\Phi}$ of estimators defined by keeping or killing empirical coefficients in the basis \mathcal{B} . Such estimators $\hat{f}(y; w)$ are given in the basis \mathcal{B} by

$$\theta_i(\hat{f}, \mathcal{B}) = w_i \theta_i(y, \mathcal{B}), \quad i = 1, \dots, N,$$

where each weight w_i is either 0 or 1. Such estimators have long been considered in the context of Fourier series estimation, where the basis is the Fourier basis, the coefficients are Fourier coefficients, and the w_i are 1 only for $1 \leq i \leq k$ for some frequency cutoff k . Estimators of this form have also been considered by Donoho and Johnstone (1994c) in the context where \mathcal{B} is a wavelet basis; in that setting the unit weights are ideally chosen at sites of important spatial variability.

Formally then $\hat{\Phi} = \{\hat{f}(\cdot; w) : w \in \{0, 1\}^N\}$ is the collection of all keep-or-kill estimators in the fixed basis \mathcal{B} . Donoho and Johnstone (1994c) studied the nonlinear estimator \hat{f}^* , defined in the basis \mathcal{B} by hard thresholding

$$\theta_i(\hat{f}^*(y), \mathcal{B}) = \eta_{\sqrt{\lambda_N}}(\theta_i(y, \mathcal{B})), \quad i = 1, \dots, N,$$

where $\eta_t(y) = 1_{\{|y| > t\}}(y) \operatorname{sgn}(y)$ is the hard thresholding nonlinearity and $\lambda_N = \sigma^2 2 \log(N)$. They showed that \hat{f}^* obeys the oracle inequality

$$R(\hat{f}^*, f) \leq (2 \log(N) + 1)(\sigma^2 + \mathcal{R}^*(\hat{\Phi}, f)) \quad \forall f,$$

as soon as $N \geq 4$. In short, simple thresholding comes within log terms of ideal keep-or-kill behavior.

The reader will find it instructive to note that the estimator \hat{f}^* can also be defined as the solution of the optimization problem

$$\min_{w \in \{0, 1\}^N} \|y - \hat{f}(y; w)\|^2 + \lambda_N \#\{i : w_i \neq 0\}.$$

This is, of course, a complexity penalized RSS, with penalty term λ_N . Thus the near-ideal estimator is the solution of a minimum CPRSS principle.

7.3. *Example: best-basis denoising.* Suppose now we are operating in a library \mathcal{L} of orthogonal bases \mathcal{B} and consider the family $\hat{\Phi}$ of estimators defined by keeping or killing empirical coefficients in *some* basis $\mathcal{B} \in \mathcal{L}$. Such estimators $\hat{f}(y; w, \mathcal{B})$ are of the form

$$\theta_i(\hat{f}, \mathcal{B}) = w_i \theta_i(y, \mathcal{B}), \quad i = 1, \dots, N,$$

where each weight w_i is either 0 or 1.

Formally $\hat{\Phi} = \{\hat{f}(\cdot; w, \mathcal{B}) : w \in \{0, 1\}^N, \mathcal{B} \in \mathcal{L}\}$. For obvious reasons, we call $\mathcal{R}^*(\hat{\Phi}, f)$ also $\mathcal{R}^*(\text{IDEAL BASIS}, f)$. Donoho and Johnstone (1994b) developed a nonlinear estimator \hat{f}^* , with near-ideal properties; it is precisely the best-basis denoising estimator defined in Section 4; see (4.1)–(4.3). In detail they supposed that among all bases in the library there are at most M

distinct elements. They suppose that we pick $\xi > 8$, and set $t = \sqrt{2 \log_e(M)}$; then with $\lambda = (\xi(1 + t))^2$, they prove a result almost as strong as the following, which we prove in the Appendix.

THEOREM 7.1. *For an appropriate constant $A(\xi)$, the BOB estimator obeys the oracle inequality*

$$(7.1) \quad R(\hat{f}^*, f) \leq A(\xi) \lambda (\sigma^2 + \mathcal{R}^*(\hat{\Phi}, f)) \quad \forall f.$$

In short, empirical best basis (with an appropriate entropy) comes within log terms of ideal keep-or-kill behavior in an ideal basis. In the specific case of the library of anisotropic Haar bases, $M = 4N$, and so for a fixed choice of ξ , (7.1) becomes

$$R(\hat{f}^*, f) \leq \text{Const} \cdot \log(N) (\sigma^2 + \mathcal{R}^*(\text{IDEAL BASIS}, f)) \quad \forall f.$$

7.4. Example: CART. Oracle inequalities for CART are now easy to state. Suppose now we are operating in the library \mathcal{L} of anisotropic Haar bases and consider the family $\hat{\Phi}_{\text{Tree}}$ of *hereditary linear estimators*, that is, estimators defined by keeping or killing the empirical coefficients in some basis $\mathcal{B} \in \mathcal{L}$, where the coefficients that are kept fall in a tree pattern T . Such estimators $\hat{f}(y; T, \mathcal{B})$ are of the form

$$\theta_i(\hat{f}, \mathcal{B}) = w_i \theta_i(y, \mathcal{B}), \quad i = 1, \dots, N,$$

where each weight w_i is either 0 or 1, and the nonzero w form a tree.

Formally let $\hat{\Phi}_{\text{Tree}} = \{\hat{f}(\cdot; T, \mathcal{B}) : T \subset T^* \mathcal{B} \in \mathcal{L}\}$ be the collection of all hereditary linear estimators in any anisotropic Haar basis. The ideal risk $\mathcal{R}^*(\hat{\Phi}_{\text{Tree}}, f)$ is just the risk of CART applied in an ideal partition selected by an oracle \mathcal{P} . So call this $\mathcal{R}^*(\text{IDEAL CART}, f)$.

Consider now the dyadic CART estimator \hat{f}^* defined with λ exactly as specified in the best-basis denoising setting of the Section 13. So for $\xi > 8$, set $\lambda = (\xi(1 + \sqrt{2 \log_e(4N)}))^2$. We prove the following in the Appendix.

THEOREM 7.2. *For all $N \geq 1$, the dyadic CART estimator obeys the oracle inequality*

$$R(\hat{f}^{*,\lambda}, f) \leq \text{Const} \cdot \log(N) (\sigma^2 + \mathcal{R}^*(\text{IDEAL CART}, f)) \quad \forall f.$$

In short, empirical dyadic CART (with an appropriate penalization) comes within log terms of ideal dyadic CART.

8. Anisotropic smoothness spaces. We now change gears slightly and consider harmonic analysis questions. Specifically we are going to show that anisotropic Haar bases are particularly well adapted to dealing with classes of functions having anisotropic smoothness.

We denote now by f a function $f(x, y)$ defined on $[0, 1]^2$, rather than an array of pixel values. We consider objects of possibly different smoothnesses

in different directions. Define the finite difference operators $(D_h^1 f)(x, y) = f(x + h, y) - f(x, y)$ and $(D_h^2 f)(x, y) = f(x, y + h) - f(x, y)$. For δ_1, δ_2 satisfying $0 \leq \delta_i \leq 1$, define the anisotropic smoothness class

$$\mathcal{F}_p^{\delta_1, \delta_2}(C) = \left\{ f: \|f\|_p \leq C, \|D_h^1 f\|_{L^p(Q_h^1)} \leq Ch^{\delta_1}, h \in (0, 1), \right. \\ \left. \|D_h^2 f\|_{L^p(Q_h^2)} \leq Ch^{\delta_2}, h \in (0, 1) \right\},$$

where $Q_h^1 = [0, 1 - h] \times [0, 1]$ and $Q_h^2 = [0, 1] \times [0, 1 - h]$. This contains objects of genuinely anisotropic smoothness whenever $\delta_1 \neq \delta_2$. The usual smoothness spaces (Hölder, Sobolev, Triebel, etc.) involve equal degrees of smoothness in different directions and are sometimes called “isotropic,” so that classes like $\mathcal{F}_p^{\delta_1, \delta_2}(C)$ would be called “anisotropic.” Spaces of this kind were introduced by Nikol’skii; for information see Nikol’skii (1969) and Temlyakov (1993).

8.1. *Spatially uniform anisotropic bases.*

DEFINITION 8.1. A *sequential partitioning of j into two parts* is a pair of series of integers $j_1(j), j_2(j), j = 0, 1, 2, \dots$ obeying

- (i) Initialization: $j_1(0) = j_2(0) = 0$;
- (ii) Partition: $j_1(j) + j_2(j) = j$;
- (iii) Sequential allocation;

$$j_1(j) = j_1(j - 1) + b_1(j), \quad b_1 \in \{0, 1\}; \\ j_2(j) = j_2(j - 1) + b_2(j), \quad b_2 = 1 - b_1.$$

We can think of two boxes and a sequential scheme where at each stage we put a ball in one of the two boxes. The expression $j_i(j)$ represents the number of balls in box i at stage j , and $b_1 = 1 - b_2$ represents the constraint that only one ball is put into the boxes at each stage.

DEFINITION 8.2. Consider a sequential partition of j into two parts. The *spatially uniform alternating partition* subordinate to this partition— $\text{SUAP}(j_1, j_2)$ —is a complete dyadic recursive partition formed in a homogeneous fashion: at stage 1, the square $[0, 1]^2$ is split horizontally if $b_1(1) = 1$, and vertically if $b_2(1) = 1$; at stage 2, each of the two resulting rectangles is split in two, horizontally if $b_1(2) = 1$, vertically if $b_2(2) = 1$; and at stage j , each of the 2^{j-1} rectangles of volume 2^{-j+1} formed at the previous stage is split vertically if $b_1(j) = 1$, horizontally if $b_2(j) = 1$.

The recursive partition $\text{SUAP}(j_1, j_2)$ defines a series of collections of rectangles: $\mathcal{R}(0)$ consists of the root rectangle, $\mathcal{R}(1)$ consists of the two children of the root, $\mathcal{R}(2)$ of the four children of the rectangles in $\mathcal{R}(1)$, and so on. In general, $\mathcal{R}(j)$ consists of 2^j rectangles of area 2^{-j} each.

This sequence of rectangles defines an orthogonal basis of $L^2([0, 1]^2)$ in a fashion similar to the discrete case, with fairly obvious changes due to the

change in setting. Let now I denote a dyadic subinterval of $[0, 1]$ and $\tilde{\chi}_I(x)$ be the “same function” as $\chi_I(i)$, under the correspondence $x_i \leftrightarrow i/n$ and under the different choice of normalizing measure $\tilde{\chi}_I(x) = 1_I(x)l(I)^{-1/2}$ where $l(I)$ denotes the length of I . Similarly, let $\tilde{h}_I(x) = (1_{I^1}(x) - 1_{I^0}(x))l(I)^{-1/2}$. Then set $\varphi_R^1 = \tilde{h}_{I_1}(x)\tilde{\chi}_{I_2}(y)$, $\varphi_R^2 = \tilde{\chi}_{I_1}(x)\tilde{h}_{I_2}(y)$. Then set

$$\begin{aligned} \xi_0 &= \chi_{[0,1]^2}, & \xi_{0,0} &= \varphi_{[0,1]^2}^{b_2(1)}, \\ \xi_{1,R} &= \varphi_R^{b_2(2)} & \text{for } R \in \mathcal{R}(1), \\ \xi_{2,R} &= \varphi_R^{b_2(3)} & \text{for } R \in \mathcal{R}(2), \end{aligned}$$

and in general

$$(8.1) \quad \xi_{j,R} = \varphi_R^{b_2(j)} \quad \text{for } R \in \mathcal{R}(j);$$

call this the *spatially homogeneous anisotropic basis* SHAB (j_1, j_2) .

