

MULTISCALE GEOMETRIC METHODS FOR DATA SETS II: GEOMETRIC MULTI-RESOLUTION ANALYSIS

WILLIAM K. ALLARD, GUANGLIANG CHEN, AND MAURO MAGGIONI

ABSTRACT. Data sets are often modeled as samples from a probability distribution in \mathbb{R}^D , for D large. It is often assumed that the data has some interesting low-dimensional structure, for example that of a d -dimensional manifold \mathcal{M} , with d much smaller than D . When \mathcal{M} is simply a linear subspace, one may exploit this assumption for encoding efficiently the data by projecting onto a dictionary of d vectors in \mathbb{R}^D (for example found by SVD), at a cost $(n + D)d$ for n data points. When \mathcal{M} is nonlinear, there are no “explicit” and algorithmically efficient constructions of dictionaries that achieve a similar efficiency: typically one uses either random dictionaries, or dictionaries obtained by black-box global optimization. In this paper we construct data-dependent multi-scale dictionaries that aim at efficiently encoding and manipulating the data. Their construction is fast, and so are the algorithms that map data points to dictionary coefficients and vice versa, in contrast with L^1 -type sparsity-seeking algorithms, but alike adaptive nonlinear approximation in classical multiscale analysis. In addition, data points are guaranteed to have a compressible representation in terms of the dictionary, depending on the assumptions on the geometry of the underlying probability distribution.

1. INTRODUCTION

We construct Geometric Multi-Resolution Analyses for analyzing intrinsically low-dimensional point clouds in high-dimensional spaces, modeled as samples from a probability distribution supported on d -dimensional set \mathcal{M} (in particular, a manifold) embedded in \mathbb{R}^D , in the regime $d \ll D$. This setting has been recognized as important in various applications, ranging from the analysis of sounds, images (RGB or hyperspectral, [1]), to gene arrays, EEG signals [2], and other types of manifold-valued data [3], and has been at the center of much investigation in the applied mathematics [4, 5, 6] and machine learning communities during the past several years. This has lead to a flurry of research on several problems, old and new, such as estimating the intrinsic dimensionality of point clouds [7, 8, 9, 10, 11, 12], parametrizing sampled manifolds [4, 13, 14, 15, 16, 17, 18, 19, 20], constructing dictionaries tuned to the data [21, 22] or for functions on the data [23, 24, 25, 26], and their applications to machine learning and function approximation [27, 28, 29, 30].

We focus on obtaining multi-scale representations in order to organize the data in a natural fashion, and obtain efficient data structures for data storage, transmission, manipulation, at different levels of precision that may be requested or needed for

Date: September 7, 2011.

Key words and phrases. Multiscale Analysis. Wavelets. Data Sets. Point Clouds. Frames. Sparse Approximation. Dictionary Learning.

The authors thank E. Monson for useful discussions.

GC was partially supported by ONR N00014-07-1-0625 and NSF CCF 0808847.

MM is grateful for partial support from DARPA, NSF, ONR, and the Sloan Foundation.

particular tasks. This work ties with a significant amount of recent work in different directions: (a) Harmonic analysis and efficient representations of signals; (b) Data-adaptive signal representations in high dimensional spaces and dictionary learning; (c) Hierarchical structures for organization of data sets; (d) Geometric analysis of low-dimensional sets in high-dimensional spaces.

Harmonic analysis and efficient representations of signals. Representations of classes of signals and data have been an important branch of research in multiple disciplines. In harmonic analysis, a linear infinite-dimensional function space \mathcal{F} typically models the class of signals of interest, and linear representations in the form $f = \sum_i \alpha_i \phi_i$, for $f \in \mathcal{F}$ in terms of a dictionary of atoms $\Phi := \{\phi_i\} \subseteq \mathcal{F}$ are studied. Such dictionaries may be bases or frames, and are constructed so that the sequence of coefficients $\{\alpha_i\}_i$ has desirable properties, such as some form of sparsity, or a distribution highly concentrated at zero. Requiring sparsity of the representation is very natural from the viewpoints of statistics, signal processing, and interpretation of the representation. This, in part, motivated the construction of Fourier-like bases, wavelets, wedgelets, ridgelets, curvelets etc... [31, 32, 33], just to name a few. Several such dictionaries are proven to provide optimal representations (in a suitably defined sense) for certain classes of function spaces (e.g. some simple models for images) and/or for operators on such spaces. While orthogonal dictionaries were originally preferred (e.g. [34]), a trend developed towards over-complete dictionaries (e.g. frames [34, 35] and references therein) and libraries of dictionaries (e.g. wavelet and cosine packets [31], multiple dictionaries [36], fusion frames [37]), for which the set of coefficients $(\alpha_i)_i$ needed to represent a signal f is typically non-unique. Fast transforms, crucial in applications, have often been considered a fundamental hallmark of several of the transforms above, and was usually achieved through a multi-scale organization of the dictionaries.

Data-adaptive signal representation and dictionary learning. A more recent trend [33, 38, 21, 39, 40, 22], motivated by the desire to model classes of signals that are not well-modeled by the linear structure of function spaces, has been that of *constructing data-adapted dictionaries*: an algorithm is allowed to see samples from a class of signals \mathcal{F} (not necessarily a linear function space), and constructs a dictionary $\Phi := \{\phi_i\}_i$ that optimizes some functional, such as the sparsity of the coefficients for signals in \mathcal{F} . The problem becomes being able to construct the dictionary Φ , typically highly over-complete, so that, given $f \in \mathcal{F}$, a rapid computation of the “best” (e.g. sparsest) coefficients $(\alpha_i)_i$ so that $f = \sum_i \alpha_i \phi_i$ is possible, and $(\alpha_i)_i$ is sparse. The problem of constructing Φ with the properties above, given a sample $\{f_n\}_n \subseteq \mathcal{F}$, is often called *dictionary learning*, and has been at the forefront of much recent research in harmonic analysis, approximation theory, imaging, vision, and machine learning: see [38, 21, 39, 40, 22] and references therein for constructions and applications.

There are several parameters in this problem: given training data from \mathcal{F} , one seeks Φ with I elements, such that every element in the training set may be represented, up to a certain precision ϵ , by at most m elements of the dictionary. The smaller I and m are, for a given ϵ , the better the dictionary.

Several current approaches may be summarized as follows [41]: consider a finite training set of signals $X_n = \{x_i\}_{i=1}^n \subset \mathbb{R}^D$, which we may represent by a $\mathbb{R}^{D \times n}$

matrix, and optimize the cost function

$$(1.1) \quad f_n(\Phi) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, \Phi)$$

where $\Phi \in \mathbb{R}^{D \times I}$ is the dictionary, and ℓ a loss function, for example

$$(1.2) \quad \ell(x, \Phi) := \min_{\alpha \in \mathbb{R}^I} \frac{1}{2} \|x - \Phi\alpha\|_{\mathbb{R}^D}^2 + \lambda \|\alpha\|_1$$

where λ is a regularization parameter. This is basis pursuit [33] or lasso [42]. One typically adds constraints on the size of the columns of Φ , for example $\|\phi_i\|_{\mathbb{R}^D} \leq 1$ for all i , which we can write as $\Phi \in \mathcal{C}$ for some convex set \mathcal{C} . The overall problem may then be written as a matrix factorization problem with a sparsity penalty:

$$(1.3) \quad \min_{\Phi \in \mathcal{C}, \alpha \in \mathbb{R}^{I \times n}} \frac{1}{2} \|X_n - \Phi\alpha\|_F^2 + \lambda \|\alpha\|_{1,1},$$

where $\|\alpha\|_{1,1} := \sum_{i_1, i_2} |\alpha_{i_1, i_2}|$. While for a fixed Φ the problem of minimizing over α is convex, and for fixed α the problem of minimizing over Φ 's is also convex, the joint minimization problem is non-convex, and alternate minimization methods are often employed. Overall, this requires minimizing a non-convex function over a very high-dimensional space. We refer the reader to [41] and references therein for techniques for attacking this optimization problem.

Constructions of such dictionaries (e.g. K-SVD [21], k -flats [22], optimization-based methods [41], Bayesian methods [39]) generally involve optimization or heuristic algorithms which are computationally intensive, do not shed light on the relationships between the dictionary size I , the sparsity of α , and the precision ϵ , and the resulting dictionary Φ is typically unstructured, and finding computationally, or analyzing mathematically, the sparse set of coefficients α may be challenging.

In this paper we construct data-dependent dictionaries based on a Geometric Multi-Resolution Analysis of the data. This approach is motivated by the intrinsically low-dimensional structure of many data sets, and is inspired by multi-scale geometric analysis techniques in geometric measure theory such as those in [43, 44], as well as by techniques in multi-scale approximation for functions in high-dimension [45, 46]. These dictionaries are structured in a multi-scale fashion (a structure that we call Geometric Multi-Resolution Analysis) and can be computed efficiently; the expansion of a data point on the dictionary elements is guaranteed to have a certain degree of sparsity m , and may be computed by a fast algorithm; the growth of the number of dictionary elements I as a function of ϵ is controlled depending on geometric properties of the data. We call the elements of these dictionaries *geometric wavelets*, since in some respects they generalize wavelets from vectors that analyze functions in linear spaces to affine vectors that analyze point clouds with possibly nonlinear geometry. The multi-scale analysis associated with geometric wavelets shares some similarities with that of standard wavelets (e.g. fast transforms, a version of two-scale relations, etc...), but is in fact quite different in many crucial respects. It is nonlinear, as it adapts to arbitrary nonlinear manifolds modeling the data space \mathcal{F} , albeit every scale-to-scale step is linear; translations or dilations do not play any role here, while they are often considered crucial in classical wavelet constructions. Geometric wavelets may allow the design of new algorithms for manipulating point clouds similar to those used for wavelets to manipulate functions.

The rest of the paper is organized as follows. In Sec. 2 we describe how to construct the geometric wavelets in a multi-scale fashion. We then present our algorithms in Sec. 3 and illustrate them on a few data sets, both synthetic and real-world, in Sec. 4. Sec. 5 introduces an orthogonal version of the construction; more variations or optimizations of the construction are postponed to Sec. 6. The next two sections discuss how to represent and compress data efficiently (Sec. 7) and computational costs (Sec. 8). A naive attempt at modeling distributions is performed in Sec. 9. Finally, the paper is concluded in Sec. 10 by pointing out some future directions.

2. CONSTRUCTION OF GEOMETRIC MULTI-RESOLUTION ANALYSES

Let (\mathcal{M}, ρ, μ) be a metric measure space with μ a Borel probability measure and $\mathcal{M} \subseteq \mathbb{R}^D$. In this paper we restrict our attention, in the theoretical sections, to the case when (\mathcal{M}, ρ, μ) is a smooth compact Riemannian manifold of dimension d isometrically embedded in \mathbb{R}^D , endowed with the natural volume measure; in the numerical examples, (\mathcal{M}, ρ, μ) will be a finite discrete metric space with counting measure, not necessarily obtained by sampling a manifold as above. We will be interested in the case when the “dimension” d of \mathcal{M} is much smaller than the dimension of the ambient space \mathbb{R}^D . While d is typically unknown in practice, efficient (multi-scale, geometric) algorithms for its estimation are available (see [8], which also contains many references to previous work on this problem), under additional assumptions on the geometry of \mathcal{M} .

Our construction of a Geometric Multi-Resolution Analyses (GMRA) consists of three steps:

1. A multi-scale geometric *tree decomposition* of \mathcal{M} into subsets $\{C_{j,k}\}_{k \in \mathcal{K}_j, j \in \mathbb{Z}}$.
2. A d -dimensional *affine approximation* in each dyadic cell $C_{j,k}$, yielding a sequence of approximating piecewise linear sets $\{\mathcal{M}_j\}$, one for each scale j .
3. A construction of low-dimensional *affine difference operators* that efficiently encode the differences between \mathcal{M}_j and \mathcal{M}_{j+1} .

This construction parallels, in a geometric setting, that of classical multi-scale wavelet analysis [34, 47, 48, 49, 50]: the nonlinear space \mathcal{M} replaces the classical function spaces, the piecewise affine approximation at each scale substitutes the linear projection on scaling function spaces, and the difference operators play the role of the classical linear wavelet projections. We show that when \mathcal{M} is a smooth manifold, guarantees on the approximation rates of \mathcal{M} by the \mathcal{M}_j may be derived (see Theorem 2.3 in Sec. 2.4), implying compressibility of the GMRA representation of the data.

We construct bases for the various affine operators involved, producing a hierarchically organized dictionary that is adapted to the data, which we expect to be useful in the applications discussed in the introduction.

2.1. Tree decomposition. Let $B_r^\mathcal{M}(x)$ be the ρ -ball inside \mathcal{M} of radius $r > 0$ centered at $x \in \mathcal{M}$. We start by a spatial multi-scale decomposition of the data set \mathcal{M} .

Definition 2.1. *A tree decomposition of a d -dimensional metric measure space (\mathcal{M}, ρ, μ) is a family of open sets in \mathcal{M} , $\{C_{j,k}\}_{k \in \mathcal{K}_j, j \in \mathbb{Z}}$, called dyadic cells, such that*

- (i) *for every $j \in \mathbb{Z}$, $\mu(\mathcal{M} \setminus \cup_{k \in \mathcal{K}_j} C_{j,k}) = 0$;*

- (ii) for $j' \geq j$ and $k' \in \mathcal{K}_{j'}$, either $C_{j',k'} \subseteq C_{j,k}$ or $\mu(C_{j',k'} \cap C_{j,k}) = 0$;
- (iii) for $j < j'$ and $k' \in \mathcal{K}_{j'}$, there exists a unique $k \in \mathcal{K}_j$ such that $C_{j',k'} \subseteq C_{j,k}$;
- (iv) each $C_{j,k}$ contains a point $c_{j,k}$ such that $B_{c_1 \cdot 2^{-j}}^{\mathcal{M}}(c_{j,k}) \subseteq C_{j,k} \subseteq B_{2^{-j}}^{\mathcal{M}}(c_{j,k})$, for a constant c_1 depending on intrinsic geometric properties of \mathcal{M} . In particular, we have $\mu(C_{j,k}) \sim 2^{-dj}$.

The construction of such tree decompositions is possible on spaces of homogeneous type [51, 52, 53]. Let \mathcal{T} be the tree structure associated to the decomposition above: for any $j \in \mathbb{Z}$ and $k \in \mathcal{K}_j$, we let $\text{children}(j, k) = \{k' \in \mathcal{K}_{j+1} : C_{j+1,k'} \subseteq C_{j,k}\}$. Note that $C_{j,k}$ is the disjoint union of its children $C_{j+1,k'}, k' \in \text{children}(j, k)$, due to (ii). We assume that $\mu(\mathcal{M}) \sim 1$ such that there is only one cell at the root of the tree with scale $\log_{2^d} \mu(\mathcal{M}) = 0$ (thus we will only consider $j \geq 0$). For every $x \in \mathcal{M}$, with abuse of notation we use (j, x) to represent the unique $(j, k(x)), k(x) \in \mathcal{K}_j$ such that $x \in C_{j,k(x)}$. The family of dyadic cells $\{C_{j,k}\}_{k \in \mathcal{K}_j}$ at scale j generates a σ -algebra \mathcal{F}_j . Functions measurable with respect to this σ -algebra are piecewise constant on each cell.

In this paper we will construct dyadic cells on i.i.d. μ -distributed samples $\{x_i\}_{i=1}^n$ from \mathcal{M} according to the following variation of the construction of diffusion maps [4, 54]: we connect each x_i to its k -nearest neighbors (default value is $k = 50$), with weights $W_{ij} = K(x_i, x_j) = e^{-\|x_i - x_j\|^2/\epsilon_i \epsilon_j}$, where ϵ_i is the distance between x_i and its $k/2$ -nearest neighbor, to obtain a weighted graph on the samples x_i (this construction is used and motivated in [55]). We then make use of METIS [56] to produce the multi-scale partitions $\{C_{j,k}\}$ and the dyadic tree \mathcal{T} above. In a future publication we will discuss how to use a variation of cover trees [57], which has guarantees in terms of both the quality of the decomposition and computational costs, and has the additional advantage of being easily updatable with new samples.

We may also construct the cells $C_{j,k}$ by intersecting Euclidean dyadic cubes in \mathbb{R}^D with \mathcal{M} : if \mathcal{M} is sufficiently regular and so is its embedding in \mathbb{R}^D (e.g. \mathcal{M} a smooth compact isometrically embedded manifold, or a dense set of samples, distributed according to volume measure, from it), then the properties in Definition 2.1 are satisfied for j large enough. In this case, a careful numerical implementation is needed in order to not be penalized by the ambient dimensionality (e.g. [58] and references therein).

Definition 2.2. *We define*

$$(2.1) \quad D(\mathcal{M}) = \{y \in \mathbb{R}^D : \exists! x \in \mathcal{M} \text{ such that } \|x - y\| = \min_{x' \in \mathcal{M}} \|x' - y\|\},$$

$$(2.2) \quad \text{tub}_r(\mathcal{M}) = \{y \in \mathbb{R}^D : d(y, \mathcal{M}) < r\}$$

and, following H. Federer [59],

$$(2.3) \quad \text{reach}(\mathcal{M}) = \sup\{r \geq 0 : \text{tub}_r(\mathcal{M}) \subset D(\mathcal{M})\}.$$

For $x \in \text{reach}(\mathcal{M})$, let x^* be the point in \mathcal{M} closest to x .

One may think of $\text{reach}(\mathcal{M})$ as the largest radius of a non-self-intersecting tube around \mathcal{M} , which depends on the embedding of \mathcal{M} in \mathbb{R}^D . This notion has appeared under different names, such as “condition number of a manifold”, in recent manifold learning literature [60, 61], as a key measure of the complexity of \mathcal{M} embedded in \mathbb{R}^D . In our setting, we require positive $\text{reach}(\mathcal{M})$ only in order to obtain uniform estimates, but for local (or pointwise) estimates only require $\text{reach}(B_z^{\mathcal{M}}(r))$, or $\text{reach}(\mathcal{M} \cap \mathbb{B}_z^D(r))$, for all r 's sufficiently small (depending on z).

2.2. Multiscale singular value decompositions and geometric scaling functions. The tools we build upon are classical in multi-scale geometric measure theory [62, 63, 53], especially in its intersection with harmonic analysis, and it is also related to adaptive approximation in high dimensions, see for example [45, 46] and references therein. An introduction to the use of such ideas for the estimation of intrinsic dimension of point clouds is in [8] and references therein (see [7, 64] for previous short accounts).

We will associate several gadgets to each dyadic cell $C_{j,k}$, starting with some geometric objects: the mean

$$(2.4) \quad c_{j,k} := \mathbb{E}_\mu[x|x \in C_{j,k}] = \frac{1}{\mu(C_{j,k})} \int_{C_{j,k}} x d\mu(x) \in \mathbb{R}^D$$

and the covariance operator restricted to $C_{j,k}$

$$(2.5) \quad \text{cov}_{j,k} = \mathbb{E}_\mu[(x - c_{j,k})(x - c_{j,k})^*|x \in C_{j,k}] \in \mathbb{R}^{D \times D}.$$

Here and in what follows points in \mathbb{R}^D are identified with D -dimensional column vectors. For a prescribed $d_{j,k}$ (e.g. $d_{j,k} = d$), let the rank- $d_{j,k}$ Singular Value Decomposition (SVD) [65] of $\text{cov}_{j,k}$ be

$$(2.6) \quad \text{cov}_{j,k} \approx \Phi_{j,k} \Sigma_{j,k} \Phi_{j,k}^*,$$

where $\Phi_{j,k}$ is an orthonormal $D \times d_{j,k}$ matrix and Σ is a diagonal $d_{j,k} \times d_{j,k}$ matrix. The linear projection operator onto the subspace $\langle \Phi_{j,k} \rangle$ spanned by the columns of $\Phi_{j,k}$ will be denoted by $P_{j,k}$. We let

$$(2.7) \quad \mathbb{V}_{j,k} := V_{j,k} + c_{j,k}, \quad V_{j,k} = \langle \Phi_{j,k} \rangle,$$

where $\langle A \rangle$ denotes the span of the columns of A , so that $\mathbb{V}_{j,k}$ is the affine subspace of dimension $d_{j,k}$ parallel to $V_{j,k}$ and passing through $c_{j,k}$. It is an approximate tangent space to \mathcal{M} at location $c_{j,k}$ and scale 2^{-j} ; and in fact it provides the best $d_{j,k}$ -dimensional planar approximation to \mathcal{M} in the least square sense:

$$(2.8) \quad \mathbb{V}_{j,k} = \operatorname{argmin}_\Pi \int_{C_{j,k}} \|x - \mathbb{P}_\Pi(x)\|^2 d\mu(x),$$

where Π is taken on the set of all affine $d_{j,k}$ -planes, and \mathbb{P}_Π is the orthogonal projection onto the affine plane Π . We think of $\{\Phi_{j,k}\}_{k \in \mathcal{K}_j}$ as the geometric analogue of a family of scaling functions at scale j , and therefore call *geometric scaling functions*. Let $\mathbb{P}_{j,k}$ be the associated affine projection

$$(2.9) \quad \mathbb{P}_{j,k}(x) := P_{j,k}(x - c_{j,k}) + c_{j,k} = \Phi_{j,k} \Phi_{j,k}^*(x - c_{j,k}) + c_{j,k}, \quad x \in C_{j,k}.$$

Then $\mathbb{P}_{j,k}(C_{j,k})$ is the projection of $C_{j,k}$ onto its local linear approximation, at least for $2^{-j} \lesssim \text{reach}(\mathcal{M})$.

We let

$$(2.10) \quad \mathcal{M}_j := \{\mathbb{P}_{j,k}(C_{j,k})\}_{k \in \mathcal{K}_j}$$

be a coarse approximation of \mathcal{M} at scale j , the geometric analogue to what the projection of a function onto a scaling function subspace is in wavelet theory. Under general conditions, $\mathcal{M}_j \rightarrow \mathcal{M}$ in the Hausdorff distance, as $j \rightarrow +\infty$. It is natural to define the nonlinear projection of \mathcal{M} onto \mathcal{M}_j by

$$(2.11) \quad x_j \equiv P_{\mathcal{M}_j}(x) := \mathbb{P}_{j,k}(x), \quad x \in C_{j,k}.$$

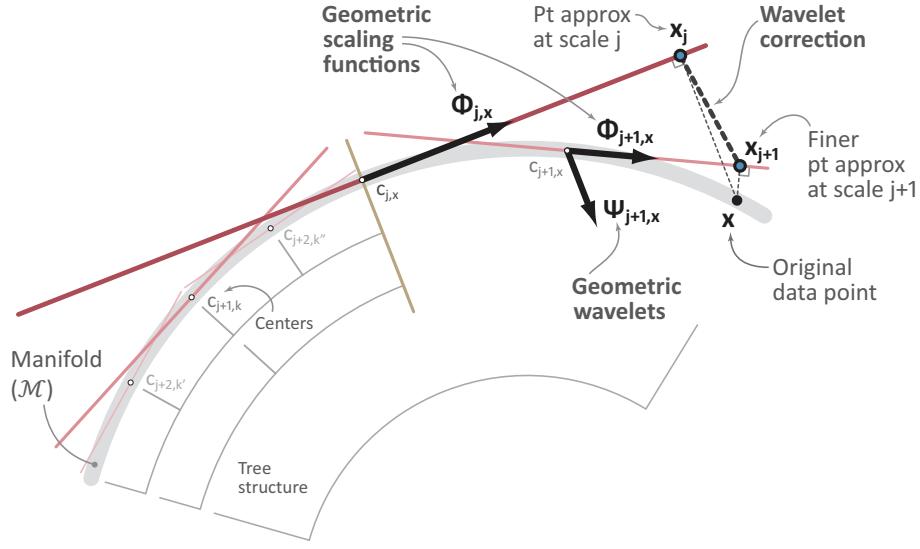


FIGURE 1. An illustration of the geometric wavelet decomposition. The centers $c_{j,x}$'s are represented as lying on \mathcal{M} while in fact they are only close (to second order) to \mathcal{M} , and the corresponding planes $V_{j,x}$ are represented as tangent planes, albeit they are only an approximation to them. Art by E. Monson.

2.3. Geometric wavelets. We would like to efficiently encode the difference needed to go from \mathcal{M}_j to \mathcal{M}_{j+1} , for $j \geq 0$. Fix $x \in \mathcal{M}$: the difference $x_{j+1} - x_j$ is a high-dimensional vector in \mathbb{R}^D , in general not contained in \mathcal{M}_{j+1} . However it may be decomposed into a sum of vectors in certain well-chosen low-dimensional spaces, which are shared across multiple points, in a multi-scale fashion. Recall that we use the notation (j, x) to denote the unique pair (j, k) , with $k \in \mathcal{K}_j$, such that $x \in C_{j,k}$. We proceed as follows: for $j \leq J-1$ we let

$$\begin{aligned}
Q_{\mathcal{M}_{j+1}}(x) &:= x_{j+1} - x_j \\
&= (x_{j+1} - \mathbb{P}_{j,x}(x_{j+1})) + (\mathbb{P}_{j,x}(x_{j+1}) - \mathbb{P}_{j,x}(x)) \\
&= (I - P_{j,x})(x_{j+1} - c_{j,x}) + P_{j,x}(x_{j+1} - x) \\
(2.12) \quad &= (I - P_{j,x}) \underbrace{(x_{j+1} - c_{j+1,x} + c_{j+1,x} - c_{j,x})}_{\in V_{j+1,x}} - P_{j,x}(x - x_{j+1}).
\end{aligned}$$

Let $W_{j+1,x}$ be the geometric wavelet subspace defined by

$$(2.13) \quad W_{j+1,x} := (I - P_{j,x}) V_{j+1,x},$$

$\Psi_{j+1,x}$ an orthonormal basis for $W_{j+1,x}$, that we will call a *geometric wavelet basis*, and $Q_{j+1,x}$ the orthogonal projection onto $W_{j+1,x}$. Clearly $\dim W_{j+1,x} \leq \dim V_{j+1,x} = d_{j+1,x}$. If we define the quantities

$$(2.14) \quad t_{j+1,x} := c_{j+1,x} - c_{j,x};$$

$$(2.15) \quad w_{j+1,x} := (I - P_{j,x}) t_{j+1,x};$$

$$(2.16) \quad Q_{j+1,x}(x) := Q_{j+1,x}(x - c_{j+1,x}) + w_{j+1,x},$$

then we may rewrite (2.12) as

$$\begin{aligned}
Q_{\mathcal{M}_{j+1}}(x) &= \underbrace{Q_{j+1,x}(x_{j+1} - c_{j+1,x})}_{\in W_{j+1,x}} + w_{j+1,x} - P_{j,x} \left(x - x_J + \sum_{l=j+1}^{J-1} (x_{l+1} - x_l) \right) \\
&= \mathbb{Q}_{j+1,x}(x_{j+1}) - P_{j,x} \sum_{l=j+1}^{J-1} (x_{l+1} - x_l) - P_{j,x}(x - x_J) \\
(2.17) \quad &= \mathbb{Q}_{j+1,x}(x_{j+1}) - P_{j,x} \sum_{l=j+1}^{J-1} Q_{\mathcal{M}_{l+1}}(x) - P_{j,x}(x - x_J).
\end{aligned}$$

Here $J \geq j+1$ is the index of the finest scale (and the last term vanishes as $J \rightarrow +\infty$, under general conditions). In terms of the geometric scaling functions and wavelets, the above may be written as

$$\begin{aligned}
x_{j+1} - x_j &= \Psi_{j+1,x} \Psi_{j+1,x}^*(x_{j+1} - c_{j+1,x}) + w_{j+1,x} - \Phi_{j,x} \Phi_{j,x}^* \sum_{l=j+1}^{J-1} Q_{\mathcal{M}_{l+1}}(x) \\
(2.18) \quad &\quad - \Phi_{j,x} \Phi_{j,x}^*(x - x_J).
\end{aligned}$$

This equation splits the difference $x_{j+1} - x_j$ into a component in $W_{j+1,x}$, a second component that only depends on the cell $(j+1, x)$ (but not on the point x *per se*), accounting for the translation of centers and lying in the orthogonal complement of $V_{j,x}$ but not necessarily in $W_{j+1,x}$, and a sum of terms which are projections on $V_{j,x}$ of differences in the same form $x_{l+1} - x_l$, but at finer scales. By construction we have the two-scale equation

$$(2.19) \quad P_{\mathcal{M}_{j+1}}(x) = P_{\mathcal{M}_j}(x) + Q_{\mathcal{M}_{j+1}}(x), \quad x \in \mathcal{M}$$

which can be iterated across scales, leading to a multi-scale decomposition along low-dimensional subspaces, with efficient encoding and algorithms. We think of $P_{j,k}$ as being attached to the node (j, k) of \mathcal{T} , and the $Q_{j+1,k'}$ as being attached to the edge connecting the node $(j+1, k')$ to its parent.

We say that the set of multi-scale piecewise affine operators $\{P_{\mathcal{M}_j}\}$ and $\{Q_{\mathcal{M}_{j+1}}\}$ form a *Geometric Multi-Resolution Analysis*, or GMRA for short.

2.4. Approximation for manifolds. We analyze the error of approximation to a d -dimensional manifold in \mathbb{R}^D by using geometric wavelets representation. The following result fully explains of the examples in Sec. 4.1.

Theorem 2.3. *Let (\mathcal{M}, ρ, μ) be a compact $C^{1+\alpha}$ Riemannian manifold of dimension d isometrically embedded in \mathbb{R}^D , with $\alpha \in (0, 1]$, and μ absolutely continuous with respect to the volume measure on \mathcal{M} . Let $\{P_{\mathcal{M}_j}, Q_{\mathcal{M}_{j+1}}\}$ be a GMRA for (\mathcal{M}, ρ, μ) . For any $x \in \mathcal{M}$, there exists a scale $j_0 = j_0(x)$ such that for any $j \geq j_0$ and any $p > 0$, if we let $d\mu_{j,x} := \mu(C_{j,x})^{-1} d\mu$,*

$$\begin{aligned}
\| \|z - P_{\mathcal{M}_j}(z)\|_{\mathbb{R}^D} \|_{L^p(C_{j,x}, d\mu_{j,x}(z))} &= \left\| z - P_{\mathcal{M}_{j_0}}(z) - \sum_{l=j_0}^{j-1} Q_{\mathcal{M}_{l+1}}(z) \right\|_{\mathbb{R}^D} \\
(2.20) \quad &\leq \|\kappa\|_{L^\infty(C_{j,x})} 2^{-(1+\alpha)j} + o(2^{-(1+\alpha)j}).
\end{aligned}$$

If $\alpha < 1$, $\kappa(x)$ depends on the $C^{1+\alpha}$ norm of a coordinate chart from $T_x(\mathcal{M})$ to $C_{j,x} \subseteq \mathcal{M}$.

If $\alpha = 1$, $\kappa(x) = \min(\kappa_1(x), \kappa_2(x))$, with

(2.21)

$$\kappa_1(x) := \frac{1}{2} \max_{i \in \{1, \dots, D-d\}} \|H_i(x)\|;$$

(2.22)

$$\kappa_2^2(x) := \max_{w \in \mathbb{S}^{D-d}} \frac{d(d+1)}{4(d+2)(d+4)} \left[\left\| \sum_{l=1}^{D-d} w_l H_l(x) \right\|_F^2 - \frac{1}{d+2} \left(\sum_{l=1}^{D-d} w_l \text{Tr}(H_l(x)) \right)^2 \right],$$

and the $D-d$ matrices $H_l(x)$ are the d -dimensional Hessians of \mathcal{M} at x .

This theorem describes the asymptotic decay of the geometric wavelet coefficients as a function of scale, and in particular it implies the compressibility of such coefficients. The decay depends on the smoothness of the manifold, and for C^2 manifolds it is quadratic in the scale; it saturates at C^2 , and for smoother manifolds we would have to use higher order geometric wavelets. We do not consider them here as the data sets we consider do not seem to benefit from higher order constructions. More quantitatively, the asymptotic rate is affected by the constant κ , which combines the distortion of $d\mu$ compared to the volume measure, and a notion of L^2 curvature. Depending on the size of κ , which in general varies from location to location, it gives an error estimate for an adaptive thresholding scheme that would threshold small coefficients in the geometric wavelet expansion (see the third example in Section 4.1).

Observe that κ_2 can be smaller than κ_1 (by a constant factor) or larger (by factors depending on d^2), depending on the spectral properties and commutativity relations between the Hessians H_l . κ_2^2 may be unexpectedly small, in the sense that it may scale as $d^{-2}r^4$ as a function of d and r , as observed in [8], because of concentration of measure phenomena.

Finally, we note that similar bounds may be obtained in $L^p(C_{j,x}, d\text{vol})$ simply by changing measure from $d\mu$ to $d\text{vol}$ and paying the price of replacing the constant κ by $\left\| \frac{d\mu}{d\text{vol}} \right\|_{L^\infty(C_{j,x})} \kappa$. This may also be achieved algorithmically with simple standard renormalizations (e.g. [4]).

The proof is postponed to the Appendix.

It is clear how to generalize the Theorem to unions of manifolds with generic intersections, at scales small enough around a point so that $C_{j,x}$ does not include intersections. Moreover, since the results are local, sets more general than manifolds may be considered as well: this is subject of a future report.

2.5. Non-manifold data and measures of approximation error. When constructing a GMRA for point-cloud data not sampled from manifolds, we may choose the dimension $d_{j,k}$ of the local linear approximating plane $\mathbb{V}_{j,k}$ by a criterion based on local approximation errors. Note that this affects neither the construction of geometric scaling functions, nor that of the wavelet subspaces and bases.