The coefficients of f in this basis are

$$(8.2) \quad \tilde{f} = \text{ave}_{[0,1]^2}, \quad \alpha_R = \langle \xi_{j,R}, f \rangle, \quad R \in \mathcal{R}(j).$$

8.2. *Best basis for a functional class.* Donoho (1993, 1996) described a notion of best-orthogonal-basis for a functional class \mathcal{F} , which describes the kinds of bases in which certain kinds of de-noising and data compression can best be conducted. For this notion, a best basis for a functional class \mathcal{F} is any basis in which the rearranged coefficients of members of \mathcal{F} decay fastest. According to this definition, one-dimensional wavelet bases are best bases for classes like bounded variation, Sobolev, Triebel and Besov classes; Wilson bases are best bases for modulation spaces; Fourier bases are best bases for L^2 -Sobolev spaces, and so on. More generally, spaces with an unconditional basis have the unconditional basis as best basis. Certain spaces have best bases which are not unconditional bases. Kashin (1985) showed, in our language, that sinusoids give a best orthogonal basis for appropriate classes of Hölder continuous functions. Kashin and Temlyakov (1994), among other things, discussed best bases for spaces with bounded mixed derivatives.

For a vector θ in sequence space, let $|\theta|_{(i)}$ denote the rearranged magnitudes of the coefficients, sorted in decreasing order $|\theta|_{(1)} \geq |\theta|_{(2)} \geq \dots$. The *weak l^τ norm* measures the decay of these by

$$\|\theta\|_{\text{wl}^\tau} = \sup_{k \geq 1} k^{1/\tau} |\theta|_{(k)}.$$

This measures decay of the coefficients since $\|\theta\|_{\text{wl}^\tau} \leq C$ implies $|\theta|_{(k)} \leq Ck^{-1/\tau}$, $k = 1, 2, \dots$.

Now, with $\theta = (\theta_i(f, \mathcal{B}))$, the coefficients of f in an orthogonal basis \mathcal{B} , a given functional class \mathcal{F} maps to a coefficient body $\Theta(\mathcal{F}, \mathcal{B}) = \{(\theta_i(f, \mathcal{B}))_i; f \in \mathcal{F}\}$. For such a set Θ , we say $\Theta \subset \text{wl}^\tau$ if

$$\sup\{\|\theta\|_{\text{wl}^\tau}; \theta \in \Theta\} < \infty.$$

DEFINITION 8.3. We call the critical exponent of Θ the number $\tau^*(\Theta)$ obtained as the infimum of all τ for which $\Theta \subset \text{wl}^\tau$.

Assuming that \mathcal{F} is a subset of L^2 , $0 \leq \tau^*(\Theta) \leq 2$.

From this point of view a best basis for \mathcal{F} is any basis \mathcal{B}^* which minimizes the critical exponent

$$(8.3) \quad \tau^*(\Theta(\mathcal{F}, \mathcal{B}^*)) = \min_{\mathcal{B}} \tau^*(\Theta(\mathcal{F}, \mathcal{B})).$$

In such a basis the rearranged coefficients will be the most rapidly decaying among all ortho bases.

8.3. *Best anisotropic bases.* With this background, it is interesting to ask about the decay properties of coefficients in different spatially homogeneous anisotropic bases. We will identify a basis \mathcal{B} within the class of anisotropic Haar bases satisfying (8.3) among all bases. The key fact is this upper bound.

LEMMA 8.4. Let $f \in \mathcal{F}_p^{\delta_1, \delta_2}(C)$, where $1/p < \rho + 1/2$, where $\rho = \delta_1 \delta_2 / (\delta_1 + \delta_2)$. If $b_1(j) = 1$,

$$(8.4) \quad \left(\sum_{R \in \mathcal{R}(j)} |\alpha_R|^p \right)^{1/p} \leq C 2^{-j_1 \delta_1} (2^{-j})^{1/2-1/p},$$

while if $b_2(j) = 1$,

$$(8.5) \quad \left(\sum_{R \in \mathcal{R}(j)} |\alpha_R|^p \right)^{1/p} \leq C 2^{-j_2 \delta_2} (2^{-j})^{1/2-1/p}.$$

Now a choice of a spatially uniform anisotropic partition which would make optimal use of these expressions as a function of j would arrange things so that the decrease of the largest of the two expressions went fastest in j . Thus optimal use of Lemma 8.4 leads to the problem of constructing a sequential partition of j into parts that optimize the rate of decay of

$$(8.6) \quad \max(2^{-j_2(j)\delta_1}, 2^{-j_2(j)\delta_2})$$

as a function of j .

There is an obvious limit on how well this can be done. Consider optimizing (8.6) subject only to the constraints $j_1(j) + j_2(j) = j$ and $j_i \geq 0$, that is, without imposing the requirement that the j_i be integers, or be sequentially chosen. The solution is $j_1(j) = \delta_2 / (\delta_1 + \delta_2) j$ and $j_2(j) = \delta_1 / (\delta_1 + \delta_2) j$, achieving an optimally small value of

$$(8.7) \quad 2^{-j\rho}, \quad \rho = \frac{\delta_1 \delta_2}{\delta_1 + \delta_2},$$

in (8.6). We cannot hope to do better than this, once we reimpose the constraints associated with a sequential partition. But we can come close.

DEFINITION 8.5. For a given pair of “exponents” δ_1, δ_2 obeying $0 < \delta_i \leq 1$, we call an *optimal sequential partition* of j a sequential partition (j_1^*, j_2^*) obtained as follows:

- (i) Start from $j_1^*(0) = j_2^*(0) = 0$;
- (ii) At stage $j + 1$, allocate $b_1(j)$ and $b_2(j)$ as follows:
 - (a) If $j_1^*(j - 1)\delta_1 = j_2^*(j - 1)\delta_2$ allocate the ball to whichever box has the smaller exponent: $b_1(j) = 1$ if $\delta_1 < \delta_2$;
 - (b) If $j_1^*(j - 1)\delta_1 \neq j_2^*(j - 1)\delta_2$ allocate the ball to whichever box has the smaller product $j_i^*(j - 1)\delta_i$.

This so-called optimal sequential partitioning of j is a greedy stepwise minimization of objective (8.6). It turns out that it is near optimal, even among nonsequential partitions.

LEMMA 8.6. For $0 < \delta_1, \delta_2 \leq 1$ and $\rho = (\delta_1 \delta_2) / (\delta_1 + \delta_2)$,

(8.8)
$$\max(2^{-j_1^*(j)\delta_1}, 2^{-j_2^*(j)\delta_2}) \leq 2 \cdot 2^{-j\rho}.$$

Due to (8.7), this is essentially optimal within the class of sequential partitions of j . Inspired by this, we propose the following definition.

DEFINITION 8.7. We call the *best anisotropic basis* $BAB(\delta_1, \delta_2)$ the anisotropic basis of $L^2[0, 1]^2$, defined using $SHAB(j_1^*, j_2^*)$.

Combining Lemma 8.6 with Lemma 8.4 above, we have the following corollary.

COROLLARY 8.8. If we use the $BAB(\delta_1, \delta_2)$ then for $\rho = (\delta_1 \delta_2) / (\delta_1 + \delta_2)$ and $1/p < \rho + 1/2$,

(8.9)
$$\left(\sum_{R \in \mathcal{R}(j)} |\alpha_R|^p \right)^{1/p} \leq 2C 2^{-j(\rho + 1/2 - 1/p)}.$$

8.4. *Optimality of BAB.* Armed with Corollary 8.8, it is possible to justify Definition 8.7 and prove that $BAB(\delta_1, \delta_2)$ is an optimal basis in the sense of Section 8.2.

THEOREM 8.9. Let \mathcal{L} denote the collection of all orthogonal bases for $L^2[0, 1]^2$.

(8.10)
$$\tau^*\left(\Theta(\mathcal{F}_p^{\delta_1, \delta_2}(C), BAB(\delta_1, \delta_2))\right) = \min_{\mathcal{B} \in \mathcal{L}} \tau^*\left(\Theta(\mathcal{F}_p^{\delta_1, \delta_2}(C), \mathcal{B})\right).$$

The proof is a consequence of three lemmas, all of which are proved in the Appendix. The first gives an evaluation of the critical exponent for $BAB(\delta_1, \delta_2)$.

LEMMA 8.10. If $1/p < \rho + 1/2$, then

(8.11)
$$\tau^*\left(\Theta(\mathcal{F}_p^{\delta_1, \delta_2}(C), BAB(\delta_1, \delta_2))\right) = 2/(2\rho + 1).$$

For comparison we briefly mention results obtainable in other bases. Suppose we use the isotropic basis $\text{SHAB}(j_1^+, j_2^+)$ defined by $j_1^+(j) = \lfloor (j+1)/2 \rfloor$, $j_2^+(j) = \lfloor j/2 \rfloor$. This has an equal frequency of splitting in each direction. By a side calculation,

$$\tau^*(\Theta(\mathcal{F}_p^{\delta_1, \delta_2}(C), \text{SHAB}(j_1^+, j_2^+))) = 2/(2\rho^+ + 1),$$

where $\rho^+ = \min(\delta_1, \delta_2)/2$. As

$$\frac{\rho}{\rho^+} = \frac{\max(\delta_1, \delta_2)}{\text{ave}(\delta_1, \delta_2)},$$

analysis in the isotropic basis yields coefficients with slower decay rate than optimal whenever the smoothness class is genuinely anisotropic, that is, whenever $\delta_1 \neq \delta_2$.

The optimality of the exponent in (8.11) among *all* bases follows from a lower bound technique developed at greater length in Donoho (1996). First, a definition: an *orthogonal hypercube* \mathcal{H} of dimension m and side ε is a collection of all sums $g_0 + \sum_{i=1}^m \alpha_i g_i$ where the g_i are orthonormal functions and the $|\alpha_i| \leq \varepsilon$.

LEMMA 8.11. *Suppose \mathcal{F} contains a sequence \mathcal{H}_j of orthogonal hypercubes of dimension m_j and side ε_j where $\varepsilon_j \rightarrow 0$, $m_j \rightarrow \infty$,*

$$m_j^{1/\tau} \varepsilon_j \geq C_0 > 0.$$

Let \mathcal{L} denote any collection of orthogonal bases.

$$\inf_{\mathcal{B} \in \mathcal{L}} \tau^*(\Theta(\mathcal{F}, \mathcal{B})) \geq \tau.$$

LEMMA 8.12. *Each class $\mathcal{F}_p^{\delta_1, \delta_2}(C)$ contains a sequence \mathcal{H}_j of orthogonal hypercubes of dimension $m_j = 2^j$ and side ε_j where $\varepsilon_j \rightarrow 0$, $m_j \rightarrow \infty$,*

$$m_j^{1/\tau} \varepsilon_j \geq KC,$$

with K a fixed constant, and $\tau = 2/(2\rho + 1)$.

A related result was developed by Kashin (1985), who showed, in our language, that sinusoids make a best orthogonal basis for Hölder- α spaces. Kashin's approach to that result uses something like Lemma 8.11, proved by different means; he uses isoperimetric inequalities for Rademacher sums, while our argument uses Khinchine's inequality for such sums. Kashin and Temlyakov (1994) develop results about what we would call best-ortho-bases for spaces with bounded mixed derivatives; their approach also implies a lemma like 8.11; their proof is based on volume estimates.

9. Near-minimaxity of BOB. As a result of the harmonic analysis in Section 8 and the ideas in Donoho (1993) we know that $\text{BAB}(\delta_1, \delta_2)$ is a kind of best basis in which to apply ideal keep-or-kill estimates. This is the key stepping-stone to our main result.

In this section we show that the risk for ideal keep-or-kill in $\text{BAB}(\delta_1, \delta_2)$ is within constants of the minimax risk over each $\mathcal{F}_p^{\delta_1, \delta_2}(C)$. From the oracle inequality of Section 7.2, we know that empirical basis selection, as in (4.1)–(4.3), which empirically selects a basis and applies thresholding within it, will always be nearly as good as ideal keep-or-kill in $\text{BAB}(\delta_1, \delta_2)$ —even though it makes no assumptions on δ_1 or δ_2 . This means that empirical best-basis denoising obeys a near-minimaxity result like Theorem 1.1.

THEOREM 9.1. *Best-basis denoising, defined in Section 4, with λ defined as in Section 7.2, comes within logarithmic factors of minimax over each functional class $\mathcal{F}_p^{\delta_1, \delta_2}(C)$, $0 < \delta_1, \delta_2 \leq 1$, $C > 0$, $1/p < \rho + 1/2$. If $\hat{f}^{*,\lambda}$ denotes the best-basis denoising estimator*

$$(9.1) \quad \sup_{\mathcal{F}} \text{MSE}(\hat{f}^{*,\lambda}, \bar{f}) \leq \text{Const}(\delta_1, \delta_2, p) \log(n) M^*(\sigma, n; \mathcal{F}) \quad \text{as } n \rightarrow \infty,$$

for each $\mathcal{F} \in \mathcal{AS}$.

The key arguments to prove Theorem 9.1 are given in Sections 9.1 and 9.4. Our main result—Theorem 1.1—will be proved in Section 10 by using some of those results a second time.

9.1. Lower bound on the minimax risk. We first study the minimax risk and show that it obeys the lower bound

$$(9.2) \quad M^*(\sigma, n; \mathcal{F}_p^{\delta_1, \delta_2}(C)) \geq K(\delta_1, \delta_2)(C^2)^{1-r} (\varepsilon^2)^r \quad \text{as } \varepsilon = \sigma/\sqrt{N} \rightarrow 0,$$

where $r = 2\rho/(2\rho + 1)$.

We use the method of cubical subproblems. In a modified definition in this section, by *orthogonal hypercube* \mathcal{H} of dimension m and side δ , we mean a collection of all sums $g_0 + \sum_{k=1}^m \alpha_k g_k$ where the $g_k = g_k(i_1, i_2)$ are $n \times n$ arrays, orthonormal with respect to the specially normalized l_N^2 norm

$$\frac{1}{\sqrt{N}} \sum_{i_1, i_2} g_k(i_1, i_2) g_{k'}(i_1, i_2) = \mathbf{1}_{\{k=k'\}}$$

and all the $|\alpha_i| \leq \delta$. The following lemma may be proved as in Donoho and Johnstone (1994a).

LEMMA 9.2. *Let $\varepsilon = \sigma/\sqrt{N}$. Suppose a class \mathcal{F} contains an orthogonal hypercube of sidelength $\varepsilon \leq \delta < (11/10)\varepsilon$ and dimension $m(\varepsilon)$. Then, for an absolute constant $A > 1/10$,*

$$(9.3) \quad M^*(\sigma, n; \mathcal{F}) \geq Am(\varepsilon)\varepsilon^2.$$

To make effective use of this, we seek cubes of sufficiently high dimension and prescribed sidelength. The following lemma is proved in the Appendix.

LEMMA 9.3. *Let $\varepsilon = \sigma/\sqrt{N}$. Each class $\mathcal{F}_p^{\delta_1, \delta_2}(C)$ contains orthogonal hypercubes (orthogonal with respect to l_N^2 norm) of sidelength $\delta = \varepsilon(1 + o(1))$ and dimension $m(\varepsilon, C)$ where*

$$(9.4) \quad m(\varepsilon, C) \geq K(\delta_1, \delta_2)(C/\varepsilon)^{2/(2\rho+1)}, \quad 0 < \varepsilon < \varepsilon_0,$$

and

$$(9.5) \quad \rho = \frac{\delta_1 \delta_2}{\delta_1 + \delta_2}.$$

Combining these two lemmas gives the lower bound (9.2).

9.2. *Equivalent estimation problems.* Sections 2–7 of this paper work in a setting of $n \times n$ arrays. Section 8 works in a setting of functions on the continuum unit square. Theorem 9.1 is based on a combination of both points of view.

From the viewpoint of Sections 2–7, one would naturally consider applying CART and BOB estimators to data $y_{\mathbf{i}}, \mathbf{i} \in [0, n]^2$. Suppose instead that we define the rescaled data

$$\tilde{y}_{\mathbf{i}} = N^{-1/2} y_{\mathbf{i}}, \quad \mathbf{i} \in [0, n]^2$$

and also define $\varepsilon = \sigma/n = \sigma/\sqrt{N}$. The results we get in applying (appropriately calibrated) CART or BOB to such data are (obviously) proportional to the results we get in applying the same techniques to the unscaled data.

There is a connection between these rescaled data and data about the function f on the continuum square. Let R denote both a dyadic rectangle of $[0, n]^2$ and the same rectangle on the continuum square $[0, 1]^2$. Recall that $\varphi_R^1(x, y)$ denotes a function on the continuum square $[0, 1]^2$ normalized to $L^2[0, 1]^2$ -norm 1, and $\phi_R^1(i_1, i_2) = h_{I_1^1}(i_1)\chi_{I_2^1}(i_2)$ is the same function, only on the grid $[0, n]^2$ and normalized to $l^2(N)$ norm 1. (Here “same” means that we identify the discrete interval $\{0, 1, 2, \dots, n-1\}$ as being the “same” as the continuous interval $[0, 1)$, and $\{0, 1, 2, \dots, n/2-1\}$ as being the “same” as the continuous interval $[0, 1/2)$, etc.) Then

$$\langle \tilde{y}_{\mathbf{i}}, \phi_R^1 \rangle_{l^2(N)} = \langle f, \varphi_R^1 \rangle_{L^2[0, 1]^2} + \varepsilon z_R^1,$$

where the z_R^1 are $N(0, 1)$, and independent in rectangles which are disjoint. Similar relationships hold between ϕ_R^2 and φ_R^2 .

Hence the discrete-basis analysis of rescaled data $\tilde{y}_{\mathbf{i}}$ has the interpretation of giving noisy measurements about the continuum coefficients of f and vice versa. Moreover, suppose that \mathcal{P}^* is a complete dyadic recursive partition of the discrete grid $[0, n]^2$ and we consider only the coefficients attached to rectangles in the nonterminal nodes of this partition. The partial reconstruction of f from just those coefficients is simply the collection of f 's pixel level

averages; formally, if we put

$$f(i_1, i_2) = \sum_{R \in NT(\mathcal{P}^*)} \alpha_R \phi_R^{s(R)}$$

and

$$\tilde{f}(x, y) = \sum_{R \in NT(\mathcal{P}^*)} \alpha_R \varphi_R^{s(R)}(x, y),$$

then $\tilde{f}(x, y)$ takes the value $f(i_1, i_2)$ throughout the rectangle $[i_1/n, (i_1 + 1)/n) \times [i_2/n, (i_2 + 1)/n)$.

Consider the problem of estimating $(\alpha_R^s(R))_{R \in NT(\mathcal{P}^*)}$ from noisy data $\langle f, \varphi_R^{s(R)} \rangle_{L^2[0,1]^2} + \varepsilon z_R^{s(R)}$. By Parseval, the squared l^2 risk

$$(9.6) \quad \varepsilon^2 + E \sum_{R \in NT(\mathcal{P}^*)} (\hat{\alpha}_R^{s(R)} - \alpha_R^{s(R)})^2 = N^{-1} E \|\hat{f} - \tilde{f}\|_{l^2(N)}^2,$$

and so the mean-squared error in the coefficient domain gives us the mean-squared error for recovery of pixel-level averages in the other domain.

9.3. Discrete and continuous partitionings. Consider now $BAB(\delta_1, \delta_2)$ for a given δ_1, δ_2 pair. This corresponds to an infinite sequence of families $\mathcal{R}(j)$, each family partitioning the continuum square $[0, 1]^2$ by congruent rectangles of area 2^{-j} .

Such a sequence of partitions of $[0, 1]^2$ usually cannot be interpreted as providing also a sequence of partitions for the discrete square $0 \leq i_1, i_2 < n$. A sequence of partitions for the discrete square also has the particular constraint that out of the first $\log_2(N)$ splits, exactly half will be vertical and half horizontal. Put another way, if we consider some BAB , those rectangles which are not too narrow in any direction, that is, where each sidelength exceeds $1/n$, also correspond to rectangles in a complete dyadic recursive partition of the discrete square $[0, n]^2$. But there exist BAB [for example those with $\min(\delta_1, \delta_2)$ close to zero and $\max(\delta_1, \delta_2)$ close to one] which, at some level j between $\log_2(n)$ and $\log_2(N)$, have already split in a certain direction more than $\log_2(n)$ times. Consequently, the continuum BAB is not quite available in the analysis of finite data sets.

On the other hand, in the analysis of finite data sets, there are available bases which achieve the same estimates of coefficient decay as in the continuum case.

DEFINITION 9.4. For a given pair of exponents (δ_1, δ_2) , and whole number J , we call a *balanced finite optimal sequential partition*, an application of the optimal sequential partitioning rule of Definition 8.5, with two extra rules:

- (iii) The process stops at stage $2J$. There are at most $2J$ “balls”;
- (iv) The process must preserve $j_i^*(j) \leq J$. Once a certain “box” has “ J balls,” all remaining allocations of “balls” are to the “other box.”

LEMMA 9.5. *If $0 < \delta_1, \delta_2 \leq 1$, let $j_i^*(j)$ denote the result of a balanced finite optimal sequential partitioning. Let $\mathcal{R}^*(j)$ denote the associated collection of rectangles. With $\rho = (\delta_1 \delta_2) / (\delta_1 + \delta_2)$ and $1/p < \rho + 1/2$,*

$$(9.7) \quad \left(\sum_{R \in \mathcal{R}^*(j)} |\alpha_R|^p \right)^{1/p} \leq 2C 2^{-j(\rho + 1/2 - 1/p)}.$$

The proof is simply to inspect the proof of Corollary 8.7 and notice that the constraint preventing allocation of “balls” to certain “boxes” means that in certain expressions one can replace terms like

$$\max(2^{-j_1^*(j)\delta_1}, 2^{-j_2^*(j)\delta_2})$$

by the even smaller

$$\min(2^{-j_1^*(j)\delta_1}, 2^{-j_2^*(j)\delta_2}).$$

9.4. *Upper bound on ideal risk.* We now study the ideal risk and show that it obeys an upper bound similar in form to the lower bound of Section 9.1. Starting now, let $BAB^*(\delta_1, \delta_2)$ denote the modified basis described in the previous subsection.