A simple measure for absolute error of approximation at scale j is:

$$\begin{aligned}
 \mathcal{E}_j^2 &= \int_{\mathcal{M}} \|P_{\mathcal{M}_j}(x) - x\|_{\mathbb{R}^D}^2 d\mu(x) = \sum_{k \in \mathcal{K}_j} \int_{C_{j,k}} \|P_{j,k}(x) - x\|_{\mathbb{R}^D}^2 d\mu|_{C_{j,k}}(x) \\
 &= \sum_{k \in \mathcal{K}_j} \mu(C_{j,k}) \frac{1}{\mu(C_{j,k})} \int_{C_{j,k}} \|P_{j,k}(x) - x\|_{\mathbb{R}^D}^2 d\mu|_{C_{j,k}}(x) \\
 (2.23) \quad &= \sum_{k \in \mathcal{K}_j} \mu(C_{j,k}) \sum_{l \geq d_{j,k}+1} \lambda_l(\text{cov}_{j,k}).
 \end{aligned}$$

We can therefore control \mathcal{E}_j by choosing $d_{j,k}$ based on the spectrum of $\text{cov}_{j,k}$. If we perform relative thresholding of $\text{cov}_{j,k}$, i.e. choose the smallest $d_{j,k}$ for which

$$(2.24) \quad \sum_{l \geq d_{j,k}+1} \lambda_l(\text{cov}_{j,k}) \leq \epsilon_j \sum_{l \geq 1} \lambda_l(\text{cov}_{j,k}),$$

for some choice of ϵ_j (e.g. $\epsilon_j = (c\theta^j) \vee \epsilon$ for some $\theta \in (0, 1)$ and $\epsilon > 0$), then we may upper bound the above as follows:

$$(2.25) \quad \mathcal{E}_j^2 \leq \sum_{k \in \mathcal{K}_j} \mu(C_{j,k}) \epsilon_j \|C_{j,k}\|_F^2 \leq \epsilon_j \|\mathcal{M}\|_F,$$

where $C_{j,k}$ and \mathcal{M} are thought of as matrices containing points in columns, and for a partitioned matrix $A = [A_1, A_2, \dots, A_r]$ and discrete probability measure μ on $\{1, \dots, r\}$ we define

$$(2.26) \quad \|\|A\|\|_F^2 := \sum_{i=1}^r \mu(\{i\}) \|A_i\|_F^2.$$

If we perform absolute thresholding of $\text{cov}_{j,k}$, i.e. choose the smallest $d_{j,k}$ for which $\sum_{l \geq d_{j,k}+1} \lambda_l(\text{cov}_{j,k}) \leq \epsilon_j$, then we have the rough bound

$$(2.27) \quad \mathcal{E}_j^2 \leq \sum_{k \in \mathcal{K}_j} \mu(C_{j,k}) \epsilon_j \leq \epsilon_j \cdot \mu(\mathcal{M}).$$

Of course, in the case of a d -dimensional \mathcal{C}^2 manifold \mathcal{M} with volume measure, if we choose $d_{j,k} = d$, by Theorem 2.3 we have

$$(2.28) \quad \mathcal{E}_j \lesssim \sum_{k \in \mathcal{K}_j} \mu(C_{j,k}) \|\kappa\|_\infty 2^{-2j} = \mu(\mathcal{M}) \|\kappa\|_\infty 2^{-2j}.$$

3. ALGORITHMS

We present in this section algorithms implementing the construction of the GMRA and the corresponding Geometric Wavelet Transform (GWT).

3.1. Construction of Geometric Multi-Resolution Analysis. The first step in the construction of the geometric wavelets is to perform a geometric nested partition of the data set, forming a tree structure. For this end, one may consider various methods listed below:

- (I). Use of METIS [56]: a multiscale variation of iterative spectral partitioning.

We construct a weighted graph as done for the construction of diffusion maps [4, 54]: we add an edge between each data point and its k nearest neighbors, and assign to any such edge between x_i and x_j the weight

```

GMRA = GeometricMultiResolutionAnalysis ( $X_n, \tau_0, \epsilon$ )
// Input:
//  $X_n$ : a set of  $n$  samples from  $\mathcal{M}$ 
//  $\tau_0$ : some method for choosing local dimensions
//  $\epsilon$ : precision
// Output:
// A tree  $\mathcal{T}$  of dyadic cells  $\{C_{j,k}\}$ , their local means  $\{c_{j,k}\}$  and bases  $\{\Phi_{j,k}\}$ ,
// together with a family of geometric wavelets  $\{\Psi_{j,k}\}, \{w_{j,k}\}$ 

Construct the dyadic cells  $C_{j,k}$  with centers  $\{c_{j,k}\}$  and form a tree  $\mathcal{T}$ .
 $J \leftarrow$  finest scale with the  $\epsilon$ -approximation property.
Let  $\text{cov}_{J,k} = |C_{J,k}|^{-1} \sum_{x \in C_{J,k}} (x - c_{J,k})(x - c_{J,k})^*$ , for  $k \in \mathcal{K}_J$ , and compute
SVD( $\text{cov}_{J,k} = \Phi_{J,k} \Sigma_{J,k} \Phi_{J,k}^*$  (where the dimension of  $\Phi_{J,k}$  is determined by  $\tau_0$ ).
for  $j = J - 1$  down to 0
    for  $k \in \mathcal{K}_j$ 
        Compute  $\text{cov}_{j,k}$  and  $\Phi_{j,k}$  as above.
        For each  $k' \in \text{children}(j, k)$ , construct the wavelet bases  $\Psi_{j+1,k'}$  and
        translations  $w_{j+1,k'}$ , according to (2.16),(2.13).
    end
end
For convenience, set  $\Psi_{0,k} := \Phi_{0,k}$  and  $w_{0,k} := c_{0,k}$  for  $k \in \mathcal{K}_0$ .

```

FIGURE 2. Pseudo-code for the construction of geometric wavelets

$e^{-\|x_i - x_j\|^2/\sigma}$. Here k and σ are parameters whose selection we do not discuss here (but see [55] for a discussion in the context of molecular dynamics data). In practice, we choose k between 10 and 50, and choose σ adaptively at each point x_i as the distance between x_i and its $\lfloor k/2 \rfloor$ nearest neighbor.

- (II). Use of cover trees [57].
- (III). Use of iterated PCA: at scale 1, compute the top d principal components of data, and partition the data based on the sign of the $(d + 1)$ -st singular vector. Repeat on each of the two partitions.
- (IV). Iterated k -means: at scale 1 partition the data based on k -means clustering, then iterate on each of the elements of the partition.

Each construction has pros and cons, in terms of performance and guarantees. For (I) we refer the reader to [56], for (II) to [57] (which also discussed several other constructions), for (III) and (IV) to [66]. Only (II) guarantees the needed properties for the cells $C_{j,k}$. However constructed, we denote by $\{C_{j,k}\}$ the family of resulting dyadic cells, and let \mathcal{T} be the associated tree structure, as in Definition 2.1.

In Fig. 2 we display pseudo-code for the construction of a GMRA for a data set X_n given a precision $\epsilon > 0$ and a method τ_0 for choosing local dimensions (e.g., using thresholds or a fixed dimension). The code first constructs a family of multi-scale dyadic cells (with local centers $c_{j,k}$ and bases $\Phi_{j,k}$), and then computes the geometric wavelets $\Psi_{j,k}$ and translations $w_{j,k}$ at all scales. In practice, we use METIS [56] to construct a dyadic (not 2^d -adic) tree \mathcal{T} and the associated cells $C_{j,k}$.

3.2. The Fast Geometric Wavelet Transform and its Inverse. For simplicity of presentation, we shall assume $x = x_J$; otherwise, we may first project x onto the local linear approximation of the cell $C_{J,x}$ and use x_J instead of x from now on.

```

 $\{q_{j,x}\} = \text{FGWT(GMRA, } x)$ 

// Input: GMRA structure,  $x \in \mathcal{M}$ 
// Output: A sequence  $\{q_{j,x}\}$  of wavelet coefficients

 $p_{J,x} = \Phi_{J,x}^*(x - c_{J,x})$ 
for  $j = J$  down to 1
     $q_{j,x} = (\Psi_{j,x}^* \Phi_{j,x}) p_{j,x}$ 
     $p_{j-1,x} = (\Phi_{j-1,x}^* \Phi_{J,x}) p_{J,x} + \Phi_{j-1,x}^*(c_{J,x} - c_{j-1,x})$ 
end
 $q_{0,x} = p_{0,x}$  (for convenience)

```

FIGURE 3. Pseudo-code for the Forward Geometric Wavelet Transform

```

 $\hat{x} = \text{IGWT(GMRA, } \{q_{j,x}\})$ 

// Input: GMRA structure, wavelet coefficients  $\{q_{j,x}\}$ 
// Output: Approximation  $\hat{x}$  at scale  $J$ 

 $Q_{J,x} = \Psi_{J,x} q_{J,x} + w_{J,x}$ 
for  $j = J-1$  down to 1
     $Q_j(x) = \Psi_{j,x} q_{j,x} + w_{j,x} + \Phi_{j-1,x} \Phi_{j-1,x}^* \sum_{\ell > j} Q_\ell(x)$ 
end
 $\hat{x} = \Psi_{0,x} q_{0,x} + w_{0,x} + \sum_{j>0} Q_j(x)$ 

```

FIGURE 4. Pseudo-code for the Inverse Geometric Wavelet Transform

That is, we will define $x_{j,J} = P_{\mathcal{M}_j}(x_J)$, for all $j < J$, and encode the differences $x_{j+1,J} - x_{j,J}$ using the geometric wavelets. Note also that $\|x_{j,J} - x_j\| \leq \|x - x_J\|$ at all scales.

The geometric scaling and wavelet coefficients $\{p_{j,x}\}, \{q_{j+1,x}\}$, for $j \geq 0$, of a point $x \in \mathcal{M}$ are chosen to satisfy the equations

$$(3.1) \quad P_{\mathcal{M}_j}(x) = \Phi_{j,x} p_{j,x} + c_{j,x};$$

$$(3.2) \quad Q_{\mathcal{M}_{j+1}}(x) = \Psi_{j+1,x} q_{j+1,x} + w_{j+1,x} - P_{j,x} \sum_{l=j+1}^{J-1} Q_{\mathcal{M}_{l+1}}(x).$$

The computation of the coefficients, from fine to coarse, is simple and fast: since we assume $x = x_J$, we have

$$(3.3) \quad \begin{aligned} p_{j,x} &= \Phi_{j,x}^*(x_J - c_{j,x}) = \Phi_{j,x}^*(\Phi_{J,x} p_{J,x} + c_{J,x} - c_{j,x}) \\ &= (\Phi_{j,x}^* \Phi_{J,x}) p_{J,x} + \Phi_{j,x}^*(c_{J,x} - c_{j,x}). \end{aligned}$$

Moreover the wavelet coefficients $q_{j+1,x}$ (defined in (3.2)) are obtained from (2.18):

$$(3.4) \quad q_{j+1,x} = \Psi_{j+1,x}^*(x_{j+1} - c_{j+1,x}) = (\Psi_{j+1,x}^* \Phi_{j+1,x}) p_{j+1,x}.$$

Note that $\Phi_{j,x}^* \Phi_{J,x}$ and $\Psi_{j+1,x}^* \Phi_{j+1,x}$ are both small matrices (at most $d_{j,x} \times d_{j,x}$), and are the only matrices we need to compute and store (once for all, and only up to a specified precision) in order to compute all the wavelet coefficients $q_{j+1,x}$ and the scaling coefficients $p_{j,x}$, given $p_{J,x}$ at the finest scale.

In Figs. 3 and 4 we display pseudo-codes for the computation of the Forward and Inverse Geometric Wavelet Transforms (F/IGWT). The input to FGWT is a

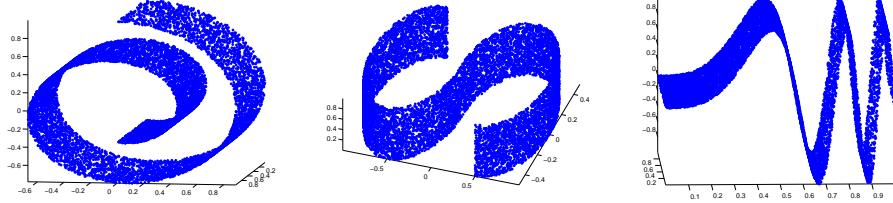


FIGURE 5. Toy data sets for geometric wavelets transform.

GMRA object, as returned by `GeometricMultiResolutionAnalysis`, and a point $x \in \mathcal{M}$. Its output is the wavelet coefficients of the point x at all scales, which are then used by IGWT for reconstruction of the point at all scales.

For any $x \in \mathcal{M}_J$, the set of coefficients

$$(3.5) \quad q_x = (q_{J,x}; q_{J-1,x}; \dots; q_{1,x}; p_{0,x})$$

is called the discrete *geometric wavelet transform* (GWT) of x . Letting $d_{j,x}^w = \text{rank}(\Psi_{j+1,x})$, the length of the transform is $d + \sum_{j>0} d_{j,x}^w$, which is bounded by $(J+1)d$ in the case of samples from a d -dimensional manifold (due to $d_{j,x}^w \leq d$).

Remark 3.1. Note that for the variation of the GMRA without adding tangential corrections (see Sec. 6.2), the algorithms above (as well as those in Sec. 5) can be simplified. First, in Fig. 2 we will not need to store the local bases functions $\{\Phi_{j,k}\}$. Second, the steps in Figs. 3 and 4 can be modified not to involve $\{\Phi_{j,k}\}$, similarly as in Figs. 17 and 18 of next section.

4. EXAMPLES

We conduct numerical experiments in this section to demonstrate the performance of the algorithm (i.e., Figs. 2, 3, 4).

4.1. Low-dimensional smooth manifolds. To illustrate the construction presented so far, we consider simple synthetic datasets: a *SwissRoll*, an *S-Manifold* and an *Oscillating2DWave*, all two-dimensional manifolds but embedded in \mathbb{R}^{50} (see Fig. 5). We apply the algorithm to construct the GMRA and obtain the forward geometric wavelet transform of the sampled data (10000 points, without noise) in Fig. 6. We use the manifold dimension $d_{j,k} = d = 2$ at each node of the tree when constructing scaling functions, and choose the smallest finest scale for achieving an absolute precision .001 in each case. We compute the average magnitude of the wavelet coefficients at each scale and plot it as a function of scale in Fig. 6. The reconstructed manifolds obtained by the inverse geometric wavelets transform (at selected scales) are shown in Fig. 7, together with a plot of relative approximation errors,

$$(4.1) \quad \mathcal{E}_{j,2}^{\text{rel}} = \sqrt{\frac{1}{n} \sum_{x \in X_n} \left(\frac{\|x - P_{j,x}(x)\|}{\|x\|} \right)^2},$$

where X_n is the training data of n samples. Both the approximation error and the magnitude of the wavelet coefficients decrease quadratically with respect to scale as expected.

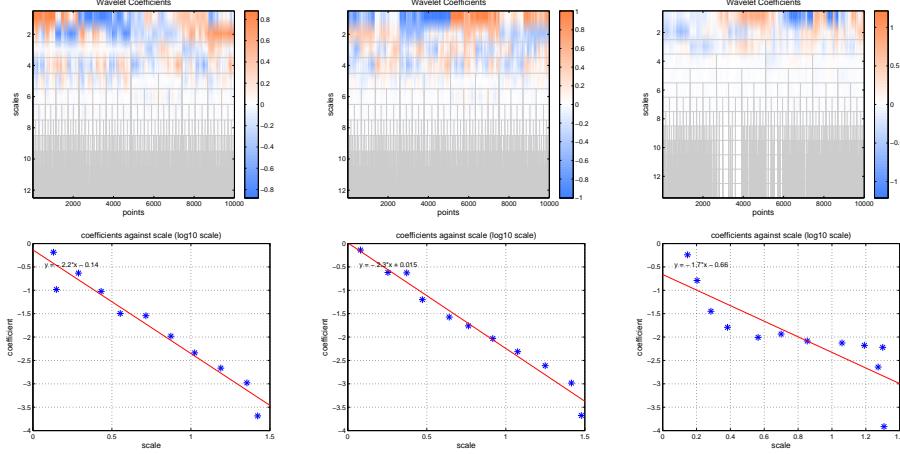


FIGURE 6. Top row: Wavelet coefficients obtained by the algorithm for the three data sets in Fig. 5. The horizontal axis indexes the points (arranged according to the tree), and the vertical axis multi-indexes the wavelet coefficients, from coarse (top) to fine (bottom) scales: the block of entries $(x, j), x \in C_{j,k}$ displays $\log_{10} |q_{j,x}|$, where $q_{j,x}$ is the vector of geometric wavelet coefficients of x at scale j (see Sec. 3). In particular, each row indexes multiple wavelet elements, one for each $k \in \mathcal{K}_j$. Bottom row: magnitude of wavelet coefficients decreasing quadratically as a function of scale.

We threshold the wavelet coefficients to study the compressibility of the wavelet coefficients and the rate of change of the approximation errors (using compressed wavelet coefficients). For this end, we use a smaller precision 10^{-5} so that the algorithm can examine a larger interval of thresholds. We first threshold the wavelet coefficients of the *Oscillating2DWave* data at the level .01 and plot in Fig. 8 the reduced matrix of wavelet coefficients and the corresponding best reconstruction of the manifold (i.e., at the finest scale). Next, we threshold the wavelet coefficients of all three data sets at different levels (from 10^{-5} to 1) and plot in Fig. 9 the compression and error curves.

4.2. Real data.

4.2.1. MNIST Handwritten Digits. We first consider the MNIST data set of images of handwritten digits¹, each of size 28×28 . We use the digits 0 and 1, and randomly sample for each digit 3000 images from the database. Fig. 10 displays a small subset of the sample images of the two digits, as well as all 6000 sample images projected onto the top three PCA dimensions. We apply the algorithm to construct the geometric wavelets and show the wavelet coefficients and the reconstruction errors at all scales in Fig. 11. We select local dimensions for scaling functions by keeping 50% and 95% of the variance, respectively, at the nonleaf and leaf nodes. We observe that the magnitudes of the coefficients stops decaying after a certain scale. This indicates that the data is not on a smooth manifold. We expect optimization

¹Available at <http://yann.lecun.com/exdb/mnist/>.