LEMMA 9.6. *Let $\mathcal{R}(\text{KEEP-KILL}, f; \varepsilon)$ be the ideal risk for keep-kill estimation in $BAB^*(\delta_1, \delta_2)$. Then with $r = 2\rho / (2\rho + 1)$,*

$$(9.8) \quad \sup_{f \in \mathcal{F}_p^{\delta_1, \delta_2}(C)} \mathcal{R}(\text{KEEP-KILL}, f; \varepsilon) \leq B(\delta_1, \delta_2, p)(C^2)^{1-r} (\varepsilon^2)^r, \quad 0 < \varepsilon < \varepsilon_0.$$

PROOF. As in Donoho, Johnstone, Kerkyacharian and Picard (1995), consider the optimization problem

$$m_j(\varepsilon; \gamma) = \max \|\theta\|_{l^2}^2 \text{ subject to } \|\theta\|_{l^r} \leq \varepsilon, \quad \|\theta\|_{l^p} \leq \gamma, \quad \theta \in R^{2^j}.$$

By Parseval (9.6) the best possible risk for a purely keep-kill estimate is $\varepsilon^2 + \sum \min(\alpha_R^2, \varepsilon^2)$. Also, by Lemma 9.5, there are constants $\gamma_j = \gamma_j(C)$ so that for $f \in \mathcal{F}_p^{\delta_1, \delta_2}(C)$,

$$\left(\sum_{\mathcal{R}^*(j)} |\alpha_R|^p \right)^{1/p} \leq \gamma_j.$$

The largest risk of ideal keep-kill is thus

$$\begin{aligned} & \max_{f \in \mathcal{F}_p^{\delta_1, \delta_2}(C)} \sum_j \sum_{R \in \mathcal{R}^*(j)} \min(\alpha_R^2, \varepsilon^2) \\ & \leq \max_j \sum_{\mathcal{R}^*(j)} \min(\alpha_R^2, \varepsilon^2) \text{ subject to } \left(\sum_{\mathcal{R}^*(j)} |\alpha_R|^p \right)^{1/p} \leq \gamma_j \\ & = \sum_j m_j(\varepsilon, \gamma_j). \end{aligned}$$

Now Donoho, Johnstone, Kerkycharian and Picard (1995) give the explicit expression

$$(9.9) \quad m_j(\varepsilon; \gamma) = \min(2^j \varepsilon^2, \gamma^p \varepsilon^{2-p}, \gamma^2)$$

and applying this, we have

$$(9.10) \quad \sum_j m_j(\varepsilon, \gamma_j) \leq (C^2)^{1-r} (\varepsilon^2)^r K(\delta_1, \delta_2, p). \quad \square$$

This is the risk of an ideal denoising by a keep-or-kill estimator not obeying hereditary constraints.

9.5. *Near minimaxity of best-basis denoising.* We have so far shown that the ideal risk is within constant factors of the minimax risk. Invoking now the oracle inequality of Theorem 7.1, the worst-case risk of the BOB estimator \hat{f} does not exceed the ideal risk—and hence the minimax risk—by more than a logarithmic factor. This completes the proof of Theorem 9.1.

10. Near minimaxity of CART. We now are in a position to complete the proof of Theorem 1.1. We do this by showing that ideal dyadic CART is essentially as good as ideal best-basis denoising.

LEMMA 10.1. *Let $\mathcal{R}(\text{KEEP-KILL}, f; \varepsilon)$ be the ideal risk for keep-kill estimation in $\text{BAB}^*(\delta_1, \delta_2)$. Let $\mathcal{R}(\text{HERED}, f; \varepsilon)$ be the ideal risk for hereditary estimation in $\text{BAB}^*(\delta_1, \delta_2)$. If $1/p < \rho + 1/2$,*

$$(10.1) \quad \begin{aligned} & \sup_{f \in \mathcal{F}_p^{\delta_1, \delta_2}(C)} \mathcal{R}(\text{HERED}, f; \varepsilon) \\ & \leq B(\delta_1, \delta_2, p) \sup_{f \in \mathcal{F}_p^{\delta_1, \delta_2}(C)} \mathcal{R}(\text{KEEP-KILL}, f; \varepsilon). \end{aligned}$$

Once this lemma is established, it follows from Sections 9.1 and 9.4 that the risk of ideal dyadic CART is within constant factors of the minimax risk. Now the oracle inequality for dyadic CART, Theorem 7.2, shows that the performance of empirical dyadic CART comes within logarithmic factors of the ideal risk for dyadic CART. Theorem 1.1 therefore follows as soon as Lemma 10.1 is established.

To prove the lemma, note that the ideal keep-or-kill estimator for a function f has nonzero coefficients at sites

$$(10.2) \quad \mathcal{S}(f) = \{(j, R) : |\alpha_{j,R}(f)| \geq \varepsilon\}.$$

This can be modified to a hereditary linear estimator by expanding \mathcal{S} slightly.

DEFINITION 10.2. Let \mathcal{S} be a collection of dyadic rectangles, for example (10.2). Then the *hereditary cover* of \mathcal{S} , denoted \mathcal{S}^* , is the collection of all such rectangles and their ancestors in the partitioning associated with the basis $\text{BAB}^*(\delta_1, \delta_2)$.

The keep-kill estimator $\hat{\alpha}[\mathcal{S}^*]$ with nonzero coefficients at sites \mathcal{S}^* is a hereditary linear estimator. The ideal risk obeys

$$\begin{aligned} E\|\hat{\alpha}[\mathcal{S}^*] - \alpha\|_2^2 &= \sum_{(j,R) \notin \mathcal{S}^*} \alpha_{j,R}^2 + \varepsilon^2(\#\mathcal{S}^* + 1) \\ &\leq \sum_{(j,R) \notin \mathcal{S}} \alpha_{j,R}^2 + \varepsilon^2(\#\mathcal{S} + 1) \quad (\text{as } \mathcal{S} \subset \mathcal{S}^*). \end{aligned}$$

Suppose we could bound $\#\mathcal{S}^* \leq A(\#\mathcal{S})$ for some constant $A \geq 1$. Then we would have

$$\begin{aligned} E\|\hat{\alpha}[\mathcal{S}^*] - \alpha\|_2^2 &\leq \sum_{(j,R) \notin \mathcal{S}} \alpha_{j,R}^2 + \varepsilon^2(A\#\mathcal{S} + 1) \quad (\text{as } \#\mathcal{S}^* \leq A\#\mathcal{S}) \\ &\leq A \left(\sum_{(j,R) \notin \mathcal{S}} \alpha_{j,R}^2 + \varepsilon^2(\#\mathcal{S} + 1) \right) \\ &= AE\|\hat{\alpha}[\mathcal{S}] - \alpha\|_2^2. \end{aligned}$$

It would then follow that risk bounds derived for keep-or-kill estimators would give rise to proportional risk bounds for hereditary linear estimators derived from their hereditary covers.

While the relation $\#\mathcal{S}^* \leq A(\#\mathcal{S})$ does not hold for every f , a weaker inequality of the same form holds, where one compares the largest possible size of $\#\mathcal{S}(f)$ for an $f \in \mathcal{F}_p^{\delta_1, \delta_2}(C)$ with the largest possible size of $\#\mathcal{S}^*$. Lemma 10.3 establishes this inequality; retracing the logic of the last few displays shows that it immediately implies Lemma 10.1, with $B = A$.

LEMMA 10.3. *Define*

$$(10.3) \quad N(\delta_1, \delta_2, p, C) = \sup\{\#\mathcal{S}(f) : f \in \mathcal{F}_p^{\delta_1, \delta_2}(C)\}$$

the largest number of coefficients used by an ideal keep-kill estimator in treating functions from $\mathcal{F}_p^{\delta_1, \delta_2}(C)(\delta_1, \delta_2)$. Similarly, let

$$(10.4) \quad N^*(\delta_1, \delta_2, p, C) = \sup\{\#\mathcal{S}^*(f) : f \in \mathcal{F}_p^{\delta_1, \delta_2}(C)\}$$

be the size of the largest corresponding hereditary cover. If $1/p < \rho + 1/2$, then for a finite positive constant $A = A(\delta_1, \delta_2, p)$,

$$N^* \leq A(\delta_1, \delta_2, p)N.$$

PROOF. If $\theta = (\theta_i)_{i=1}^d$ is a vector of dimension d satisfying $\|\theta\|_{l^p} \leq \gamma$, then $\#\{i : |\theta_i| \geq \varepsilon\} \varepsilon^p \leq \gamma^p$ so

$$(10.5) \quad \#\{i : |\theta_i| \geq \varepsilon\} \leq (\gamma/\varepsilon)^p$$

and of course

$$(10.6) \quad \#\{i : |\theta_i| \geq \varepsilon\} \leq d.$$

Consider now the application of this to the vector $\theta_i = (\alpha_{j,R})_i$ which has $d = 2^j$, with $\gamma = \gamma_j(\mathcal{F}_p^{\delta_1, \delta_2}(C))$. Then $\#\{i : |\theta_i| \geq \varepsilon\} \leq \min(2^j, (\gamma_j/\varepsilon)^p)$. The

first term 2^j is sharp for $0 \leq j \leq j_0$, where $j_0 = j_0(\varepsilon, C; \delta_1, \delta_2)$ is the real root of $2^j = (C2^{-j(\rho+1/2-1/p)}/\varepsilon)^p$. By a calculation, $j_0 = \log_2(C/\varepsilon)/(\rho + 1/2)$.

For notational convenience, stratify the set \mathcal{S} as

$$\mathcal{S} = \mathcal{S}_0 \cup \mathcal{S}_1 \cup \dots \mathcal{S}_j \cup \dots,$$

where $\mathcal{S}_j = \{(j', R) \in \mathcal{S}, j' = j\}$ and

$$(10.7) \quad \#\mathcal{S}_j = 2^j, \quad 0 \leq j \leq j_0.$$

Also we have

$$(10.8) \quad \begin{aligned} \#\mathcal{S}_j &\leq 2^{j_0} 2^{-\beta(j-j_0)}, \quad j \geq j_0, \\ \beta &= \beta(\delta_1, \delta_2, p) = p(\rho + 1/2 - 1/p) > 0. \end{aligned}$$

Now consider the cover \mathcal{S}^{**} defined by:

$$\begin{aligned} \mathcal{S}^{**} &= \{(j, R), 0 \leq j \leq j_0, \\ &\quad (j, R), j > j_0 \text{ and } (j, R) \text{ has a descendant in } \mathcal{S}\}. \end{aligned}$$

By construction, \mathcal{S}^{**} contains the hereditary cover (it contains terms at $j < j_0$ which the hereditary cover might not), and so bounds on the size of \mathcal{S}^{**} apply to \mathcal{S}^* also. Now

$$(10.9) \quad \#\mathcal{S}^{**} \leq 2^{j_0+1} + \sum_{j>j_0} A(j, R)\#\mathcal{S}_j,$$

where $A(j, R)$ is the number of ancestors (j', R') of a term (j, R) at level $j_0 < j' \leq j \leq J$. As $A(j, R) \leq (j - j_0)$, (10.8) gives

$$\begin{aligned} \#\mathcal{S}^{**} &\leq 2^{j_0+1} + \sum_{j>j_0} (j - j_0)2^{j_0}2^{-\beta(j-j_0)} \\ &= 2^{j_0} \left(2 + \sum_{j>j_0} (j - j_0)2^{-\beta(j-j_0)} \right). \end{aligned}$$

We conclude that

$$N^*(\delta_1, \delta_2, p) \leq 2^{j_0} B_1,$$

for some constant $B_1(\delta_1, \delta_2, p)$. On the other hand, by constructing a hypercube at level $\lfloor j_0 \rfloor$ using the approach of Lemmas 8.12 and 9.3, we obtain, for a constant $B_2(\delta_1, \delta_2, p)$,

$$N(\delta_1, \delta_2, p) \geq B_2 2^{j_0}.$$

Hence we may take $A = B_1/B_2$. \square

11. Discussion. We collect here some final remarks.

11.1. *Clarifications.* We would like to point out clearly that the way that the term CART is generally construed—as greedy growing of an exhaustive partition followed by optimal pruning in the implicit basis—is not what we have studied in this paper. Also, the data structure we have assumed—regu-

lar equispaced data on a two-dimensional rectangular lattice—is unlike the irregularly scattered data often assumed in CART studies. It would be interesting to know what properties can be established for the typical greedy growing nondyadic CART algorithm in the irregularly scattered data case.