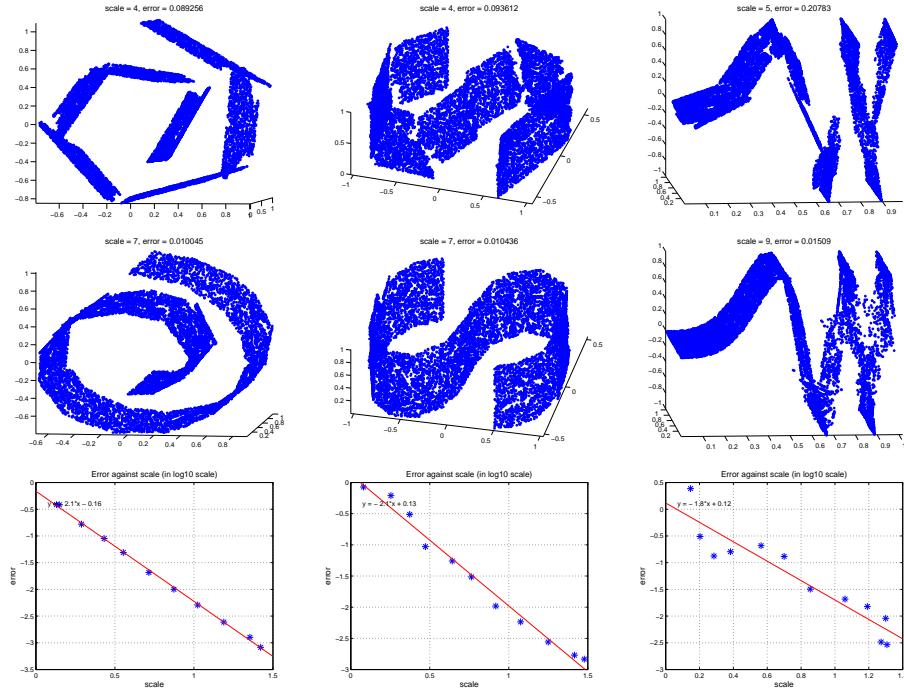


FIGURE 7. Top and Middle: Reconstructions by the algorithm of the three toy data sets in Fig. 5 at two selected scales. Bottom: Reconstruction errors as a function of scale.

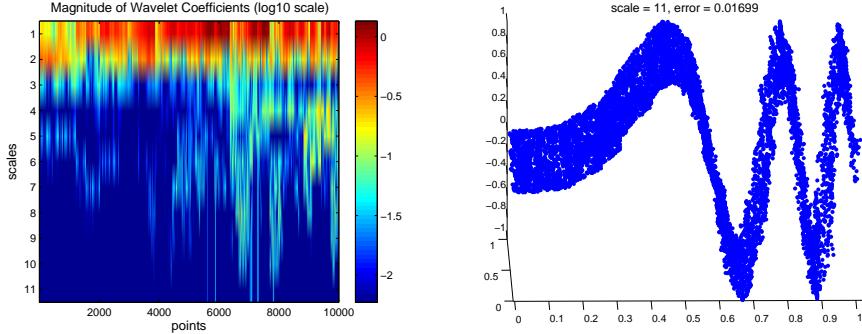


FIGURE 8. We threshold the wavelet coefficients of the *Oscillating2DWave* data at the level of .01 and prune the dyadic tree accordingly. The figure, from left to right, respectively shows the reduced matrix of wavelet coefficients (only their magnitudes), and the corresponding best approximation of the manifold.

of the tree and of the wavelet dimensions in future work to lead to a more efficient representation in this case.

We then fix a data point (or equivalently an image), for each digit, and show in Fig. 12 its reconstructed coordinates at all scales and the corresponding dictionary

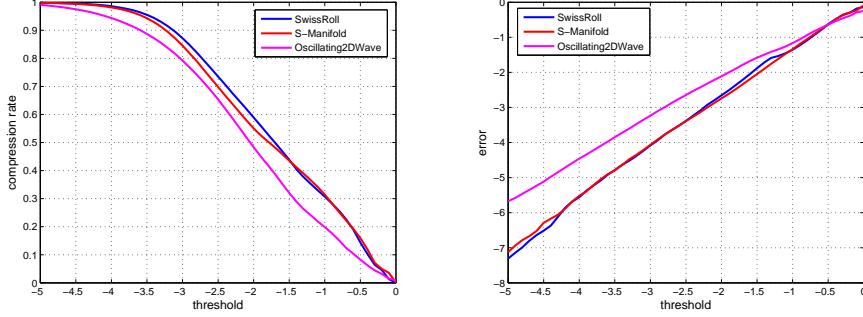


FIGURE 9. Left: the compression ratio of the matrix of the wavelet coefficients shown in Fig. 6. Right: the corresponding approximation errors. The linearity is consistent with Theorem 2.3, and essentially says that thresholding level δ generates approximation errors of order at most $O(\delta)$.

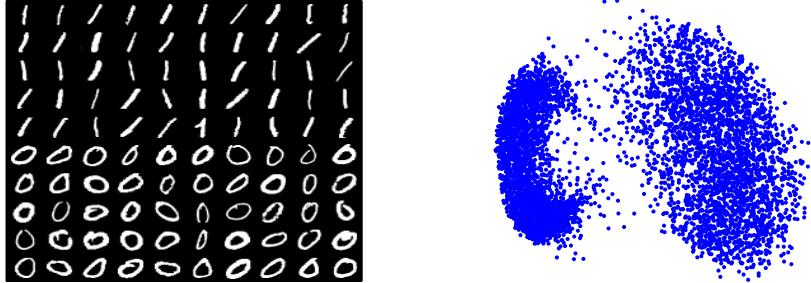


FIGURE 10. Some examples of the MNIST digits 1 and 0 (left) and 6000 sample images shown in top three PCA dimensions (right)

elements (all of which are also images). We see that at every scale we have a handwritten digit, which is an approximation to the fixed image, and those digits are refined successively to approximate the original data point. The elements of the dictionary quickly fix the orientation and the thickness, and then they add other distinguishing features of the image being approximated.

4.2.2. Human Face Images. We consider the cropped face images in both the Yale Face Database B² and the Extended Yale Face Database B³, which are available for 38 human subjects each seen in frontal pose and under 64 illumination conditions. (Note that the original images have large background variations, sometimes even for one fixed human subject, so we decide not to use them and solely focus on the faces.) Among these 2432 images, 18 of them are corrupted, which we discard. Fig. 13 displays a random subset of the 2414 face images. Since the images have large size (192×168), to reduce computational complexity we first project the

²<http://cvc.yale.edu/projects/yalefacesB/yalefacesB.html>

³<http://vision.ucsd.edu/~leekc/ExtYaleDatabase/ExtYaleB.html>

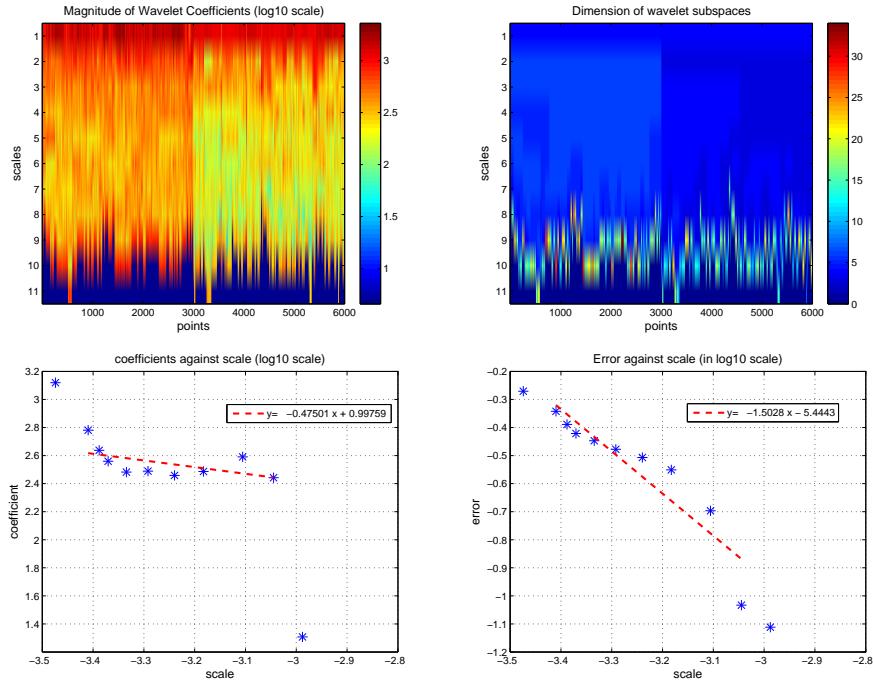


FIGURE 11. Top left: geometric wavelet representation of the MNIST digits 1 and 0. As usual, the vertical axis multi-indexes the wavelet coefficients, from coarse (top) to fine (bottom) scales: the block of entries at $(x, j), x \in C_{j,k}$ is $\log_{10} |q_{j,x}|$, where $q_{j,x}$ is the vector of geometric wavelet coefficients of x at scale j (see Sec. 3). In particular, each row indexes multiple wavelet elements, one for each $k \in \mathcal{K}_j$. Top right: dimensions of the wavelet subspaces (with the same convention as in the previous plot). Bottom: magnitude of coefficients (left) and reconstruction error (right) as functions of scale. The red lines are fitted omitting the first and last points (in each plot) in order to more closely approximate the linear part of the curve.

images into the first 500 dimensions by SVD, keeping about 99.5% variance. We apply the algorithm to the compressed data to construct the geometric wavelets and show the wavelet coefficients, dimensions and reconstruction errors at all scales in Fig. 14. Again, we have kept 50% and 95% of the variance, respectively, at the nonleaf and leaf nodes when constructing scaling functions. Note that both the magnitudes of the wavelet coefficients and the approximation errors have similar patterns with those for the MNIST digits (see Fig. 11), indicating again a lack of manifold structure in this data set. We also fix an image and show in Fig. 15 its reconstructed coordinates at all scales and the corresponding wavelet bases (all of which are also images).

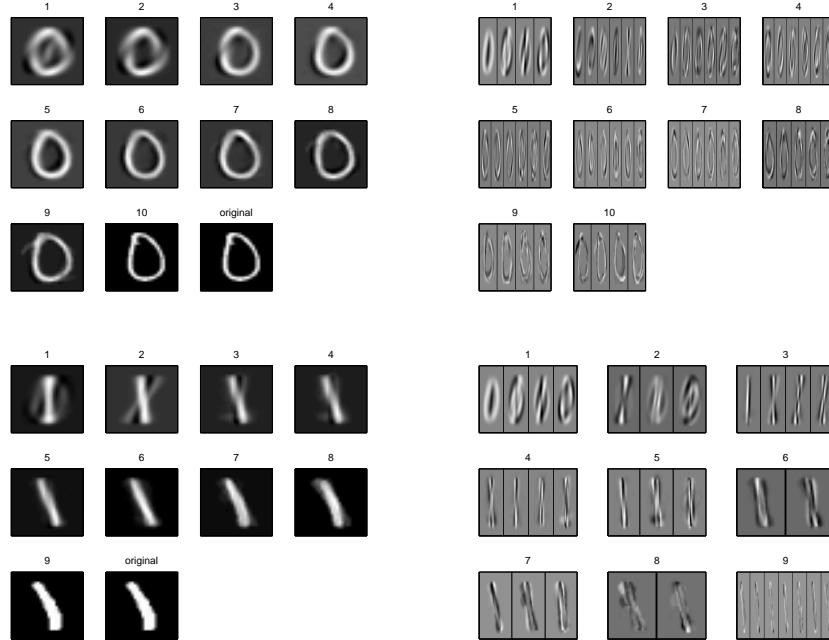


FIGURE 12. Left column: in each figure we plot coarse-to-fine geometric wavelet approximations of the original data point (represented in the last image). Right column: elements of the wavelet dictionary (ordered from coarsest to finest scales) used in the expansion on the left.

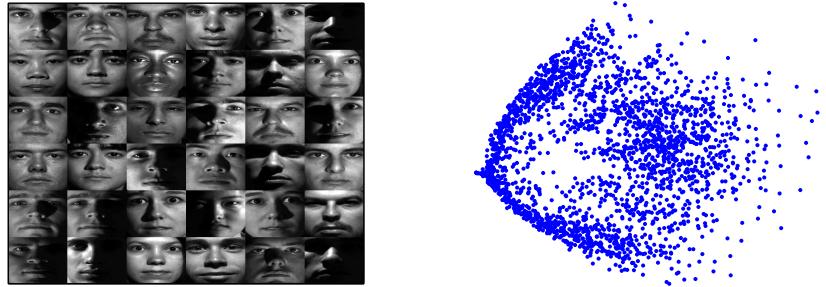


FIGURE 13. Left: A random subset of the 2414 face images (38 human subjects in frontal pose under 64 illumination conditions); Right: the entire data set shown in top three PCA dimensions.

5. ORTHOGONAL GEOMETRIC MULTI-RESOLUTION ANALYSIS

Neither the vectors $Q_{\mathcal{M}_{j+1}}(x)$, nor any of the terms that comprise them, are in general orthogonal across scales. On the one hand, this is natural since \mathcal{M} is

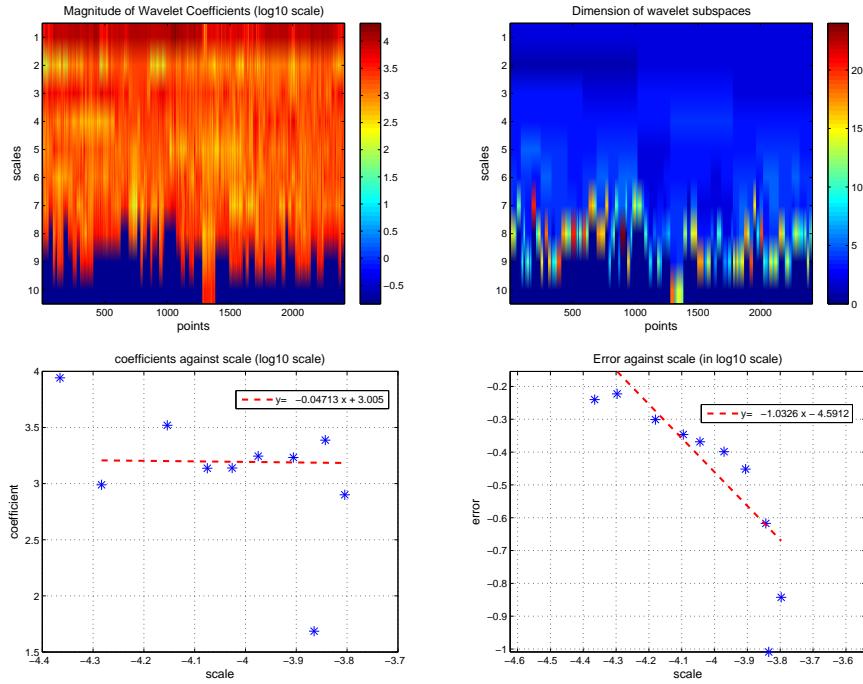


FIGURE 14. Top left: magnitudes of the wavelet coefficients of the cropped faces (2414 images) arranged in a tree. Top right: dimensions of the wavelet subspaces. Bottom: magnitude of coefficients (left) and reconstruction error (right) as functions of scale. The red lines are fitted omitting the first and last points (in each plot) in order to more closely approximate the linear part of the curve.

nonlinear, and the lack of orthogonality here is a consequence of that. On the other hand, the $Q_{\mathcal{M}_{j+1}}(x)$ may be almost parallel across scales or, for example, the subspaces $W_{j+1,x}$ may share directions across scales. If that was the case, we could more efficiently encode the dictionary by not encoding shared directions twice. A different construction of geometric wavelets achieves this. We describe this modification with a coarse-to-fine algorithm, which seems most natural. We start at scales 0 and 1, letting

$$(5.1) \quad S_{0,x} = V_{0,x} \quad , \quad S_{1,x} = S_{0,x} \oplus W_{1,x} \quad , \quad U_{1,x} = W_{1,x},$$

and for $j \geq 1$,

$$(5.2) \quad U_{j+1,x} = P_{S_{j,x}^\perp}(W_{j+1,x}) \quad , \quad S_{j+1,x} = S_{j,x} \oplus U_{j+1,x}$$

Observe that the sequence of subspaces $S_{j,x}$ is increasing: $S_{0,x} \subseteq S_{1,x} \subseteq \dots \subseteq S_{j,x} \subseteq \dots$ and the subspace $U_{j+1,x}$ is exactly the orthogonal complement of $S_{j,x}$ into $S_{j+1,x}$. This is a situation analogous to that of classical wavelet theory. Also, we may write

$$(5.3) \quad W_{j+1,x} = U_{j+1,x} \oplus P_{S_{j,x}}(W_{j+1,x})$$

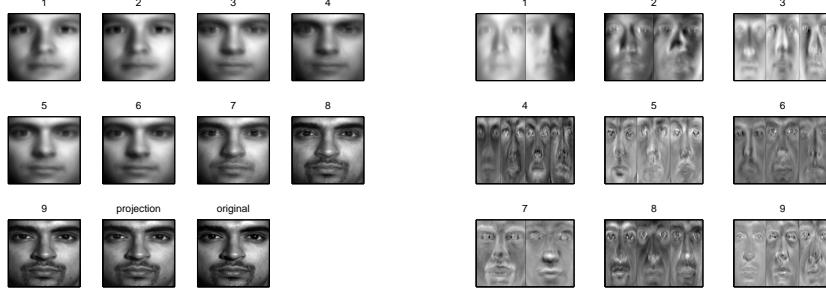


FIGURE 15. Left: in images 1-9 we plot coarse-to-fine geometric wavelet approximations of the projection and the original data point (represented in the last two images). Right: elements of the wavelet dictionary (ordered from coarse to fine in 1-9) used in the expansion on the left.

where the direct sum is orthogonal. At each scale j we do not need to construct a new wavelet basis for each $W_{j+1,x}$, but we only need to construct a new basis for $U_{j+1,x}$, and express $Q_{j+1,x}(x)$ in terms of this new basis, and the wavelet and scaling function bases constructed at the previous scales. This reduces the cost of encoding the wavelet dictionary as soon as $\dim(U_{j+1,x}) < \dim(W_{j+1,x})$ which, as we shall see, may occur in both artificial and real world examples. From a geometrical perspective, this roughly corresponds to the normal space to \mathcal{M} at a point not varying much at fine scales.