To minimize misunderstanding, let us be clear about the intersection between CART and BOB. CART is a general methodology used for classification and discrimination or for regression. It can be used on regular or irregularly spaced data and it can construct optimal or greedy partitions within the general framework. Best-ortho-basis is a general methodology for adaptation of orthogonal bases to specific problems in applied mathematics. It can be used in constructing adaptive time frequency bases, and also (as we have seen in this paper) in constructing adaptive bases for functions on Cartesian product domains. We have shown that the methods have something in common, but, strictly speaking, only intersect under a very specific choice of problem areas and entropy. Further discussion about patent lawsuits is unwarranted and pointless.

11.2. *Extensions.* Somewhat more general results are implicit in the results established here.

First, one can consider classes $\mathcal{F}_{p_1, p_2}^{\delta_1, \delta_2}(C_1, C_2)$, $0 < \delta_i < 1$ and $p_i, C_i > 0$, consisting of functions obeying

$$\|D_h^i f\|_{p_i} \leq C_i h^{\delta_i}, \quad h > 0, i = 1, 2.$$

The classes we have considered here in this paper are the special cases $C_1 = C_2 = C$ and $p_1 = p_2 = p$. Parallel results hold for these more general classes, and by essentially the same arguments, with a bit more bookkeeping. We avoided the study of these more general classes only to simplify exposition.

Second, the log terms we have established in Theorems 1.1 and 9.1 can be replaced by somewhat smaller log terms. More specifically, in cases where the minimax risk scales like N^{-r} , the method of proof given here actually shows that the worst-case risk of dyadic CART is within a factor $O(\log(n)^r)$ of minimax. As $0 < r < 1$, this is an improvement in the size of the log term.

Third, one can obtain results for higher-order smoothness classes $\mathcal{F}_{p_1, p_2}^{\alpha_1, \alpha_2}(C_1, C_2)$ with $\alpha_i = m_i + \delta_i$, with m_i whole numbers and δ_i fractions, consisting of functions f obeying

$$\|D_h^i f^{(m_i)}\|_{p_i} \leq C_i h^{\delta_i}, \quad h > 0, i = 1, 2.$$

Such classes can be addressed using recursive partitioning methods with piecewise polynomial fits. Instead of using a piecewise constant reconstruction, one uses a piecewise polynomial of some fixed degree D on each piece. The analysis of such procedures is entirely parallel to the analysis in this article; one simply replaces the library of all anisotropic Haar functions by the library of all anisotropic Alpert functions. Alpert functions are piecewise polynomials of degree D deriving from Legendre polynomials in the same way that the Haar wavelets derive from the indicator function $1_{[0,1]}$. Using

such procedures, everything in this paper generalizes straightforwardly. We avoided discussing the potentially more powerful procedures based on polynomial fitting in order to focus attention on the relation between BOB and CART. Useful background on Alpert bases can be had in Alpert, Beylkin, Coifman and Rokhlin (1993) and Donoho, Dyn, Levin and Yu (1996).

Fourth, DeVore (1994) has informed the author that heuristic methods of data compression and denoising based on the library of systematic tensor products of smooth wavelets can be used to get results which in practice are very effective.

11.3. *Linear estimators.* As mentioned in the introduction, dyadic CART can outperform kernel estimators and related linear procedures in a minimax sense. Let $M_L^*(\sigma, n; \mathcal{F})$ denote the *minimax linear risk*, defined as in (1.3), but where only linear estimators are allowed. For example, the class of linear estimates includes anisotropic kernel smoothing procedures, where the bandwidth is allowed to differ in each of the two directions. The following result shows that such linear estimates, even when tuned optimally for a class $\mathcal{F}_p^{\delta_1, \delta_2}(C)$ where $p < 2$, can be outperformed by the dyadic CART estimator at the level of rates.

THEOREM 11.1. *Let $1 < p < 2$ and $\rho' = \rho + 1/2 - 1/p > 0$. For M_L^* the minimax risk among all linear procedures, we have the lower bound*

$$M_L^*(\sigma, n; \mathcal{F}_p^{\delta_1, \delta_2}(C)) \geq \text{Const} \cdot N^{-(2\rho'/(2\rho'+1))}.$$

For comparison, of course, dyadic CART achieves

$$\sup_{\mathcal{F}_p^{\delta_1, \delta_2}(C)} \text{MSE}(\hat{f}_n^*, f) \leq \text{Const} \cdot (\log(N)/N)^{2\rho/(2\rho+1)}.$$

As $\rho > \rho'$ in the applicable range $p < 2$, dyadic CART achieves a faster rate of convergence over $\mathcal{F}_p^{\delta_1, \delta_2}(C)$ than any linear procedure. For example, set $p = 3/2$, $\delta_1 = 1$, $\delta_2 = 2/3$, so $\rho = 2/5$. Then $\rho' = 7/30$, so linear estimates converge no faster over this class than a rate $N^{-r'}$ where $r' = 2\rho'/(2\rho'+1) = 7/22 < 1/3$; while dyadic CART achieves a rate at most a logarithmic factor worse than $r = 2\rho/(2\rho+1) = 4/9$.

This phenomenon parallels results in other settings where linear estimates have been shown not to achieve minimax rates. Examples of such settings include, for squared L^2 -loss, classes of bounded variation, Sobolev spaces with $p < 2$, and Besov and Triebel classes with $p < 2$. These settings have been treated in work of Nemirovskii, Tsybakov and Polyak (1985), Nemirovskii (1986), and Donoho and Johnstone (1994d); see Donoho, Johnstone, Kerkyacharian and Picard (1995) for references.

The phenomenon derives from a geometric property of the classes $\mathcal{F}_p^{\delta_1, \delta_2}(C)$, the lack of 2-convexity, which means that the quadratic hull of $\mathcal{F}_p^{\delta_1, \delta_2}(C)$ is essentially larger than $\mathcal{F}_p^{\delta_1, \delta_2}(C)$ itself; compare Donoho, Liu, and MacGibbon (1990). In its simplest form, the phenomenon appears as follows [Donoho and

Johnstone (1994d)]. Suppose we have the problem of estimating a d -vector $(\xi_i; 1 \leq i \leq d)$ from observations $v_i = \xi_i + \varepsilon z_i$, where the z_i are i.i.d. $N(0, 1)$. The vector ξ is known to lie in a d -dimensional l^p -ball. Let $0 < p < 2$. *The minimax linear risk (min over linear procedures only, max over the unit l_d^p ball) for estimating ξ with squared l^2 -norm loss, is the same as the minimax linear risk over the standard Euclidean ball (min over linear procedures only, max over the unit l_d^2 ball).*

This observation lies at the heart of Theorem 11.1. By combining (8.9) and a construction similar to the one underlying Lemma 9.3, one can show that the class $\mathcal{F}_p^{\delta_1, \delta_2}(C)$ contains, for each $j \geq 0$, an l_d^p ball of radius $r_j = \text{Const} \cdot 2^{-j(\rho + 1/2 - 1/p)}$ and dimension $d_j = 2^j$. Each such ball furnishes a finite-dimensional parametric subfamily of the original functional class, in which the function to be estimated is isometrically identified with a parameter vector in a d_j -dimensional space obeying a constraint on the l_d^p norm. Estimation of a function in such a restricted subfamily, from data (1.1), can be reduced, by a sufficiency argument, to estimating the parameter vector ξ using observations $v = \xi + \varepsilon z$ of dimension d_j , where the noise is a white Gaussian noise and $\varepsilon = \sigma/\sqrt{N}$. Call this problem a finite-dimensional subproblem; its minimax linear risk is a lower bound on the minimax linear risk of the full problem.

Applying the italicized observation about minimax linear risk in such a finite-dimensional subproblem, the l_d^p constraint can be relaxed to an l_d^2 constraint without affecting the minimax linear risk in the subproblem. This relaxation implies a geometrically larger subproblem. Of course, the minimax linear risk of the enlarged subproblem is not smaller than the minimax risk (all measurable procedures allowed). We may use Lemma 9.2 to obtain a lower bound of the form $\text{Const} \cdot d_j \varepsilon^2$ for the minimax risk of such enlarged subproblems, for each j in a certain range.

Selecting, at each sample size N , the most difficult such subproblem in this range of j , we get a lower bound of the form indicated in the theorem.

In short, the proof has two ideas. First, the minimax linear risk over $\mathcal{F}_p^{\delta_1, \delta_2}(C)$ behaves as if this class contained $l_{d_j}^2$ balls of radius $r'_j = \text{Const} \cdot 2^{-j\rho'}$, and dimension $d_j = 2^j$, $\rho' = \rho + 1/2 - 1/p$. Second, such a class could not have a minimax risk better than $\text{Const} \cdot (\varepsilon^2)^{2\rho'/2\rho'+1}$.

11.4. Important related work. We also mention some related work that may be of interest to the reader.

Complexity bounds and oracle inequalities. Of course there is a heavy reliance of this paper on Donoho and Johnstone (1994a,b). But let us also clearly point out that the general idea of oracle inequalities is clearly present in Foster and George (1994), who used a slightly different oracle less suited for our purposes here. Our underlying proof technique—the complexity bound underlying the proofs of Theorems 7.1 and 7.2—is very closely related to the minimum complexity formalism of Barron and Cover (1991), and subsequent work by Birgé and Massart (1997).

Density estimation. This paper grew out of a discussion with Engel, who wondered how to generalize the results of Engel (1994) to higher dimensions. Engel (personal communication) has reported progress on obtaining results on the behavior of a procedure like dyadic CART in the setting of density estimation.

Anisotropic smoothness spaces. Neumann and von Sachs (1995) have also recently studied anisotropic smoothness classes and have shown that wavelet thresholding in a tensor wavelet basis is nearly minimax for higher-order anisotropic smoothness classes. This shows that nonadaptive basis methods could also be used for obtaining nearly minimax results; the full adaptivity of CART is not really necessary for minimaxity alone.

Time-frequency analysis. Important related ideas are contained in two recent manuscripts associated with Coifman's group at Yale. The article of Thiele and Villemoes (1996) independently uses fast dyadic recursive partitioning of the kind discussed here, only in a setting where the two dimensions are time and frequency. The thesis of Bennett (1997) independently uses fast dyadic recursive partitioning of the kind discussed here, only in a setting where the basis functions are anisotropic Walsh functions rather than anisotropic Haar functions.

APPENDIX

A.1. Proof of Theorems 7.1 and 7.2. We prove a more general fact, concerning estimation in overcomplete dictionaries. The proof we give is a light modification of arguments in Donoho and Johnstone (1994a,b).

A1.1. Constrained Minimum Complexity Estimates. Suppose we have an $N \times 1$ vector \mathbf{y} and a dictionary of $N \times 1$ vectors φ_γ . We wish to approximate \mathbf{y} as a superposition of dictionary elements $\mathbf{y} \approx \sum_{i=1}^m \beta_\gamma \varphi_\gamma$.

We construct a matrix Φ which is N by p , where p is the total number of dictionary elements. Let each column of the Φ matrix represent one dictionary element. Note that in the case of most interest to us, $p \gg N$, as Φ contains more than just a single basis. For example, in the setting of this paper, \mathcal{D} is the dictionary of all anisotropic Haar functions, which has approximately $p = 4N$ elements.