Finally, we note that we can define new projections of a point x into these subspaces $S_{j,x}$:

$$(5.4) \quad s_{j,x} = P_{S_{j,x}}(x - c_{j,x}) + c_{j,x}.$$

Note that since $V_{j,x} \subseteq S_{j,x}$, $s_{j,x}$ is a better approximation than x_j to x at scale j (in the least squares sense). Also,

$$(5.5) \quad s_{j+1,x} - s_{j,x} = U_{j+1,x}U_{j+1,x}^*(x - c_{j+1,x}) + (I - P_{S_{j,x}})(c_{j+1,x} - c_{j,x}).$$

We display in Figs. 16, 17, 18 pseudo-codes for the orthogonal GMRA and the corresponding forward and inverse transforms. The reader may want to compare with the corresponding routines for the regular GMRA construction, displayed in Figs. 2, 3, 4. Note that as the name suggests, the wavelet bases $\Psi_{j,k}$ along any path down the tree are mutually orthogonal. Moreover, the local scaling function at each node of such a path is effectively the union of the wavelet bases of the node itself and its ancestors. Therefore, the Orthogonal GMRA tree will have small height if the data set has a globally low dimensional structure, i.e., there is small number of normal directions in which the manifold curves.

Example: A connection to Fourier analysis. Suppose we consider the classical space of band-limited functions of band B :

$$(5.6) \quad BF_B = \{f : \text{supp. } \hat{f} \subseteq [-B\pi, B\pi]\}.$$

It is well-known that classical classes of smooth functions (e.g. $W^{k,2}$) are characterized by their L^2 -energy in dyadic spectral bands of the form $[-2^{j+1}\pi, -2^j\pi] \cup$

```

OrthoGMRA = OrthogonalGMRA ( $X_n, \tau_0, \epsilon$ )
// Input:
//  $X_n$ : a set of  $n$  samples from  $\mathcal{M}$ 
//  $\tau_0$ : some method for choosing local dimensions
//  $\epsilon$ : precision
// Output:
// A tree  $\mathcal{T}$  of dyadic cells  $\{C_{j,k}\}$  with their local means  $\{c_{j,k}\}$ , and a family of
orthogonal geometric wavelets  $\{U_{j,k}\}$ , and corresponding translations  $\{w_{j,k}\}$ 

Construct the cells  $C_{j,k}$ , and form a dyadic tree  $\mathcal{T}$  with local centers  $c_{j,k}$ .
Let  $\text{cov}_{0,k} = |C_{0,k}|^{-1} \sum_{x \in C_{0,k}} (x - c_{0,k})(x - c_{0,k})^*$ , for  $k \in \mathcal{K}_0$ , and compute
SVD( $\text{cov}_{0,k}$ ) =  $\Phi_{0,k} \Sigma_{0,k} \Phi_{0,k}^*$  (where the dimension of  $\Phi_{0,k}$  is determined by  $\tau_0$ ).
Set  $j = 0$  and  $\Psi_{0,k} := \Phi_{0,k}$ ,  $w_{0,k} := c_{0,k}$ 
Let  $J$  be the maximum scale of the tree
while  $j < J$ 
    for  $k \in \mathcal{K}_j$ 
        Let  $\Phi_{j,k}^{(cum)} = [\Psi_{\ell,k''}]_{0 \leq \ell \leq j}$  be the union of all wavelet bases of the cell  $C_{j,k}$ 
        and its ancestors. If the subspace spanned by  $\Phi_{j,k}^{(cum)}$  can approximate the
        cell within the given precision  $\epsilon$ , then remove all the offspring of  $C_{j,k}$  from
        the tree. Otherwise, do the following.
            Compute  $\text{cov}_{j+1,k'}$  and  $\Phi_{j+1,k'}$ , for all  $k' \in \text{children}(j, k)$ , as above
            For each  $k' \in \text{children}(j, k)$ , construct the wavelet bases  $U_{j+1,k'}$  as the
            complement of  $\Phi_{j+1,k'}$  in  $\Phi_{j,k}^{(cum)}$ . The translation  $w_{j+1,k'}$  is the projec-
            tion of  $c_{j+1,k'} - c_{j,k}$  into the space orthogonal to that spanned by the
             $\Phi_{j,k}^{(cum)}$ .
    end
     $j = j + 1$ 
end

```

FIGURE 16. Pseudo-code for the construction of an Orthogonal Geometric Multi-Resolution Analysis.

```

 $\{q_{j,x}\} = \text{orthoFGWT}(\text{orthoGMRA}, x)$ 
// Input: orthoGMRA structure,  $x \in \mathcal{M}$ 
// Output: A sequence  $\{q_{j,x}\}$  of wavelet coefficients

 $r = x$ 
for  $j = J$  down to 0
     $q_{j,x} = U_{j,x}^*(r - c_{j,x})$ 
     $r = r - (U_{j,x} \cdot q_{j,x} + w_{j,x})$ 
end

```

FIGURE 17. Pseudo-code for the Forward Orthogonal Geometric Wavelet Transform

$[2^j\pi, 2^{j+1}\pi]$, i.e. by the L^2 -size of their projection onto $BF_{2^{j+1}} \ominus BF_{2^j}$ (some care is in fact needed in smoothing these frequency cutoffs, but this issue is not relevant for our purposes here). If we observe samples from such smoothness spaces, which kind of dictionary would result from our GMRA construction? We consider

```

 $\hat{x} = \text{orthoIGWT}(\text{orthoGMRA}, \{q_{j,x}\})$ 
// Input: orthoGMRA structure, wavelet coefficients  $\{q_{j,x}\}$ 
// Output: Approximation  $\hat{x}$  at scale  $J$ 

 $\hat{x} = 0$ 
for  $j = 0$  to  $J$ 
     $\hat{x} = \hat{x} + U_{j,x} q_{j,x} + w_{j,x}$ 
end

```

FIGURE 18. Pseudo-code for the Inverse Orthogonal Geometric Wavelet Transform

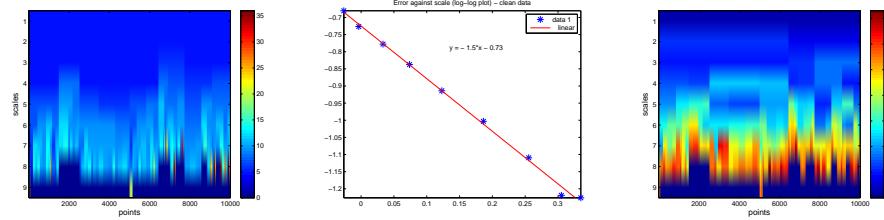


FIGURE 19. We construct an Orthogonal Geometric Multi-Resolution Analysis (see Sec. 5) on a random sample of 10000 band-limited functions. Left: dimension of the GMRA wavelet subspaces. Center: approximation error as a function of scale. Right: dominant frequency in each GMRA subspace, showing that frequencies are sorted from low (top, coarse GMRA scales) to high (bottom, fine GMRA scales). This implies that the geometric scaling function subspaces roughly corresponds to a Littlewood-Paley decomposition, and the GWT of a function f corresponds to a rough standard wavelet transform.

the following example: we generate random smooth (band-limited!) functions as follows:

$$(5.7) \quad f_\omega(x) = \sum_{j=0}^J a_j(\omega) \cos(jx)$$

with a_j random Gaussian (or bounded) with mean $2^{-\lfloor \frac{j}{2} \rfloor \alpha}$ and standard deviation $2^{-\lfloor \frac{j}{2} \rfloor \alpha} \cdot \frac{1}{5}$. These functions are smooth and have comparable norms in a wide variety of smoothness spaces, e.g. $W^{2,2}$, so that they may be thought of as approximately random samples from the unit ball in such space, intersected with band-limited functions. We construct a GMRA on a random sample from this family of functions and see that it organizes this family of functions in a Littlewood-Paley type of decomposition: the scaling function subspace at scale j roughly corresponds to $BF_{2^{j+1}} \ominus BF_{2^j}$, and the GMRA of a point is essentially a block Fourier transform, where coefficients in the same dyadic band are grouped together. This is as expected since the geometry of this data set is that of an ellipsoid with axes of equal length in each dyadic frequency band, and decreasing length as j increases. It follows that the coefficients in the FGWT of a function f measure the energy of f in dyadic

bands in frequency, and is therefore an approximate FFT of sorts. Finally, observe that the cost of the FGWT of a point f is comparable to the cost of the Fast Fourier Transform.

6. VARIATIONS, GREEDY ALGORITHMS, AND OPTIMIZATIONS

We discuss several techniques for reducing the encoding cost of the geometric wavelet dictionary and/or speeding up the decay of the geometric wavelet coefficients.

6.1. Splitting of the wavelet subspaces. Fix a cell $C_{j,k}$. For any $k' \in \text{children}(j, k)$, we may reduce the cost of encoding the subspace $W_{j+1,k'}$ by splitting it into a part that depends only on (j, k) and another on $(j + 1, k')$:

$$(6.1) \quad W_{j,k}^\cap := \cap_{k' \in \text{children}(j,k)} W_{j+1,k'}$$

and $W_{j+1,k'}^\perp$ be the orthogonal complement of $W_{j,k}^\cap$ in $W_{j+1,k'}$. We may choose orthonormal bases $\Psi_{j,k}^\cap$ and $\Psi_{j+1,k'}^\perp$ for $W_{j,k}^\cap$ and $W_{j+1,k'}^\perp$ respectively, and let $Q_{j,k}^\cap$, $Q_{j+1,k'}^\perp$ be the associated orthogonal projections. For the data in $C_{j+1,k'}$, we have therefore constructed the geometric wavelet basis

$$(6.2) \quad \Psi_{j+1,k'} = [\Psi_{j,k}^\cap | \Psi_{j+1,k'}^\perp],$$

together with orthogonal splitting of the projector

$$(6.3) \quad Q_{j+1,k'} = Q_{j,k}^\cap + Q_{j+1,k'}^\perp,$$

where the first term in the right-hand side only depends on the parent (j, k) , and the children-dependent information necessary to go from coarse to fine is encoded in the second term. This is particularly useful when $\dim(W_{j,k}^\cap)$ is large relative to $\dim(W_{j+1,k'})$.

6.2. A fine-to-coarse strategy with no tangential corrections. In this variation, instead of the sequence of approximations $x_j = \mathbb{P}_{V_{j,x}}(x)$ to a point $x \in \mathcal{M}$, we will use the sequence $\tilde{x}_j = \mathbb{P}_{V_{j,x}}(\tilde{x}_{j+1})$, for $j < J$, and $\tilde{x}_J := x_J$. The collection $\widetilde{\mathcal{M}}_j$ of \tilde{x}_j for all $x \in \mathcal{M}$ is a coarse approximation to the manifold \mathcal{M} at scale j . This roughly corresponds to considering only the first term in (2.17), disregarding the tangential corrections. The advantage of this strategy is that the tangent planes and the corresponding dictionary of geometric scaling functions do not need to be encoded. The disadvantage is that the point \tilde{x}_j does not have the same clear-cut interpretation as x_j has, as it is not anymore the orthogonal projection of x onto the best (in the least square sense) plane approximating $C_{j,x}$. Moreover, \tilde{x}_j really depends on J : if one starts the transform at a different finest scale, the sequence changes. Notwithstanding this, if we choose J so that $\|x_J - x\| < \epsilon$, for some precision $\epsilon > 0$, then this sequence does provide an efficient multi-scale encoding of x_J (and thus of x up to precision ϵ).

The claims above become clear as we derive the equations for the transform:

$$(6.4) \quad \begin{aligned} Q_{\widetilde{\mathcal{M}}_j}(\tilde{x}_{j+1}) &:= \tilde{x}_{j+1} - \tilde{x}_j = \tilde{x}_{j+1} - P_{j,x}(\tilde{x}_{j+1} - c_{j,x}) - c_{j,x} \\ &= (I - \Phi_{j,x}\Phi_{j,x}^*)((\tilde{x}_{j+1} - c_{j+1,x}) + (c_{j+1,x} - c_{j,x})). \end{aligned}$$

Noting that $\tilde{x}_{j+1} - c_{j+1,x} \in \langle \Phi_{j+1,x} \rangle$, we obtain

$$(6.5) \quad Q_{\widetilde{\mathcal{M}}_j}(\tilde{x}_{j+1}) = \Psi_{j+1,x}\Psi_{j+1,x}^*(\tilde{x}_{j+1} - c_{j+1,x}) + w_{j+1,x},$$

where $\Psi_{j+1,x}, w_{j+1,x}$ are the same as in (2.18). By definition we still have the multi-scale equation

$$(6.6) \quad \tilde{x}_{j+1} = \tilde{x}_j + Q_{\widetilde{\mathcal{M}}_{j+1}}(\tilde{x}_{j+1})$$

for $\{\tilde{x}_j\}$ defined as above.

6.3. Out-of-sample extension. In many applications it will be important to extend the geometric wavelet expansion to points that were not sampled, and/or to points that do not lie exactly on \mathcal{M} . For example, \mathcal{M} may be composed of data points satisfying a model, but noise or outliers in the data may not lie on \mathcal{M} .

Fix $x \in \mathbb{R}^D$, and let J be the finest scale in the tree. Let $c_{J,x}$ be a closest point to x in the net $\{c_{J,k}\}_{k \in \mathcal{K}_J}$; such a point is unique if x is close enough to \mathcal{M} . For $j \leq J$, we will let (j, x) be the index of the (unique) cell at scale j that contains $c_{J,x}$. With this definition, we may calculate a geometric wavelet expansion of the point $P_{J,x}(x)$. However, $e_J(x) := x - P_{J,x}(x)$ is large if x is far from \mathcal{M} . We may encode this difference by greedily projecting it onto the family of linear subspaces $W_{J,x}, \dots, W_{1,x}$ and $V_{0,x}$, i.e. by computing

$$\begin{aligned} Q_{\mathcal{M}^\perp,J}(x) &:= Q_{J,x}(e_J(x)), \\ Q_{\mathcal{M}^\perp,J-1}(x) &:= Q_{J-1,x}(e_J(x) - Q_{\mathcal{M}^\perp,J}(x)), \\ &\dots \\ (6.7) \quad Q_{\mathcal{M}^\perp,0}(x) &:= P_{0,x}(e_J(x) - Q_{\mathcal{M}^\perp,J}(x) - \dots - Q_{\mathcal{M}^\perp,1}(x)). \end{aligned}$$

These projections encode, greedily along the multi-scale “normal” subspaces $\{Q_{j,x}\}$.

The computational complexity of this operation is comparable to that of computing two sets of wavelet coefficients, plus that of computing the nearest neighbor of x among the centers $\{c_{J,k}\}_{k \in \mathcal{K}_J}$ at the finest scale. By precomputing a tree for fast nearest neighbor computations, this essentially requires $O(\log(|\mathcal{K}_J|))$ operations. Also, observe that $|\mathcal{K}_J|$ in general does not depend on the number of points n , but on the precision in the approximation specified in the tree construction.

6.4. Spin-cycling: multiple random partitions and trees. Instead of one multi-scale partition and one associated tree, in various situations it may be advantageous to construct multiple multi-scale partitions and corresponding trees. This is because a single partition introduces somewhat arbitrary cuts and possible related artifacts in the approximation of \mathcal{M} , and in the construction of the geometric wavelets in general. Generating multiple partitions or families of approximations is a common technique in signal processing. For example, in [67] it is shown that denoising by averaging the result of thresholding on multiple shifted copies of the Haar system is as optimal (in a suitable asymptotic, minimax sense) as performing the same algorithm on a single system of smoother wavelets (and in that paper the technique was called *spin-cycling*). In the study of approximation of metric spaces by trees [68], it is well understood that using a suitable weighted average of metrics of suitably constructed trees is much more powerful than using a single tree (this may be seen already when trying to find tree metrics approximating the Euclidean metric on an interval).

In our context, it is very natural to consider a family of trees and the associated geometric wavelets, and then perform operations on either the union of such geometric wavelet systems (which would be a generalization of sorts of tight frames, in a geometric context), or perform operations on each system independently and

then average. In particular, the construction of trees via cover trees [57] is very easily randomized, while still guaranteeing that each instance of such trees is well-balanced and well-suited for our purposes. We leave a detailed investigation to a future publication.