For approximating the vector y , we consider the vector $\tilde{\beta} \in R^p$, the vector $\tilde{f} = \Phi \tilde{\beta}$ denotes a corresponding linear combination of dictionary elements. This places the approximation \tilde{f} in correspondence with the coefficient vector $\tilde{\beta}$. Owing to the possible overcompleteness of Φ , this correspondence is in general one-to-many.

Define now the *empirical complexity* functional

$$K(\tilde{f}, y) = \|\tilde{f} - y\|_2^2 + \lambda^2 \sigma^2 N(\tilde{f}),$$

where

$$N(\tilde{f}) = \min_{\tilde{f} = \Phi\tilde{\beta}} \#\{j: \tilde{\beta}_j \neq 0\}$$

is the complexity of constructing \tilde{f} from the dictionary Φ . Also, define the *theoretical complexity* functional

$$K(\tilde{f}, f) = \|\tilde{f} - f\|_2^2 + \lambda^2 \sigma^2 N(\tilde{f}).$$

Let \mathcal{E} be a collection of “allowable” coefficient vectors $\beta \in R^p$. We will be interested in approximations to y obeying these constraints and having small complexity. In a general setting, one can think of many interesting constraints to impose on allowable coefficients; for example, that coefficients should be positive, that coefficients should generate a monotone function, that nonzero coefficients are attached to pairwise orthogonal elements.

Define the \mathcal{E} -constrained minimum empirical complexity estimate

$$\hat{f}^* = \operatorname{argmin}_{\{\tilde{f} = \Phi\tilde{\beta}: \tilde{\beta} \in \mathcal{E}\}} K(\tilde{f}, y).$$

In a moment we will prove the following.

COMPLEXITY BOUND. Suppose $y = f + z$, where z is i.i.d. $N(0, 1)$. Fix $\mathcal{E} \subset \mathbf{R}^p$, fix $\zeta > 8$ and consider the \mathcal{E} -constrained minimum complexity model selection with $\lambda = \zeta(1 + \sqrt{2 \log p})$.

$$(A.1) \quad EK(\hat{f}^*, f) \leq A(\zeta) \left(\lambda^2 \sigma^2 + \min_{\{\tilde{f} = \Phi\tilde{\beta}: \tilde{\beta} \in \mathcal{E}\}} K(\tilde{f}, f) \right).$$

This shows that the empirical minimum complexity estimate is not far off from minimizing the theoretical complexity.

(In the above bound, the limitation $\zeta > 8$ is not intrinsic to the problem; Johnstone has informed the author that by refinements of the arguments below one can obtain results of the same general form for smaller values of λ corresponding roughly to any $\zeta > 1$.)

A.1.2. Relation to CART and BOB. We now explain why the complexity bound implies Theorems 7.1 and 7.2.

We begin with the observation that the empirical complexity $K(f, y)$ is just what we earlier called a complexity penalized sum of squares.

Assume now that the dictionary \mathcal{D} is the collection of all anisotropic Haar functions. Two constraint sets are particularly interesting.

First, let \mathcal{E}_{BOB} be the collection of all coefficient vectors β which arise from combinations of atoms that all belong together in some orthobasis built from the anisotropic Haar dictionary. Remember, the dictionary has $p \gg N$ atoms. So at most N elements of β can be nonzero at once under this constraint. Also, we have seen in Section 3 that each basis in the anisotropic Haar system corresponds to a certain decorated tree so this constraint says

that collections of coefficients which are allowed to be nonzero simultaneously correspond to certain collections of indices. This constraint can be made quite explicit and algorithmic, although we do not go into details here.

If we optimize the empirical complexity $K(\tilde{f}, y)$ over all \tilde{f} arising from a $\beta \in \mathcal{E}_{\text{BOB}}$ we get exactly the estimator (4.1)–(4.3). We encourage the reader to check this fact.

Second, there is the CART constraint. Let $\mathcal{E}_{\text{CART}}$ be the collection of all vectors β for which the nonzero coefficients in the corresponding β only refer to atoms which can appear together in an orthogonal basis, and for which the nonzero coefficients only occur in an hereditary pattern in that basis. We remark that $\mathcal{E}_{\text{CART}} \subset \mathcal{E}_{\text{BOB}}$.

If we optimize the empirical complexity $K(\tilde{f}, y)$ over all \tilde{f} arising from a $\beta \in \mathcal{E}_{\text{CART}}$ we get exactly the estimator (2.4)–(2.5). We again encourage the reader to check this fact.

We now make two simple observations about the minimum complexity formalism, valid for any \mathcal{E} , which the reader should verify.

K1. The theoretical complexity of \hat{f}^* upperbounds the predictive loss

$$K(\hat{f}^*, f) \geq \|\hat{f}^* - f\|_2^2.$$

K2. The minimum theoretical complexity is within a logarithmic factor of the ideal risk

$$\begin{aligned} \min_{\tilde{f} \in \mathcal{E}} K(\tilde{f}, f) &= \min_{\tilde{f} \in \mathcal{E}} \|\tilde{f} - f\|_2^2 + \lambda^2 \sigma^2 N(\tilde{f}) \\ &\leq \lambda^2 \min_{\tilde{f} \in \mathcal{E}} (\|\tilde{f} - f\|_2^2 + \sigma^2 N(\tilde{f})) \\ &= \lambda^2 \min_{\beta \in \mathcal{E}} (\|\Phi\tilde{\beta} - \Phi\beta\|_2^2 + \sigma^2 \#\{j: \tilde{\beta}_j \neq 0\}) \\ &= \lambda^2 \mathcal{R}(\text{Ideal } \mathcal{E}, f). \end{aligned}$$

These observations, translated into the cases \mathcal{E}_{BOB} and $\mathcal{E}_{\text{CART}}$, give Theorems 7.1 and 7.2, respectively.

A.1.3. *Proof of the complexity bound.* In what follows we assume the noise level $\sigma^2 = 1$. We follow, line-by-line, Donoho and Johnstone (1994a, b) who analyzed the unconstrained case $\mathcal{E} = \mathbf{R}^p$. Exactly the same analysis applies in the constrained case.

We first let f^0 denote a model of minimum theoretical complexity:

$$K(f^0, f) = \min_{\tilde{f} \in \mathcal{E}} K(\tilde{f}, f).$$

As \hat{f}^* has minimum empirical complexity,

$$K(\hat{f}^*, y) \leq K(f^0, y).$$

As $\|\hat{f}^* - y\|_2^2 = \|\hat{f}^* - f - z\|_2^2$ we can relate empirical and theoretical complexities by

$$K(\hat{f}^*, y) = K(\hat{f}^*, f) + 2\langle z, f - \hat{f}^* \rangle + \|z\|_2^2,$$

and so, combining the last two displays,

$$K(\hat{f}^*, f) \leq K(f^0, f) + 2\langle z, \hat{f}^* - f^0 \rangle.$$

Now define the random variable

$$\mathscr{W}(k) = \sup\{\langle z, m^2 - m^1 \rangle : \|m^j - f\|_2^2 \leq k, \lambda^2 N(m^j) \leq k\}.$$

Then

$$K(\hat{f}^*, f) \leq K(f^0, f) + 2\mathscr{W}(K(\hat{f}^*, f)).$$

This display shows the key idea. It turns out that $\mathscr{W}(k) \ll k$ for all large k , and so this display forces $K(\hat{f}^*, f)$ to be not much larger than $K(f^0, f)$.

Denote the minimum theoretical complexity by $K^0 = K(f^0, f)$. Define $k_j = 2^j(1 - 8/\zeta)^{-1} \max(K^0, \lambda^2)$ for $j \geq 0$. Define the event

$$B_j = \{\mathscr{W}(k) \leq 4k/\zeta \text{ for all } k \geq k_j\}.$$

On the event B_j , the inequality

$$k \leq K^0 + 2\mathscr{W}(k)$$

has no solutions for $k \geq k_j$. Hence, on event B_j ,

$$K(\hat{f}^*, f) \leq k_j.$$

It follows that

$$\begin{aligned} EK(\hat{f}^*, f) &\leq \sum_{j=0}^{\infty} k_{j+1} \text{Prob}\{K(\hat{f}^*, f) \in [k_j, k_{j+1})\} \\ &\leq \sum_{j=0}^{\infty} k_{j+1} \text{Prob}\{K(\hat{f}^*, f) \geq k_j\} \\ &\leq \sum_{j=0}^{\infty} k_{j+1} \text{Prob}\{B_j^c\}. \end{aligned}$$

By Lemma A.1 we get

$$\begin{aligned} EK(\hat{f}^*, f) &\leq k_0 \sum_{j \geq 0} 2^{j+1}/(2^j)! \\ &\leq \max(K^0, \lambda^2)(1 - 8/\zeta)^{-1}6. \end{aligned}$$

Hence, the complexity bound (A.1) holds, with $A(\zeta) = (1 - 8/\zeta)^{-1}6$.

LEMMA A.1 [Donoho and Johnstone (1994a, b)].

$$\text{Prob}\{B_j^c\} \leq 1/(2^j)!$$

The proof depends on tail bounds for chi-squared variables, which, ultimately, depend on concentration-of-measure estimates (e.g., Borell–Tsirel’son inequality).

A.2. Proof of Lemma 8.4. The proof of each display is similar, so we just discuss the first. Fix a rectangle $R = I_1 \times I_2$ with $|I_1| = 2^{-j_1}$. Let $R^{1,0}$ and $R^{1,1}$ denote the left and right halves:

$$\langle f, \phi_R^1 \rangle = 2^{j/2} \left\{ \int_{R^{1,1}} f dx dy - \int_{R^{1,0}} f dx dy \right\}.$$

Hence for the very special increment $h = 2^{-j_1}$,

$$\begin{aligned} \int_{R^{1,0}} (D_h f)(x, y) dx dy &= \int_{R^{1,0}} (f(x + h, y) - f(x, y)) dx dy \\ &= \int_{R^{1,1}} f dx dy - \int_{R^{1,0}} f dx dy = 2^{-j/2} \langle f, \phi_R^1 \rangle. \end{aligned}$$

For any sum \sum_R over rectangles R with disjoint interiors,

$$\sum_R |\langle f, \phi_R^1 \rangle|^p = 2^{jp/2} \sum_R \left| \int_{R^{1,0}} D_h^1 f \right|^p.$$

Now by assumption $1/p < \rho + 1/2$, which (as $\rho \leq 1/2$) means $p > 1$. Let $1/p + 1/p' = 1$,

$$\left| \int_{R^{1,0}} D_h^1 f \right| \leq \|D_h^1 f\|_{L^p(R^{1,0})} \|1\|_{L^{p'}(R^{1,0})},$$

so

$$\sum_R |\langle f, \phi_R^1 \rangle|^p \leq 2^{j(p/2+(1-p))} \sum_R \|D_h^1 f\|_{L^p(R^{1,0})}^p.$$

Now if \sum_R is interpreted to mean the sum over a partition of $[0, 1]^2$ by congruent rectangles, then

$$\sum_R \|D_h^1 f\|_{L^p(R)}^p = \|D_h^1 f\|_{L^p(Q_h^1)}^p,$$

and so from $\|D_h^1 f\|_{L^p(R^{1,0})} \leq \|D_h^1 f\|_{L^p(R)}$ we conclude that

$$\begin{aligned} \left(\sum_{R \in \mathcal{R}(j)} |\langle f, \phi_R^1 \rangle|^p \right)^{1/p} &\leq 2^{j(1/p-1/2)} \|D_h^1 f\|_{L^p(Q_h^1)} \\ &\leq 2^{j(1/p-1/2)} h^{\delta_1} C = 2^{-j_1 \delta_1} 2^{j(1/p-1/2)} C. \end{aligned}$$

A.3. Proof of Lemma 8.6. We assume for the proof below that δ_1 and δ_2 are mutually irrational. Very slight modifications allow us to handle the exceptional cases.