7. DATA REPRESENTATION AND COMPRESSION

A generic point cloud with n points in \mathbb{R}^D can trivially be stored in space Dn . If the point cloud lies, up to, say, a least-squares error (relative or absolute) ϵ in a linear subspace of dimension $d_\epsilon \ll D$, we could encode n points in space

$$(7.1) \quad \underbrace{Dd_\epsilon}_{\substack{\text{cost of} \\ \text{encoding basis}}} + \underbrace{nd_\epsilon}_{\substack{\text{cost of encoding} \\ n \text{ points}}} = d_\epsilon(D + n),$$

which is clearly much less than nD . In particular, if the d -dimensional point cloud lies in a d -dimensional subspace, then $d_\epsilon = d$ and

$$(7.2) \quad d(D + n).$$

Let us compute the cost of encoding with a geometric multi-resolution analysis a manifold \mathcal{M} of dimension d sampled at n points, and fix a precision $\epsilon > 0$. We are interested in the case $n \rightarrow +\infty$. The representation we use is, as in (2.20):

$$(7.3) \quad x \sim x_J = P_{\mathcal{M}_0}(x) + \sum_{j=1}^J Q_{\mathcal{M}_j}(x),$$

where we choose the smallest J such that $\|x - x_J\| < \epsilon$. In the case of a \mathcal{C}^2 manifold, $J = \log_2 \epsilon^{-\frac{1}{2}}$ because of Theorem 2.3. However, d_ϵ as defined above with global SVD may be as large as D in this context, even for $d = 1$.

Since \mathcal{M} is nonlinear, we expect the cost of encoding a point cloud sampled from \mathcal{M} to be larger than the cost (7.2) of encoding a d -dimensional flat \mathcal{M} ; however the geometric wavelet encoding is not much more expensive, having a cost:

$$(7.4) \quad \underbrace{dD + 2\epsilon^{-\frac{d}{2}}(d^\perp + 2^{-d}d^\cap + 2)}_{\substack{\text{cost of} \\ \text{encoding basis}}} D + \underbrace{nd(1 + \log_2 \epsilon^{-\frac{1}{2}})}_{\substack{\text{cost of encoding} \\ n \text{ points}}}$$

In Sec. 7.2 we compare this cost with that in 7.1 on several data sets. To see that the cost of the geometric wavelet encoding is as promised, we start by counting the geometric wavelet coefficients used in the multi-scale representation. Recall that $d_{j,x}^w = \text{rank}(\Psi_{j,x})$ is the number of wavelet coefficients at scale j for the given point x . Clearly, $d_{j,k}^w \leq d$. Then, the geometric wavelet transform of all points takes space at most

$$(7.5) \quad nd + \sum_{j=1}^J \sum_x d_{j,x}^w \leq nd + ndJ \leq nd(1 + \log_2 \epsilon^{-\frac{1}{2}}),$$

independently of D . The dependency on n, d is near optimal, and this shows that data points have a sparse, or rather, compressible, representation in terms of geometric wavelets. Next we compute the cost of the geometric wavelet dictionary, which contains the geometric wavelet bases $\Psi_{j,k}$, translations $w_{j,k}$, and cell centers $c_{j,k}$. If we add the tangential correction term as in (2.18), then we should also

include the geometric scaling functions $\Phi_{j,k}$ in the cost. Let us assume for now that we do not need the geometric scaling functions. Define

$$(7.6) \quad d_{j,k}^\cap := \text{rank}(\Psi_{j,k}^\cap),$$

$$(7.7) \quad d_{j+1,k'}^\perp := \text{rank}(\Psi_{j+1,k'}^\perp)$$

and assume that $d_{j,k}^\cap \leq d^\cap$, $d_{j+1,k'}^\perp \leq d^\perp$ for fixed constants $d^\cap, d^\perp \leq d$. The cost of encoding the wavelet bases $\{\Psi_{j,k}\}_{k \in \mathcal{K}_j, 0 \leq j \leq J}$ is at most

$$\begin{aligned} & \underbrace{dD}_{\text{cost of } \Psi_{0,k}} + \sum_{j=0}^{J-1} \underbrace{2^{dj}}_{\#\text{ cells at scale } j} \underbrace{d^\cap D}_{\text{cost of } \Psi_{j,k}^\cap} + \underbrace{2^{d(j+1)}}_{\#\text{ cells at scale } j+1} \underbrace{d^\perp D}_{\text{cost of } \Psi_{j+1,k'}^\perp} \\ (7.8) \quad & = dD + \frac{2^{dJ} - 1}{2^d - 1} (d^\cap D + 2^d d^\perp D) \leq dD + 2\epsilon^{-\frac{d}{2}} (d^\perp + 2^{-d} d^\cap) D. \end{aligned}$$

The cost of encoding $w_{j,k}, c_{j,k}$ is

$$(7.9) \quad 2 \sum_{j=0}^J 2^{dj} D \leq 2 \cdot 2^{dJ+1} \cdot D = 4D\epsilon^{-\frac{d}{2}}.$$

Therefore, the overall cost of the dictionary is

$$(7.10) \quad dD + 2\epsilon^{-\frac{d}{2}} (d^\perp + 2^{-d} d^\cap + 2) D.$$

In the case that we also need to encode the geometric scaling functions $\Phi_{j,k}$, we need an extra cost of

$$(7.11) \quad \sum_{j=0}^J 2^{dj} dD \leq 2\epsilon^{-\frac{d}{2}} dD.$$

7.1. Pruning of the geometric wavelets tree. In this section we discuss how to prune the geometric wavelets tree with the goal of minimizing the total cost for ϵ -encoding a given data set, i.e., encoding the data within the given precision $\epsilon > 0$. Since we are not interested in the intermediate approximations, we will adopt the GMRA version without adding the tangential corrections (see Sec. 6.2) and thus there is no need to encode the scaling functions. The encoding cost includes both the cost of the dictionary, defined for simplicity as the number of dictionary elements $\{\Psi_{j,k}, w_{j,k}, c_{j,k}\}$ multiplied by the ambient dimension D , and the cost of the coefficients, defined for simplicity to be the number of nonzero coefficients required to reconstruct the data up to precision ϵ .

7.1.1. Discussion. We fix an arbitrary nonleaf node $C_{j,k}$ of the partition tree \mathcal{T} and discuss how to ϵ -encode the local data in $C_{j,k}$ in order to achieve minimal encoding cost. We assume that the data in the children nodes $C_{j+1,k'}, k' \in \text{children}(j, k)$, has been optimally ϵ -encoded by some methods, with scaling functions $\Phi_{j+1,k'}$ of dimensions $d_{j+1,k'}$ and corresponding encoding costs $\varphi_{j+1,k'}$. For example, when $C_{j+1,k'}$ is a leaf node, it can be optimally ϵ -encoded by using a local PCA plane of minimal dimension $d_{j+1,k'}^\epsilon$, with the corresponding encoding cost

$$(7.12) \quad \varphi_{j+1,k'} = n_{j+1,k'} \cdot d_{j+1,k'}^\epsilon + D \cdot d_{j+1,k'}^\epsilon + D,$$

where $n_{j+1,k'}$ is the size of this node.

We consider the following ways of ϵ -encoding the data in $C_{j,k}$:

- (I) using the existing methods for the children $C_{j+1,k'}$ to encode the data in $C_{j,k}$ separately;
- (II) using only the parent node and approximating the local data by a PCA plane of minimal dimension $d_{j,k}^\epsilon$ (with basis $\Phi_{j,k}^\epsilon$);
- (III) using a multi-scale structure to encode the data in the node $C_{j,k}$, with the top $d_{j,k}^w$ PCA directions $\Phi_{j,k}^w$ being the scaling function at the parent node and $d_{j+1,k'}^w$ dimensional wavelets encoding differences between $\Phi_{j+1,k'}$ and $\Phi_{j,k}^w$. Here, $0 \leq d_{j,k}^w \leq d_{j,k}^\epsilon$.

We refer to the above methods as *children-only* encoding, *parent-only* encoding and *wavelet* encoding, respectively. We make the following comments. First, method (I) leads to the sparsest coefficients for each point, while method (II) produces the smallest dictionary. Second, in method (III), it is possible to use other combinations of the PCA directions as the scaling function for the parent, but we will not consider those in this paper. Lastly, the children-only and parent-only encoding methods can be thought of corresponding to special cases of the wavelet encoding method, i.e., when $d_{j,k}^w = 0$ and $d_{j,k}^w = d_{j,k}^\epsilon$, respectively.

We compare the encoding costs of the three methods above. Suppose there are $n_{j,k}$ points in the node $C_{j,k}$ and $n_{j+1,k'}$ points in each $C_{j+1,k'}$, so that $n_{j,k} = \sum_{k'} n_{j+1,k'}$. When we encode the data in $C_{j,k}$ with a $d_{j,k}^\epsilon$ dimensional plane, we need space

$$(7.13) \quad n_{j,k} \cdot d_{j,k}^\epsilon + D \cdot d_{j,k}^\epsilon + D.$$

If we use the children nodes to encode the data in $C_{j,k}$, the cost is

$$(7.14) \quad \sum_{k'} \varphi_{j+1,k'}.$$

The encoding cost of the wavelet encoding method has a more complex formula, and is obtained as follows. Suppose that we put at the parent node a $d_{j,k}^w$ dimensional scaling function consisting of the top $d_{j,k}^w$ principal vectors, where $0 \leq d_{j,k}^w \leq d_{j,k}^\epsilon$, and that $\Psi_{j+1,k'}$ are the corresponding wavelet bases for the children nodes. Let $d_{j,k}^\cap \geq 0$ be the dimension of the intersection of the wavelet functions, and write $d_{j+1,k'}^w = d_{j,k}^\cap + d_{j+1,k'}^\perp$. Note that the intersection only needs to be stored once for all children. Then the overall encoding cost is

$$\begin{aligned}
\varphi_{j,k}^w &= \underbrace{\sum_{k'} \varphi_{j+1,k'} - d_{j+1,k'} (n_{j+1,k'} + D)}_{\text{children excluding the scaling functions and coefficients}} + \underbrace{n_{j,k} \cdot d_{j,k}^w + D \cdot d_{j,k}^w + D}_{\text{the parent}} \\
&\quad + \underbrace{n_{j,k} \cdot d_{j,k}^\cap + D \cdot d_{j,k}^\cap}_{\text{intersection of children wavelets}} + \underbrace{\sum_{k'} n_{j+1,k'} \cdot d_{j+1,k'}^\perp + D \cdot d_{j+1,k'}^\perp + D}_{\text{children-specific wavelets}} \\
&= \underbrace{\sum_{k'} \varphi_{j+1,k'} - (d_{j+1,k'} - d_{j+1,k'}^\perp) \cdot (n_{j+1,k'} + D)}_{\text{new cost for children}} + \underbrace{(n_{j,k} + D) \cdot (d_{j,k}^w + d_{j,k}^\cap)}_{\text{parent and children intersection}} \\
&\quad + \underbrace{D + \sum_{k'} D}_{\text{parent center and wavelet translations}}
\end{aligned} \tag{7.15}$$

Once the encoding costs in (7.13), (7.14) and (7.15) (for all $0 \leq d_{j,k}^w \leq d_{j,k}^\epsilon$) are all computed, we pick the method with the smallest cost for encoding the data in $C_{j,k}$, and also update $\Phi_{j,k}, \varphi_{j,k}$ correspondingly. We propose in the next section a pruning algorithm for practical realization of the above ideas.

7.1.2. A pruning algorithm. The algorithm requires as input a data set X_n and a precision parameter $\epsilon > 0$, and outputs a forest with orthonormal matrices $\{\Phi_{j,k}\}$ and $\{\Psi_{j,k}\}$ attached to the nodes and an associated cost function $\varphi_{j,k}$ defined on every node of the forest quantifying the cost of optimally ϵ -encoding the data in that node.

Our strategy is bottom-up. That is, we start at the leaf nodes and ϵ -encode them by using local PCA planes of minimal dimensions, and let $\{\Phi_{j,k}\}$ and $\{\varphi_{j,k}\}$ be their bases and corresponding encoding costs. We then proceed to their parents and determine the optimal way of encoding them using (7.13), (7.14) and (7.15). If the parent-only encoding achieves the minimal encoding cost, then we remove all the offspring of this node from the tree, including the children. If the children-only is the best, then we separate out the children subtrees from the tree and form new trees (we also remove the parent from the original tree and discard it). Note that these new trees are already optimized, thus we will not need to examine them again. If the wavelet encoding with some $\Phi_{j,k}^w$ (and corresponding wavelet bases $\Psi_{j+1,k'}$) does the best, then we update $\Phi_{j,k} := [\Phi_{j,k}^w \Phi_{j,k}^\cap]$ and $\varphi_{j,k}$ accordingly and let $\Phi_{j+1,k'}$ store the complement of $\Phi_{j,k}^\cap$ in $\Phi_{j+1,k'}$. We repeat the above steps for higher ancestors until we reach the root of the tree. We summarize these steps in Fig. 20 below.

7.2. Comparison with SVD. In this section we compare our algorithm with Singular Value Decomposition (SVD) in terms of encoding cost for various precisions. We may think of the SVD, being a global analysis, as providing a sort of Fourier geometric analysis of the data, to be contrasted with our GMRA, a multi-scale wavelet analysis. We use the two real data sets above, together with a new data set, the Science News, which comprises about 1100 text documents, modeled as vectors in 1000 dimensions, whose i -th entry is the frequency of the i -th word in a dictionary (see [30] for detailed information about this data set). For GMRA, we now consider three different versions: (1) the regular GMRA, but with the optimization strategies discussed in Secs. 6.1 and 6.2 (2) the orthogonal GMRA (in Sec. 5) and (3) the pruning GMRA (in Sec. 7.1). For each version of the GMRA, we threshold the wavelet coefficients to study the rates of change of the approximation errors and encoding costs. We present three different costs: one for encoding the wavelet coefficients, one for the dictionary, and one for both (see Fig. 21).

We compare these curves with those of SVD, which is applied in two ways: first, we compute the SVD costs and errors using all possible PCA dimensions; second, we gradually threshold the full SVD coefficients and correspondingly compress the dictionary (i.e., discard those multiplying identically zero coefficients). The curves are superposed in the same plots (see the black curves in Fig. 21).

8. COMPUTATIONAL CONSIDERATIONS

The computational cost may be split as follows.

Construction of proximity graph: we find the k nearest neighbors of each of the n points. Using fast nearest neighbor codes (e.g. cover trees [57] and references

```

PrunGMRA = PruningGMRA ( $X_n, \epsilon$ ) x

// Input:
//  $X_n$ : a set of  $n$  samples from  $\mathcal{M}$ 
//  $\epsilon$ : precision

// Output:
// A forest  $\mathcal{F}$  of dyadic cells  $\{C_{j,k}\}$  with their local means  $\{c_{j,k}\}$  and PCA bases
//  $\{\Phi_{j,k}\}$ , and a family of geometric wavelets  $\{\Psi_{j,k}\}, \{w_{j,k}\}$ , as well as encoding costs
//  $\{\varphi_{j,k}\}$ , associated to the nodes

Construct the dyadic cells  $C_{j,k}$ , and form a tree  $\mathcal{T}$  with local centers  $c_{j,k}$ .
For every leaf node in the tree  $\mathcal{T}$ , compute the minimal dimension  $d_{j,k}^\epsilon$  and corresponding basis  $\Phi_{j,k}$  and encoding costs  $\varphi_{j,k}$  for achieving precision  $\epsilon$ 

for  $j = J - 1$  down to 1
    Find all the nonleaf nodes of the tree  $\mathcal{T}$  at scale  $j$ 
    For each of the nodes  $(j, k), k \in \mathcal{K}_j$ ,
        (1) Compute the encoding costs of the three methods, i.e., parent-only,
            children-only, and wavelet, using equations (7.13), (7.14) and (7.15).
        (2) Update  $\varphi_{j,k}$  with the minimum cost.
            if parent-only is the best,
            delete all the offspring of the node from  $\mathcal{T}$ , and let  $\Phi_{j,k} = \Phi_{j,k}^\epsilon$ 
            elseif children-only is the best,
            separate out the children subtrees from  $\mathcal{T}$  and form new trees, and also
            remove and discard the parent node
            else
            update  $\Phi_{j,k} := [\Phi_{j,k}^w \Phi_{j,k}^\cap]$  and  $\varphi_{j,k}$  accordingly and let  $\Phi_{j+1,k'}$  store the
            complement of  $\Phi_{j,k}^\cap$  in  $\Phi_{j+1,k'}$ .
            end
    end
end

```

FIGURE 20. Pseudo-code for the construction of the Pruning Geometric Wavelets

therein) the cost is $O_{d,D}(n \log n)$, with the constant being exponential in d , the intrinsic dimension of \mathcal{M} , and linear in D , the ambient dimension. The cost of computing the weights for the graph is $O(knD)$.

Graph partitioning: we use METIS [56] to create a dyadic partition, with cost $O(kn \log n)$. We may (albeit in practice we do not) compress the METIS tree into a 2^d -adic tree; however, this will not change the computational complexity below.

Computation of the $\Phi_{j,k}$'s: At scale j each cell $C_{j,k}$ of the partition has a number of points $n_{j,k} = O(2^{-jd}n)$, and there are $|\mathcal{K}_j| = O(2^{jd})$ such $C_{j,k}$'s. The cost of computing the rank- d SVD in each $C_{j,k}$ is $O(n_{j,k} D d)$, by using the algorithms of [69]. Summing over $j = 0, 1, \dots, J$ with $J \sim \log_{2^d} n$ we obtain a total cost $O(Dn \log n)$. At this point we have constructed all the $\Phi_{j,k}$'s. Observe that instead of $J \sim \log_{2^d} n$ we may stop at the coarsest scale at which a predetermined precision ϵ is reached (e.g. $J \sim \log_2 \frac{1}{\sqrt{\epsilon}}$ for a smooth manifold). In this case, the cost of this part of the algorithm only depends on ϵ and is independent of n . A similar but more complex strategy that we do not discuss here could be used also for the first two steps.

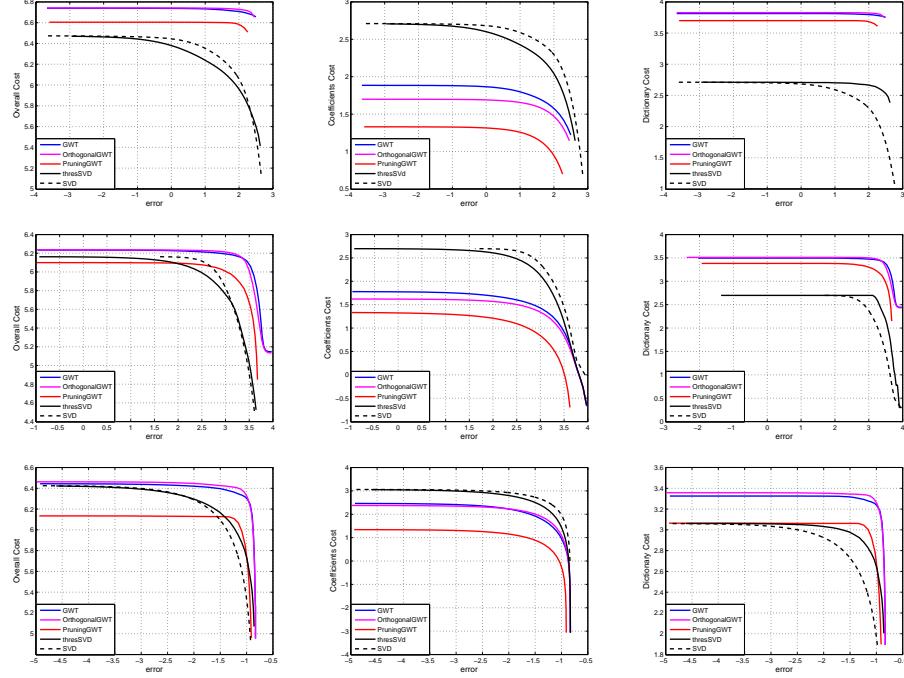


FIGURE 21. Cost-error curves for different kinds of encoding costs (left to right columns: overall, coefficients, dictionary) obtained on the three real data sets (top to bottom rows: MNIST digits, Yale Faces, and Science News) by the GMRA and SVD algorithms (represented by different curves in different colors). We see that all GMRA versions outperform SVD and its thresholding version in terms of coefficient costs (middle column), but take more space to store the dictionary (right column). This makes sense from the sparse coding perspective. Overall, the pruning GMRA algorithm does the best, while the other two GMRA versions have very close performance with both versions of SVD (see left column).

Computation of the $\Psi_{j,k}$'s: For each cell $C_{j,k}$, where $j < J$, the wavelet bases $\Psi_{j+1,k'}$, $k' \in \text{children}(j, k)$ are obtained by computing the partial SVD of a $d \times 2^d d$ matrix of rank at most d , which takes $O(D \cdot 2^d d \cdot d)$. Summing this up over all $j < J$, we get a total cost of $O(nDd^2)$.

Overall, the algorithm costs

$$(8.1) \quad O(nD(\log(n) + d^2)) + O_{d,D}(n \log n).$$

The cost of performing the FGWT of a point (or its inverse) is the sum of the costs of finding the closest leaf node, projecting onto the corresponding geometric scaling function plane, and then computing the multi-scale coefficients:

$$(8.2) \quad \underbrace{O_d(D \log n)}_{\text{cost of finding nearest } c_{J,k}} + \underbrace{\frac{dD}{\text{cost of projecting on } \Phi_{J,x}}}_{\text{cost of multi-scale transform}} + \underbrace{O(d^2 \log \epsilon^{-\frac{1}{2}})}_{\text{cost of multi-scale transform}},$$

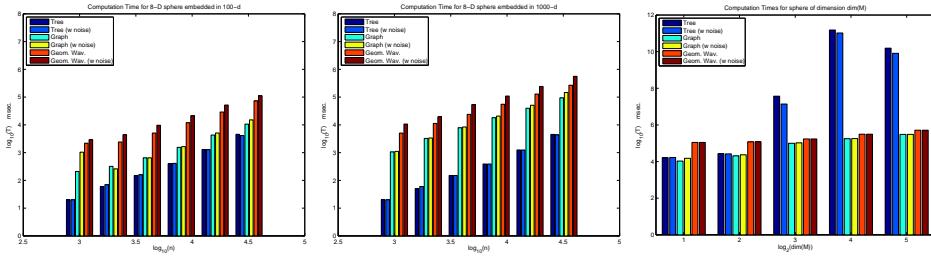


FIGURE 22. Timing experiments for the construction of geometric wavelets. We record separately the time to construct the nearest neighbor graph ('Graph'), the multi-scale partitions ('Tree'), and the geometric wavelets ('Geom. Wav.'). Left: time in *milliseconds* (on the vertical axis, in \log_{10} scale) vs. n (on the horizontal axis, also \log_{10} scale) for $\mathbb{S}^d(n, D, \sigma)$, for $n = 1000, 2000, 4000, 8000, 16000, 32000$, $d = 8$, $D = 100$, and $\sigma = 0, \frac{0.5}{\sqrt{D}}$. All the computational times grow linearly in n , with the noise increasing the computational time of each sub-computation. Center: same as left, but with $D = 1000$. A comparison with the experiment on the left shows that the increased ambient dimensionality does not cause, in this instance, almost any increase in the noiseless case, and in the noisy case the increase is a meager factor of 10, which is exactly the cost of handling vectors which are 10 times larger in distance computations, with no curse of ambient dimensionality. Right: computation times as a function of intrinsic dimension: we vary $d = 2, 4, 8, 16, 32$ (in \log_{10} scale on the horizontal axis)), and notice a mild increase in computation time, but with higher variances in the times for the computation of the multi-scale partitions.

with the O_d in the first term subsuming an exponential dependence on d . The cost of the IGWT is similar, but without the first term.

We report some results in practical performance in Fig. 22.

9. A NAÏVE ATTEMPT AT MODELING DISTRIBUTIONS

We present a simple example of how our techniques may be used to model measures supported on low-dimensional sets which are well-approximated by the multi-scale planes we constructed; results from more extensive investigations will be reported in an upcoming publication.

We sample n training points from a point cloud \mathcal{M} and, for a fixed scale j , we consider the coarse approximation \mathcal{M}_j (defined in (2.10)), and on each local linear approximating plane $V_{j,k}$ we use the training set to construct a multi-factor Gaussian model on $C_{j,k}$: let $\pi_{j,k}$ be the estimated distribution. We also estimate from the training data the probability $\pi_j(k)$ that a given point in \mathcal{M} belongs to $C_{j,k}$ (recall that j is fixed, so this is a probability distribution over the $|\mathcal{K}_j|$ labels of the planes at scale j). We may then generate new data points by drawing a $k \in \mathcal{K}_j$ according to π_j , and then drawing a point in $V_{j,k}$ from the distribution $\pi_{j,k}$: this defines a probability distribution supported on \mathcal{M}_j , that we denote by $p_{\mathcal{M}_j}$.

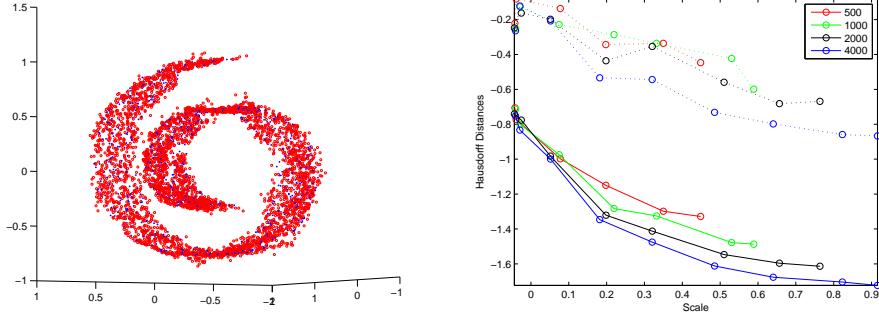


FIGURE 23. We generate a family of multi-scale models $\{p_i\}_{i=1}^4$, from 500, 1000, 2000, 4000 (corresponding to $i = 1, \dots, 4$) training samples from the swiss-roll manifold. Left: the blue points are 1000 training points, the red points are 4000 points generated according to p_2 at the finest scale $j = 6$. Right: for each $i = 1, \dots, 4$ and each scale j , we generate from p_i at scale j a point cloud of 4000 samples, and measure its Hausdorff distance (dotted lines) and “Hausdorff median distance” (continuous lines) from a randomly generated point cloud with 4000 points from the true distribution on the swiss roll. The x -axis is the scale j of the model used, and colors map the size of the training set. The construction of these models and the generation of the points clouds takes a few seconds on a standard desktop.

In this way we may generate new data points which are consistent with both the geometry of the approximating planes $V_{j,k}$ and with the distribution of the data on each such plane. In Fig. 23 we display the result of such modeling on a simple manifold. In Fig. 24 we construct $p_{\mathcal{M}_j}$ by training on 2000 handwritten 7's from the MNIST database, and on the same training set we train two other algorithms: the first one is based on projecting the data on the first a_j principal components, where a_j is chosen so that the cost of encoding the projection and the projected data is the same as the cost of encoding the GMRA up to scale j and the GMRA of the data, and then running the same multi-factor Gaussian model used above for generating $\pi_{j,k}$. This leads to a probability distribution we denote by $p_{SVD,j}$. Finally, we compare with the recently-introduced Multi-Factor Analyzer (MFA) Bayesian models from [39]. In order to test the quality of these models, we consider the following two measures. The first measure is simply the Hausdorff distance between 2000 randomly chosen samples according to each model and the training set: this is measuring how close the generated samples are to the training set. The second measure quantifies if the model captures the variability of the true data, and is computed by generating multiple point clouds of 2000 points for a fixed model, and looking at the pairwise Hausdorff distances between such point clouds, called the within-model Hausdorff distance variability.

The bias-variance tradeoff in the models $p_{\mathcal{M}_j}$ is the following: as j increases the planes better model the geometry of the data (under our usual assumptions), so that the bias of the model (and the approximation error) decreases as j increases;

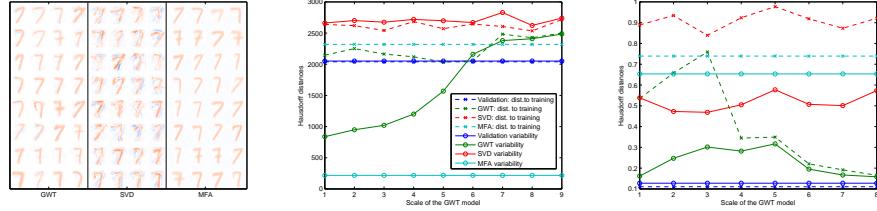


FIGURE 24. A training set of 2000 digits 7 from the MNIST data set are used to train probability models with GMRA ($p_{\mathcal{M}_j}$, one for each scale j in the GMRA of the training set), SVD (p_{SVD_j} , one for each GMRA scale, see text), and MFA p_{MFA} . Left: 32 digits drawn from $p_{\mathcal{M}_5}$, p_{SVD_5} and p_{MFA} : the quality of $p_{\mathcal{M}_5}$ and p_{MFA} is qualitatively better than that of p_{SVD_5} ; moreover $p_{\mathcal{M}_5}$ seem to capture more variability than p_{MFA} . Center: plots of the Hausdorff distance to training set and in-model Hausdorff distance variability. We see that both $p_{\mathcal{M}_j}$ and p_{MFA} have similar distance to the training set, while p_{SVD_j} , being a model in the ambient space, generates points farther from the distribution. Looking at the plots of the in-model Hausdorff distance variability, we see that such measure increases for $p_{\mathcal{M}_j}$ as a function of j (reflecting the increasing expression power of the model), while the same measure for p_{MFA} is very small, implying that MFA fails to capture the variability of the distribution, and simply generates an almost fixed set of points (in fact, local averages of points in the training set), well-scattered along the training set. Timings: construction of GMRA and model construction for all scales for GMRA took approximately 1 min, for SVD 0.3 min, for MFA about 15 hrs. Right: a similar experiment with a training set of 2000 points from a swissroll shaped manifold with no noise: the finest scale GMRA-based models perform best (in terms of both approximation and variability, the SVD-based models are once again unable to take advantage of the low-intrinsic dimension, and MFA-based models fail as well, to succeed they seem to require tuning the parameters far from the defaults, as well as a much larger training set. Timings: construction of GMRA and model construction for all scales for GMRA took approximately 4 sec, for SVD 0.5 sec, for MFA about 4 hrs.

on the other hand the sampling requirements for correctly estimating the density of $C_{j,k}$ projected on $V_{j,k}$ increases with j as less and less training points fall in $C_{j,k}$. A pruning greedy algorithm that selects, in each region of the data, the correct scale for obtaining the correct bias-variance tradeoff, depending on the samples and the geometry of the data, similar in spirit to the what has been studied in the case of multi-scale approximation of functions, will be presented in a forthcoming publication.

10. FUTURE WORK

We consider this work as a first “bare bone” construction, which may be refined in a variety of ways and opens the way to many generalizations and applications. For example:

- **User interface.** We are currently developing a user interface for interacting with the geometric wavelet representation of data sets [70].
- **Higher order approximations.** One can extend the construction presented here to piecewise quadratic, or even higher order, approximators, in order to achieve better approximation rates when the underlying set is smoother than \mathcal{C}^2 .
- **Better encoding strategies for the geometric wavelet tree.** The techniques discussed in this paper are not expected to be optimal, and better tree pruning/tuning constructions may be devised. In particular, to optimize the encoding cost of a data set, the geometric wavelet tree should be pruned and slightly modified to use a near-minimal number of dictionary elements to achieve a given approximation precision ϵ .
- **Sparsifying dictionary.** While the approximation only depends on the subspaces $\langle \Phi_{j,k} \rangle$, the sparsity of the representation of the data points will in general depend on the choice of $\Phi_{j,k}$ and $\Psi_{j,k}$, and such choice may be optimized (“locally” in space and in dimension) by existing algorithms, thereby retaining both the approximation guarantees and the advantages of running these black-box algorithms only on small number of samples and in a low-dimensional subspace.
- **Probabilistic construction.** One may cast the whole construction in a probabilistic setting, where subspaces are enriched with distributions on those subspaces, thereby allowing geometric wavelets to generate rich families of probabilistic models.

11. APPENDIX

Proof of Theorem 2.3. . The first equality follows by recursively applying the two-scale equation (2.19), so we only need to prove the upper bound. We start with the case $p = +\infty$. By compactness, for every $x \in \mathcal{M}$ and for j_0 large enough and $j \geq j_0$, there is a unique point $z_{j,x} \in \mathcal{M}$ closest to $c_{j,x}$, and $C_{j,x}$ is the graph of a $C^{1+\alpha}$ function $f := f_{j,x} : P_{T_{z_{j,x}}}(\mathcal{M}) \rightarrow C_{j,x}$, where $T_{z_{j,x}}(\mathcal{M})$ is the plane tangent to \mathcal{M} at $z_{j,x}$. Note that this is true whether we construct dyadic cells $C_{j,x}$ with respect to the manifold metric ρ , or by intersecting Euclidean dyadic cubes with \mathcal{M} . The following calculations are in the spirit of those in [8]. Since all the quantities involved are invariant under rotations and translations, up to a change of coordinates we may assume that $f(z_{j,x}) = 0$, $T_{z_{j,x}} = \langle x_1, \dots, x_d \rangle$. Assume $\alpha = 1$, i.e. the manifold is \mathcal{C}^2 . In the coordinates above the function $f =: (f_1, \dots, f_{D-d})$ above may be written

$$(11.1) \quad f_i(w) = \frac{1}{2}(w - z_{j,x})^T H_i f|_{z_{j,x}}(w - z_{j,x}) + o(\|w - z_{j,x}\|^2),$$

where H_i is the $d \times d$ Hessian of the i -th coordinate f_i of f . The calculations in [8] show that, up to higher order terms, $\mathbb{V}_{j,x}$ is parallel to $T_{z_{j,x}}$, and differs from it by a translation along the normal space $N_{c_{j,x}}$, since $\mathbb{V}_{j,x}$ passes through $c_{j,x}$ while

T_{z_jx} passes through $z_{j,x}$. Therefore we have

$$\begin{aligned}
& \left\| \|z - P_{\mathcal{M}_j}(z)\|_{\mathbb{R}^D} \right\|_{L^\infty(C_{j,x})} = \sup_{z \in C_{j,x}} \|z - \mathbb{P}_{j,x}(z)\|_{\mathbb{R}^D} \\
&= \sup_{z \in C_{j,x}} \|z - P_{T_{z_j,x}}(z - c_{j,x}) - c_{j,x}\|_{\mathbb{R}^D} \\
&\leq \sup_{z \in C_{j,x}} \|(z - z_{j,x}) - P_{T_{z_j,x}}(z - z_{j,x})\|_{\mathbb{R}^D} + \|z_{j,x} - c_{j,x}\|_{\mathbb{R}^D} \\
&\leq \sup_{w \in P_{T_{z_j,x}}(C_{j,x})} \left\| \frac{1}{2}(w - z_{j,x})^* H_i f|_{z_{j,x}} (w - z_{j,x}) + o(\|w - z_{j,x}\|^2) \right\|_{\mathbb{R}^D} \\
&\quad + \|z_{j,x} - c_{j,x}\|_{\mathbb{R}^D} \\
&\leq 2\kappa 2^{-2j} + o(2^{-2j}),
\end{aligned}$$

where $\kappa = \frac{1}{2} \max_{i \in \{1, \dots, D-d\}} \|H_i\|$ is a measure of extrinsic curvature, and where we used that $c_{j,x}$ is in the convex hull of $C_{j,x}$. A similar calculation applies to the case where $f_i \in \mathcal{C}^{1+\alpha}$, where $O(\|w - z_{j,x}\|^{1+\alpha})$ replaces the second order terms, and κ is replaced by $\max_{i \in \{1, \dots, D-d\}} \|\nabla f_i\|_{\mathcal{C}^\alpha}$.