Think of the quarterplane consisting of (x, y) with $x, y \geq 0$ as a collection of square “unit cells,” with vertices on the integer lattice. Think of the set where $x\delta_1 = y\delta_2$ as a ray S in this quarterplane, originating at $(0, 0)$.

Let $p_j = (j_1(j), j_2(j))$ denote the sequence of pairs of values $j_1\delta_1 = j_2\delta_2$ where $j_1 + j_2 = j$. Let $p_j^* = (j_1^*(j), j_2^*(j))$ denote the sequence of pairs of values obtained from the optimal sequential partitioning of definition 8.5.

Our claim, to be established below: p_j and p_j^* always belong to the same unit cell.

It follows from this claim that $j_1^*(j) > j_1(j) - 1$ and $j_2^*(j) > j_2(j) - 1$; as a result

$$\max(2^{-j_1^*(j)\delta_1}, 2^{-j_2^*(j)\delta_2}) \leq 2 \max(2^{-j_1(j)\delta_1}, 2^{-j_2(j)\delta_2}),$$

and the lemma follows.

The claim is proved by induction. Indeed, at $j = 0$, $p_j = p_j^* = 0$. So the claim is true at $j = 0$.

For the inductive step, suppose the claim is true for steps $0 \leq j \leq J$; we prove it for $J + 1$. Let C_j denote the unit cell containing p_j , where, if several cells qualify, we select a cell having p_j on the skew diagonal.

Under this convention, at each step j , p_j lies on the skew diagonal through this cell, which joins its upper left corner to its lower right corner. Supposing the claim is true at step j , p_j^* is either at the upper left corner or at the lower right corner of the cell. Note also that C_{j+1} is either above C_j , or to the right of C_j .

With this set-up, the inductive step requires two things: (1) that if p_j^* is at the lower right corner of C_j , and C_{j+1} is above C_j , then, p_{j+1}^* is above p_j^* , that is, $b_2(J + 1) = 1$; (2) that if p_j^* is at the upper left corner of C_j , and if C_{j+1} is the cell to the right of C_j , then p_{j+1}^* is to the right of p_j^* , that is, $b_1(J + 1) = 1$.

Now note that the trajectory of p_j^* is being determined by greedy minimization of the function $f(x, y) = \max(2^{-\delta_1 x}, 2^{-\delta_2 y})$ by paths through integer lattice points. Below the ray S , $(\partial/\partial x)f(x, y) = 0$. We conclude that unit moves in the x -direction are useless when one is below S . On the other hand, below S , $(\partial/\partial y)f(x, y) < 0$. So a unit move in the y -direction if it is available, is useful. Above the ray S , the situation is reversed: $(\partial/\partial y)f(x, y) = 0$. We conclude that any move in the y -direction is useless when one is above S . But a unit move in the x -direction, if it is available, is useful.

Suppose one is in case (1) of the above paragraph. Then one knows that the upper right vertex of C_j is below or on the ray S . It follows that a full unit move in the y direction is available and useful. The greedy algorithm will certainly take it, and case (1) is established.

Suppose one is in case (2) of the above paragraph. Then one knows that the upper right vertex of C_j is above or on the ray S . It follows that a full unit move in the x direction is both available and useful. The greedy algorithm will certainly take it, and case (2) is established.

A.4. Proof of Lemma 8.10. Define

$$N(\varepsilon) = \sup\{\#\{i: |\theta_i(f, \text{BAB}(\delta_1, \delta_2))| \geq \varepsilon\}: f \in \mathcal{F}_p^{\delta_1, \delta_2}(C)\}.$$

The property in question amounts to the assertion that

$$(A.2) \quad N(\varepsilon)^{\rho+1/2} \varepsilon \leq KC, \quad \forall \varepsilon > 0.$$

By Corollary 8.7, there are constants $\gamma_j = \gamma_j(C)$ so that for $f \in \mathcal{F}_p^{\delta_1, \delta_2}(C)$, the coefficients in $\text{BAB}(\delta_1, \delta_2)$ obey

$$\left(\sum_{\mathcal{R}(j)} |\alpha_R|^p \right)^{1/p} \leq \gamma_j.$$

Now define

$$n(\varepsilon, d, \gamma) = \sup\{\#\{i: |\theta_i| \geq \varepsilon\}: \theta \in \mathbf{R}^d, \|\theta\|_{l^p} \leq \gamma\}.$$

Then

$$N(\varepsilon) \leq 1 + \sum_{j \geq 0} n(\varepsilon, 2^j, \gamma_j),$$

where the γ_j are as above. Easy calculations [see (10.5) and (10.6)] yield $n(\varepsilon, d, \gamma) = \min(d, (\gamma/\varepsilon)^p)$; from $\gamma_j = C2^{-j(\rho+1/2-1/p)}$ we get (A.2).

A.5. Proof of Lemma 8.11. The proof is an application of the following fact, called the ‘‘incompressibility of Hypercubes’’ in Donoho (1993). Suppose that \mathcal{H} is an orthogonal hypercube symmetric about zero; then it can be written $\sum_j \alpha_j g_j$ where the g_j are orthogonal and the α vary throughout the cube $|\alpha_j| \leq \varepsilon$. We call any basis starting with elements g_1, g_2, \dots, g_m a natural basis for \mathcal{H} . In that basis, \mathcal{H} is rotated so that the axes cut orthogonally through its faces.

Let \mathcal{B} be a natural basis for such a \mathcal{H} and let $\Theta = \Theta(\mathcal{H}, \mathcal{B})$ be the body of coefficients of \mathcal{H} in that basis. Let U be any orthogonal matrix. Then for absolute constants $c(p)$ and $0 < p \leq 2$,

$$\sup_{\theta \in \Theta} \|U\theta\|_{l^p} \geq c(p) \sup_{\theta \in \Theta} \|\theta\|_{l^p}.$$

In Donoho (1993), this is shown to be a consequence of Khintchine’s inequality

To use this, we argue by contradiction. Suppose that the hypotheses of the lemma hold, and yet for a certain basis \mathcal{B}^* , $\tau^*(\Theta(\mathcal{F}, \mathcal{B}^*)) = \tau - \varepsilon$ where $\varepsilon > 0$. Then for $0 < \delta < \varepsilon$, we have the weak-type inclusion $\Theta(\mathcal{F}, \mathcal{B}^*) \subset w l^{\tau-\delta}$. Equally, we have the stronger inclusion $\Theta(\mathcal{F}, \mathcal{B}^*) \subset l^{\tau-\delta}$.

Let \mathcal{H}_j be the j th hypercube in the sequence posited by the theorem, and let \mathcal{B}_j be a natural basis for \mathcal{H}_j . There is an orthogonal matrix U_j so that $\theta(f, \mathcal{B}^*) = U_j \theta(f, \mathcal{B}_j)$:

$$\begin{aligned} \sup_{h \in \mathcal{H}_j} \|\theta(h, \mathcal{B}^*)\|_{l^{\tau-\delta}} &= \sup_{h \in \mathcal{H}_j} \|U_j \theta(h, \mathcal{B}_j)\|_{l^{\tau-\delta}} \\ &\geq c(\tau - \delta) \sup_{h \in \mathcal{H}_j} \|\theta(h, \mathcal{B}_j)\|_{l^{\tau-\delta}}. \end{aligned}$$

On the other hand,

$$\begin{aligned} \sup_{h \in \mathcal{H}_j} \|\theta(h, \mathcal{B}_j)\|_{l^{\tau-\delta}} &= m_j^{1/\tau-\delta} \varepsilon_j \\ &= (m_j^{1/\tau} \varepsilon_j) m_j^{1/(\tau-\delta)-1/\tau} \\ &\geq c_0 m_j^{1/(\tau-\delta)-1/\tau} \rightarrow \infty. \end{aligned}$$

Hence $\Theta(\mathcal{F}, \mathcal{B}^*) \not\subset l^{\tau-\delta}$ for any $\delta > 0$. This contradiction proves Lemma 8.10.

A.6. Proof of Lemma 8.12.

The construction. Let g be a smooth function on \mathbf{R}^2 supported inside the unit square $[0, 1]^2$, whose support contains the half-square $[1/2, 3/4]^2$. Suppose that $\|(\partial/\partial x)g\|_{L^\infty} \leq \xi$ and that $\|(\partial/\partial y)g\|_{L^\infty} \leq \xi$. Suppose also that $\|g\|_{L^2} = 1$.

Let $\mathcal{A}(j)$ be the tiling of $[0, 1]$ selected at level j by $\text{BAB}(\delta_1, \delta_1)$. As this is a spatially homogeneous basis, all tiles are congruent. For an $R \in \mathcal{A}(j)$, let g_R denote the translation and dilation of g so that it just fits inside R , that is, $\text{supp}(g_R) \subset R$ and $R/2 \subset \text{supp}(g_R)$ where $R/2$ denotes the rectangle with the same center homothetically shrunk by a factor of 50 percent.

Let $\varepsilon_j = C2^{-j(\rho+1/2)}/(6\xi)$; define

$$\mathcal{H}_j = \left\{ \sum_{R \in \mathcal{A}(j)} \alpha_R g_R : |\alpha_R| \leq \varepsilon_j \right\}.$$

PROPERTY 1. We first note that \mathcal{H}_j obeys the dimension inequality assumed in the statement of the lemma, with $K = 1/(6\xi)$. Set $\rho = \delta_1\delta_2/(\delta_1 + \delta_2)$ and $\tau = 2/(2\rho + 1)$. With $m_j = 2^j$ the dimension of \mathcal{H}_j and ε_j the sidelength, one gets

$$m_j^{1/\tau} \varepsilon_j = C_0 > 0,$$

with $C_0 = C/(6\xi)$.

PROPERTY 2. The key claim about \mathcal{H}_j is the embedding $\mathcal{H}_j \subset \mathcal{F}_p^{\delta_1, \delta_2}(C)$: for any $f \in \mathcal{H}_j$,

$$\begin{aligned} \sup_{0 < h < 1} h^{-\delta_1} \|D_h^1 f\|_{L^p(Q_h^1)} &\leq C, \\ \sup_{0 < h < 1} h^{-\delta_2} \|D_h^2 f\|_{L^p(Q_h^2)} &\leq C. \end{aligned} \tag{A.3}$$

We prove the first inequality only, starting with estimates for differences of g_R . Let R be of side $2^{-j_1} \times 2^{-j_2}$.