We now derive an $L^2(C_{j,x}, \mu_{j,x})$ estimate:

$$\begin{aligned}
& \left\| \|z - P_{\mathcal{M}_j}(z)\|_{\mathbb{R}^D} \right\|_{L^2(C_{j,x}, d\mu_{j,x}(z))}^2 \\
&= \frac{1}{\mu(C_{j,x})} \int_{C_{j,x}} \|z - \mathbb{P}_{j,x}(z)\|_{\mathbb{R}^D}^2 d\mu(z) \\
&= \min_{\Pi: \text{ an affine } d-\text{plane}} \frac{1}{\mu(C_{j,x})} \int_{C_{j,x}} \|z - P_\Pi(z)\|^2 d\mu(z) \\
&= \sum_{l=d+1}^D \lambda_l(\text{cov}_{j,x}) \\
&\leq \frac{d(d+1)}{2} \lambda_{d+1}(\text{cov}_{j,x}) + o(2^{-4j}) \\
&\leq \max_{w \in \mathbb{S}^{D-d}} \frac{d(d+1)}{4(d+2)(d+4)} \left[\left\| \sum_{l=1}^{D-d} w_l H_l \right\|_F^2 - \frac{1}{d+2} \left(\sum_{l=1}^{D-d} w_l \text{Tr}(H_l) \right)^2 \right] 2^{-4j} \\
&\quad + o(2^{-4j}),
\end{aligned}$$

where the inequality before the last follows from the fact that, up to order 2^{-4j} , there are no more than $d(d+1)/2$ curvature directions, and the last inequality follows from the bounds in [8], which formalize the fact that the eigenspace spanned by the top d vectors of $\text{cov}_{j,x}$ is, up to higher order, parallel to the tangent plane, and passing through a point $c_{j,x}$ which is second-order close to \mathcal{M} , and therefore provides a second-order approximation to \mathcal{M} at scale 2^{-j} . This latter bounds could be strengthened in obvious ways if some decay of $\lambda_l(\text{cov}_{j,x})$ for $l = d+1, \dots, d(d+1)/2$ was assumed. The estimate in (2.20) follows by interpolation between the estimate in L^2 and the one in L^∞ . \square

The measure of curvature multiplying 2^{-4j} in the last bound appeared in [8]: it may be as large as $O((D-d)\kappa^2)$, but also quite small depending on the eigenvalues of the Hessians H_l .

REFERENCES

- [1] R. Coifman, S. Lafon, M. Maggioni, Y. Keller, A. Szlam, F. Warner, S. Zucker, Geometries of sensor outputs, inference, and information processing, in: J. C. Z. E. Intelligent Integrated Microsystems; Ravindra A. Athale (Ed.), Proc. SPIE, Vol. 6232, 2006, p. 623209.
- [2] E. Causevic, R. Coifman, R. Isenhart, A. Jacquin, E. John, M. Maggioni, L. Prichep, F. Warner, QEEG-based classification with wavelet packets and microstate features for triage applications in the ER, Vol. 3, ICASSP Proc., 2006, 10.1109/ICASSP.2006.1660859.
- [3] I. U. Rahman, I. Drori, V. C. Stodden, D. L. Donoho, Multiscale representations for manifold-valued data, SIAM J. Multiscale Model. Simul. 4 (2005) 1201–1232.
- [4] R. R. Coifman, S. Lafon, A. B. Lee, M. Maggioni, B. Nadler, F. Warner, S. W. Zucker, Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps, PNAS 102 (21) (2005) 7426–7431.
- [5] R. R. Coifman, S. Lafon, A. B. Lee, M. Maggioni, B. Nadler, F. Warner, S. W. Zucker, Geometric diffusions as a tool for harmonic analysis and structure definition of data: Multiscale methods, PNAS 102 (21) (2005) 7432–7438.
- [6] R. Coifman, M. Maggioni, Geometry analysis and signal processing on digital data, emergent structures, and knowledge building, SIAM News (November 2008).
- [7] A. Little, Y.-M. Jung, M. Maggioni, Multiscale estimation of intrinsic dimensionality of data sets, in: Proc. A.A.A.I., 2009.
- [8] A. Little, M. Maggioni, L. Rosasco, Multiscale geometric methods for data sets I: Estimation of intrinsic dimension, submitted.
- [9] J. Costa, A. Hero, Learning intrinsic dimension and intrinsic entropy of high dimensional datasets, in: Proc. of EUSIPCO, Vienna, 2004.
- [10] F. Camstra, A. Vinciarelli, Intrinsic dimension estimation of data: An approach based on grassberger-procaccia's algorithm, Neural Processing Letters 14 (1) (2001) 27–34.
- [11] F. Camstra, A. Vinciarelli, Estimating the intrinsic dimension of data with a fractal-based method, IEEE P.A.M.I. 24 (10) (2002) 1404–10.
- [12] W. Cao, R. Haralick, Nonlinear manifold clustering by dimensionality, ICPR 1 (2006) 920–924.
- [13] J. B. Tenenbaum, V. D. Silva, J. C. Langford, A global geometric framework for nonlinear dimensionality reduction, Science 290 (5500) (2000) 2319–2323.
- [14] S. Roweis, L. Saul, Nonlinear dimensionality reduction by locally linear embedding, Science 290 (2000) 2323–2326.
- [15] M. Belkin, P. Niyogi, Using manifold structure for partially labelled classification, Advances in NIPS 15.
- [16] D. L. Donoho, C. Grimes, When does isomap recover natural parameterization of families of articulated images?, Tech. Rep. 2002-27, Department of Statistics, Stanford University (August 2002).
- [17] D. L. Donoho, C. Grimes, Hessian eigenmaps: new locally linear embedding techniques for high-dimensional data, Proc. Nat. Acad. Sciences (2003) 5591–5596.
- [18] Z. Zhang, H. Zha, Principal manifolds and nonlinear dimension reduction via local tangent space alignment, SIAM Journal of Scientific Computing 26 (2002) 313–338.
- [19] P. Jones, M. Maggioni, R. Schul, Manifold parametrizations by eigenfunctions of the Laplacian and heat kernels, Proc. Nat. Acad. Sci. 105 (6) (2008) 1803–1808.
- [20] P. Jones, M. Maggioni, R. Schul, Universal local manifold parametrizations via heat kernels and eigenfunctions of the Laplacian, Ann. Acad. Scient. Fen. 35 (2010) 1–44, <http://arxiv.org/abs/0709.1975>.
- [21] M. Aharon, M. Elad, A. Bruckstein, K-SVD: Design of dictionaries for sparse representation, in: PROCEEDINGS OF SPARS 05', 2005, pp. 9–12.
- [22] A. Szlam, G. Sapiro, Discriminative k -metrics, in: Proceedings of the 26th Annual International Conference on Machine Learning, 2009, pp. 1009–1016.
- [23] R. Coifman, M. Maggioni, Diffusion wavelets, Appl. Comp. Harm. Anal. 21 (1) (2006) 53–94, (Tech. Rep. YALE/DCS/TR-1303, Yale Univ., Sep. 2004).
- [24] J. Bremer, R. Coifman, M. Maggioni, A. Szlam, Diffusion wavelet packets, Appl. Comp. Harm. Anal. 21 (1) (2006) 95–112, (Tech. Rep. YALE/DCS/TR-1304, 2004).
- [25] A. Szlam, M. Maggioni, R. Coifman, J. B. Jr., Diffusion-driven multiscale analysis on manifolds and graphs: top-down and bottom-up constructions, Vol. 5914-1, SPIE, 2005, p. 59141D.

- [26] M. Maggioni, J. B. Jr., R. Coifman, A. Szlam, Biorthogonal diffusion wavelets for multiscale representations on manifolds and graphs, Vol. 5914, SPIE, 2005, p. 59141M.
- [27] S. Mahadevan, M. Maggioni, Proto-value functions: A spectral framework for solving markov decision processes, JMLR 8 (2007) 2169–2231.
- [28] M. Maggioni, S. Mahadevan, Fast direct policy evaluation using multiscale analysis of markov diffusion processes, in: ICML 2006, 2006, pp. 601–608.
- [29] A. Szlam, M. Maggioni, R. Coifman, Regularization on graphs with function-adapted diffusion processes, Jour. Mach. Learn. Res. (9) (2008) 1711–1739, (YALE/DCS/TR1365, Yale Univ, July 2006).
- [30] R. Coifman, M. Maggioni, Multiscale data analysis with diffusion wavelets, Proc. SIAM Bioinf. Workshop, Minneapolis.
- [31] R. Coifman, Y. Meyer, S. Quake, M. V. Wickerhauser, Signal processing and compression with wavelet packets, in: Progress in wavelet analysis and applications (Toulouse, 1992), Frontières, Gif, 1993, pp. 77–93.
- [32] E. Candès, D. L. Donoho, Curvelets: A surprisingly effective nonadaptive representation of objects with edges, in: L. L. S. et al. (Ed.), Curves and Surfaces, Vanderbilt University Press, Nashville, TN, 1999.
- [33] S. S. Chen, D. L. Donoho, M. A. Saunders, Atomic decomposition by basis pursuit, SIAM Journal on Scientific Computing 20 (1) (1998) 33–61.
- [34] I. Daubechies, Ten lectures on wavelets, Society for Industrial and Applied Mathematics, 1992.
- [35] O. Christensen, An introduction to frames and Riesz bases, Applied and Numerical Harmonic Analysis, Birkhäuser Boston Inc., Boston, MA, 2003.
- [36] J. I. Starck, M. Elad, D. Donoho, Image decomposition via the combination of sparse representations and a variational approach, IEEE Transactions on Image Processing 14 (2004) 1570–1582.
- [37] P. Casazza, G. Kutyniok, Frames of subspaces, Contemporary Math. 345 (2004) 87–114.
- [38] B. A. Olshausen, D. J. Field, Sparse coding with an overcomplete basis set: A strategy employed by V1?, Vision Research (37).
- [39] M. Zhou, H. Chen, J. Paisley, L. Ren, G. Sapiro, L. Carin, Non-parametric Bayesian dictionary learning for sparse image representations, in: Neural and Information Processing Systems (NIPS), 2009.
- [40] J. Mairal, F. Bach, J. Ponce, G. Sapiro, Online dictionary learning for sparse coding, in: ICML, 2009, p. 87.
- [41] J. Mairal, F. Bach, J. Ponce, G. Sapiro, Online learning for matrix factorization and sparse coding, Journ. Mach. Learn. Res. 11 (2010) 19–60.
- [42] R. Tibshirani, Regression shrinkage and selection via the lasso, J. Royal. Statist. Soc B. 58 (1) (1996) 267–288.
- [43] P. W. Jones, Rectifiable sets and the traveling salesman problem, Invent. Math. 102 (1) (1990) 1–15.
- [44] G. David, S. Semmes, Analysis of and on uniformly rectifiable sets, Vol. 38 of Mathematical Surveys and Monographs, American Mathematical Society, Providence, RI, 1993.
- [45] P. Binev, R. Devore, Fast computation in adaptive tree approximation, Numer. Math. (2004) 193–217.
- [46] P. Binev, A. Cohen, W. Dahmen, R. Devore, V. Temlyakov, Universal algorithms for learning theory part i: Piecewise constant functions, Journ. Mach. Learn. 6 (2005) 1297–1321.
- [47] S. Mallat, A theory for multiresolution signal decomposition: The wavelet representation, IEEE Trans. Pattern Anal. Mach. Intell. 11 (7) (1989) 674–693.
- [48] S. Mallat, Multiresolution approximations and wavelet orthonormal bases of $l^2(\mathbb{R})$, Trans Amer Math Soc (315) (1994) 69–87.
- [49] S. Mallat, A wavelet tour in signal processing, Academic Press, 1998.
- [50] Y. Meyer, Ondelettes et Opératateurs, Hermann, Paris, 1990.
- [51] M. Christ, A $T(b)$ theorem with remarks on analytic capacity and the Cauchy integral, Colloq. Math. 60/61 (2) (1990) 601–628.
- [52] G. David, Wavelets and singular integrals on curves and surfaces, Vol. 1465 of Lecture Notes in Mathematics, Springer-Verlag, Berlin, 1991.
- [53] G. David, Wavelets and Singular Integrals on Curves and Surfaces, Springer-Verlag, 1991.
- [54] R. Coifman, S. Lafon, Diffusion maps, Appl. Comp. Harm. Anal. 21 (1) (2006) 5–30.

- [55] M. A. Rohrdanz, W. Zheng, M. Maggioni, C. Clementi, Determination of reaction coordinates via locally scaled diffusion map, submitted.
- [56] G. Karypis, V. Kumar, A fast and high quality multilevel scheme for partitioning irregular graphs, SIAM Journal on Scientific Computing 20 (1) (1999) 359–392.
- [57] A. Beygelzimer, S. Kakade, J. Langford, Cover trees for nearest neighbor, in: ICML, 2006, pp. 97–104.
- [58] P. Binev, W. Dahmen, P. Lamby, Fast high-dimensional approximation with sparse occupancy trees, Journal of Computational and Applied Mathematics 235 (8) (2011) 2063 – 2076.
- [59] H. Federer, Curvature measures, Trans. Am. Math. Soc. 93 (3) (1959) 418–491.
- [60] R. Baraniuk, M. Wakin, Random projections of smooth manifolds, preprint.
- [61] P. Niyogi, S. Smale, S. Weinberger, Finding the homology of submanifolds with high confidence from random samples, Discrete and Computational Geometry 39 (2008) 419–441, 10.1007/s00454-008-9053-2.
- [62] P. W. Jones, The traveling salesman problem and harmonic analysis, Publ. Mat. 35 (1) (1991) 259–267, conference on Mathematical Analysis (El Escorial, 1989).
- [63] G. David, S. Semmes, Uniform Rectifiability and Quasiminimizing Sets of Arbitrary Codimension, AMS.
- [64] A. Little, J. Lee, Y.-M. Jung, M. Maggioni, Estimation of intrinsic dimensionality of samples from noisy low-dimensional manifolds in high dimensions with multiscale SVD, in: Proc. S.S.P., 2009.
- [65] G. Golub, C. V. Loan, Matrix Computations, Johns Hopkins University Press, 1989.
- [66] A. Szlam, Asymptotic regularity of subdivisions of euclidean domains by iterated PCA and iterated 2-means, Appl. Comp. Harm. Anal. 27 (3) (2009) 342–350.
- [67] R. R. Coifman, D. Donoho, Translation-invariant de-noising, Springer-Verlag, 1995, pp. 125–150.
- [68] Y. Bartal, Probabilistic approximation of metric spaces and its algorithmic applications, in: In 37th Annual Symposium on Foundations of Computer Science, 1996, pp. 184–193.
- [69] V. Rokhlin, A. Szlam, M. Tygert, A randomized algorithm for principal component analysis, SIAM Jour. Mat. Anal. Appl. 31 (3) (2009) 1100–1124.
- [70] E. Monson, G. Chen, R. Brady, M. Maggioni, Data representation and exploration with geometric wavelets, in: Visual Analytics Science and Technology (VAST), 2010 IEEE Symposium, 2010, pp. 243–244.

MATHEMATICS DEPARTMENT, DUKE UNIVERSITY, P.O. Box 90320, DURHAM, NC 27708, U.S.A.
E-mail address: wka@math.duke.edu

MATHEMATICS DEPARTMENT, DUKE UNIVERSITY, P.O. Box 90320, DURHAM, NC 27708, U.S.A.
E-mail address: gchen@math.duke.edu

MATHEMATICS AND COMPUTER SCIENCE DEPARTMENTS, DUKE UNIVERSITY, P.O. Box 90320, DURHAM, NC 27708, U.S.A.
E-mail address: mauro@math.duke.edu (corresponding author)