Let $h > 2^{-j_1}$, and let R_h denote the translation of R by “to the left by h .” Then if $(x, y) \in R_h$, $D_h^1 g_R(x, y) = g_R(x + h, y)$, while if $(x, y) \in R$, $D_h^1 g_R(x, y) = -g_R(x, y)$. Note further that R_h is not generally part of the tiling $\mathcal{A}(j)$, but instead overlaps with two tiles, R_h^- and R_h^+ , say. Let $b_R(x, y) = g_R(x + h, y)1_{R_h^-}(x, y)$, and $c_R(x, y) = g_R(x + h, y)1_{R_h^+}(x, y)$. Then $D_h^1 g_R = a_R + b_R + c_R$, where a_R is supported in R , and b_R and c_R are

supported in R_h^\pm . We have for each R ,

$$\|\alpha_R\|_\infty, \|b_R\|_\infty, \|c_R\|_\infty \leq \xi 2^{j/2}.$$

Now consider the case $0 < h \leq 2^{-j_1}$. Let $b_R(x, y) = g_R(x + h, y)1_{R_h}(x, y)$, and $c_R(x, y) = 0$ and set $R_h^+ = R_h^- = R_h$. Then $D_h^1 g_R = \alpha_R + b_R + c_R$, where α_R is supported in R and b_R and c_R are supported in R_h^\pm . We have

$$\|\alpha_R\|_\infty, \|b_R\|_\infty, \|c_R\|_\infty \leq \min(h2^{j_1}, 1)\xi 2^{j/2}.$$

Now consider increments of $f = \sum_R \alpha_R g_R$. Rearrange the terms to have common support

$$D_h^1 f = \sum_R \alpha_R \alpha_R + \alpha_{R_h^-} b_R + \alpha_{R_h^+} c_R.$$

Now

$$\|(\alpha_{R_h^-})\|_{l^p}, \|(\alpha_{R_h^+})\|_{l^p} \leq \|(\alpha_R)\|_{l^p}.$$

By assumption, $1/p < \rho + 1/2$; as $\rho \leq 1/2$ we have $p > 1$, and the triangle inequality gives

$$\begin{aligned} \|D_h^1 f\|_p &\leq \left\| \sum_R \alpha_R \alpha_R \right\|_p + \left\| \sum \alpha_{R_h^-} b_R \right\|_p + \left\| \sum_R \alpha_{R_h^+} c_R \right\|_p \\ &\leq \|(\alpha_R)\|_{l^p} \max_R \|\alpha_R\|_\infty + \|(\alpha_{R_h^-})\|_{l^p} \max_R \|b_R\|_\infty + \|(\alpha_{R_h^+})\|_{l^p} \max_R \|c_R\|_\infty \\ &\leq 3\xi \min(h2^{j_1}, 1) \|(\alpha_R)\|_{l^p}. \end{aligned}$$

Hence

$$\begin{aligned} \sup_{0 < h < 1} h^{-\delta_1} \|D_h^1 f\|_p &\leq \sup_h h^{-\delta_1} 3\xi \|(\alpha_R)\|_{l^p} \min(h2^{j_1}, 1) 2^{j(1/2-1/p)} \\ &= 3\xi \|(\alpha_R)\|_{l^p} 2^{j(1/2-1/p)} \sup_h h^{-\delta_1} \min(h2^{j_1}, 1) \\ &= 3\xi \|(\alpha_R)\|_{l^p} 2^{j(1/2-1/p)} 2^{j_1 \delta_1}. \end{aligned}$$

Now from the proof of Lemma 8.6 we know that for $\text{BAB}(\delta_1, \delta_2)$,

$$2^{j_1 \delta_1} \leq 2 \cdot 2^{j\rho},$$

we conclude that

$$\sup_{0 < h < 1} h^{-\delta_1} \|D_h^1 f\|_p \leq 6\xi 2^{j(\rho+1/2-1/p)} \|\alpha\|_p \leq C.$$

This establishes (A.3).

A.7. Proof of Lemma 9.3. Recall the proof of Lemma 8.12. Let $\varepsilon > 0$ be given, and pick j so that the values $\varepsilon_j, \varepsilon_{j+1}$ defined in that lemma satisfy

$$\varepsilon_{j+1} < \varepsilon \leq \varepsilon_j.$$

Construct the hypercube \mathcal{K}_j exactly as in Lemma 8.12, only using side-length ε in place of ε_j .

We first note that the generating elements g_R are orthogonal with respect to the sampling measure l_N^2 , because they are disjointly supported. We also note that because of the dyadic structure of the sampling and the congruency of the hypercubes, each g_R has the same l_N^2 norm as every other g_R . Call this

norm $\delta = \delta(\varepsilon)$. Finally, we note that

$$\delta = \varepsilon \left(\frac{1}{M_1 M_2} \sum g(x_i)^2 \right)^{1/2} / \|g\|_{L^2[0,1]^2},$$

where the sum is over an $M_1 \times M_2$ array of grid points, where $M_i = 2^{J-j_i(J)}$. Hence \mathcal{H}_j is an orthogonal hypercube for l_N^2 . The asymptotics of the side-length can be derived from the fact that g is nice, the grid is becoming finer as j increases, and so the indicated sum converges to the corresponding integral, whence

$$\delta = \varepsilon(1 + o(1)).$$

The hypercube \mathcal{H}_j that results has two properties: first,

$$m_j^{1/\tau} \varepsilon = C_1(\varepsilon),$$

where

$$C_1(\varepsilon) = C_0(\varepsilon/\varepsilon_j) > C_0(\varepsilon_{j+1}/\varepsilon_j) = C/6\xi 2^{-(\rho+1/2)}.$$

Hence the dimensionality of the hypercube obeys (9.4), with $K = (6\xi 2^{(\rho+1/2)})^\tau$.

Second,

$$\mathcal{H}_j \subset \mathcal{F}_p^{\delta_1, \delta_2}(C).$$

This inclusion follows exactly as in Lemma 8.12.

Acknowledgments. This paper was stimulated by some interesting conversations about CART and wavelets which the author had with Joachim Engel at Oberwolfach, March 1995. It is a pleasure also to acknowledge conversations and helpful comments of A. Cohen, R. R. Coifman, R. A. DeVore, I. M. Johnstone, and V. N. Temlyakov. Thanks to Helen Tombropoulos for an efficient and enthusiastic typing job.

REFERENCES

- ALPERT, B., BEYLKIN, G., COIFMAN, R. and ROKHLIN, V. (1993). Wavelet-like bases for the fast solution of second-kind integral equations. *SIAM J. Sci. Comput.* **14** 159–184.
- BARRON, A. and COVER, T. (1991). Minimum complexity density estimation. *IEEE Trans. Inform. Theory* **37** 1034–1054.
- BENNETT, N. (1995). Fast algorithms for best anisotropic Walsh bases, and relatives. Ph.D. dissertation, Yale Univ.
- BIRGÉ, L. and MASSART, P. (1994). From model selection to adaptive estimation. In *Festschrift for Lucien Le Cam* (D. Pollard, E. Torgersen and G. Yang, eds.) 55–58. Springer, New York.
- BREIMAN, L., FRIEDMAN, J., OLSHEN, R. and STONE, C. J. (1983). *Classification and Regression Trees*. Wadsworth, CA.
- COIFMAN, R. R., MEYER, Y., QUAKE, S. and WICKERHAUSER, M. V. (1994). Wavelet analysis and signal processing. In *Wavelets and Their Applications* (J. S. Byrnes, J. L. Byrnes, K. A. Hargreaves and K. Berry, eds.) 363–380. Kluwer, Boston.
- COIFMAN, R. R. and WICKERHAUSER, M. V. (1992). Entropy-based algorithms for best-basis selection. *IEEE Trans. Inform. Theory* **38** 713–718.
- DEVORE, R. A. (1994). Adaptive wavelet bases for image compression. In *Wavelets, Images, and Surface Fitting* (P. J. Laurent, A. Le Mehauté, L. L. Schumaker, eds.) A. K. Peters, Boston.

- DONOHO, D. L. (1993). Unconditional bases are optimal bases for data compression and for statistical estimation. *Appl. Comput. Harmon. Anal.* **1** 100–115.
- DONOHO, D. L. (1995a). Abstract statistical estimation and modern harmonic analysis. In *Proceedings 1994 Int. Cong. Math.* 997–1005. Birkhäuser, Basel.
- DONOHO, D. L. (1995b). CART and best-ortho-basis: a connection. Technical report, Dept. Statistics, Stanford Univ. Available at <ftp://stat.stanford.edu/reports/donoho/cart.ps>
- DONOHO, D. L. (1995c). De-noising by soft thresholding. *IEEE Trans. Inform. Theory* **41** 613–627.
- DONOHO, D. L. (1996). Unconditional bases and bit-level compression. *Appl. Comput. Harmon. Anal.* **3** 388–392.
- DONOHO, D. L., DYN, N., LEVIN, D. and YU, T. P. Y. (1996). Smooth multiwavelet duals of Alpert bases by moment-interpolation, with applications to recursive partitioning. Unpublished manuscript.
- DONOHO, D. L. and JOHNSTONE, I. M. (1994a). Empirical atomic decomposition. Unpublished manuscript.
- DONOHO, D. L. and JOHNSTONE, I. M. (1994b). Ideal de-noising in a basis chosen from a library of orthonormal bases. *C.R. Acad. Sci. Paris Sér. I Math.* **319** 1317–1322.
- DONOHO, D. L. and JOHNSTONE, I. M. (1994c). Ideal spatial adaptation via wavelet shrinkage. *Biometrika* **81** 425–455.
- DONOHO, D. L. and JOHNSTONE, I. M. (1994d). Minimax risk over l^q balls. *Probab. Theory Related Fields* **99** 277–303.
- DONOHO, D. L., LIU, and MACGIBBON, B. (1990). Minimax risk over hyperrectangles, and implications. *Ann. Statist.* **18** 1416–1437.
- DONOHO, D. L., JOHNSTONE, I. M., KERKYACHARIAN, G. and PICARD, D. (1995). Wavelet shrinkage: asymptopia? *J. Roy. Statist. Soc. Ser. B* **57** 301–369.
- ENGEL, J. (1994). A simple wavelet approach to nonparametric regression from recursive partitioning schemes. *J. Multivariate Anal.* **49** 242–254.
- FOSTER, D. and GEORGE, E. I. (1996). The risk inflation factor in multiple linear regression. Unpublished manuscript.
- KASHIN, B. S. (1985). On approximation properties of complete orthonormal systems. *Trudy Math. Inst. Steklov* **172**. (In Russian.) [Trans. *Proc. Steklov Inst. Math.* 1987 207–211.]
- KASHIN, B. S. and TEMLYAKOV, V. N. (1994). On best m -term approximations and the entropy of sets in the space L^1 . *Mat. Zametki* **56** 57–86. (In Russian.) [Trans. *Math. Notes* **56** 1137–1157.]
- MALLAT, S. and ZHANG, Z. (1993). Matching pursuit in a time-frequency dictionary. *IEEE Trans. Signal Processing* **41** 3397–3415.
- NEMIROVSKII, A. S. (1986). Nonparametric estimation of smooth regression functions. *Soviet J. Comput. Systems Sci.* **23** 1–11.
- NEMIROVSKII, A. S., POLYAK, B. T. and TSYBAKOV, A. B. (1985). The rate of convergence of nonparametric estimates of maximum likelihood type. *Problemy Peredachi Informat-sii* **21** 17–33. (In Russian.)
- NEUMANN, M. H. and VON SACHS, R. (1995). Wavelet thresholding in anisotropic function classes and application to adaptive estimation of evolutionary spectra. Technical report, Berichte der AG Technomathematik 132, Univ. Kaiserslautern.
- NIKOL'SKII, S. M. (1969). *Approximation of Functions of Several Variables and Imbedding Theorems*. Springer, New York.
- TEMLYAKOV, V. N. (1993). *Approximation of Periodic Functions*. Nova Press.
- THIELE, C. M. and VILLEMOS, L. F. (1996). A fast algorithm for adapted time-frequency tilings. *Appl. Comput. Harmon. Anal.* **3** 91–99.

DEPARTMENT OF STATISTICS
 SEQUOIA HALL
 STANFORD UNIVERSITY
 STANFORD, CALIFORNIA 94305-4065
 E-MAIL: donoho@stat.stanford.edu