

# Multi-Scale Vector Quantization with Reconstruction Trees

E. Cecini<sup>†</sup>, E. De Vito<sup>†</sup>, L. Rosasco<sup>‡,\*†</sup>

<sup>‡</sup> DIBRIS, Università degli Studi di Genova

<sup>†</sup> Dipartimento di Matematica, Università di Genova

<sup>\*</sup> Istituto Italiano di Tecnologia

<sup>†</sup> Massachusetts Institute of Technology

September 5, 2019

## Abstract

We propose and study a multi-scale approach to vector quantization. We develop an algorithm, dubbed reconstruction trees, inspired by decision trees. Here the objective is parsimonious reconstruction of unsupervised data, rather than classification. Contrasted to more standard vector quantization methods, such as  $K$ -means, the proposed approach leverages a family of given partitions, to quickly explore the data in a coarse to fine –multi-scale– fashion. Our main technical contribution is an analysis of the expected distortion achieved by the proposed algorithm, when the data are assumed to be sampled from a fixed unknown distribution. In this context, we derive both asymptotic and finite sample results under suitable regularity assumptions on the distribution. As a special case, we consider the setting where the data generating distribution is supported on a compact Riemannian sub-manifold. Tools from differential geometry and concentration of measure are useful in our analysis.

## 1 Introduction

Dealing with large high-dimensional data-sets is a hallmark of modern signal processing and machine learning. In this context, finding parsimonious representation from unlabeled data is often key to both reliable estimation and efficient computations, and more generally for exploratory data analysis. A classical approach to this problem is principal component analysis (PCA), relying on the assumption that data are well represented by a linear subspace. Starting from PCA a number of developments can be considered to relax the linearity assumption. For example kernel PCA is based on performing PCA after a suitable nonlinear embedding [26]. Sparse dictionary learning tries to find a set of vectors on which the data can be written as sparse linear combinations [21]. Another line of works assumes the data to be sampled from a distribution supported on a manifold and includes isomap [29], Hessian eigenmaps [8], Laplacian eigenmaps [3] and related developments such as diffusion maps [20]. A more recent and original perspective has been proposed in [14], and called geometric multi-resolution analysis (GMRA). Here the idea is to borrow and generalize ideas from multi-resolution analysis and wavelet theory [17], to derive locally linear representation organized in a multi-scale fashion. The corresponding algorithm is based on a cascade of local PCAs

and is reminding of classical decision trees for function approximation, see e.g. [12]. In this paper we further explore the ideas introduced to GMRA which is our main reference.

Indeed, we consider these ideas in the context of vector quantization, which is another classical, and extreme, example of parsimonious representation. Here, a set of centers and corresponding partition is considered, and then all data points in each cell of the partition are represented by the corresponding center. The most classical approach in this context is probably  $k$ -means, where a set of centers (means) is defined by a non-convex optimization problem over all possible partitions. Our approach offers an alternative to  $k$ -means, by following the basic idea of GMRA and decision trees, but considering local means rather than local PCA. In this view, our approach can be seen as a zero-th order version of GMRA, hence providing a piece-wise constant data approximation. Compared to  $k$ -means, the search for a partition is performed through a coarse-to-fine recursive procedure, rather than by global optimization. A strategy that we call reconstruction tree. As a byproduct the corresponding vector quantization is multi-scale, and naturally yields a multi-resolution representation of the data. Our main technical contribution is a theoretical analysis of the above multi-scale vector quantization procedure. We consider a statistical learning framework, where the data are assumed to be sampled according to some fixed unknown distribution and measure performance according to the so called expected distortion, measuring the reconstruction error with respect to the whole data distribution. Our main result is deriving corresponding finite sample bound in terms of natural geometric assumptions.

The rest of the paper is organized as follows. After describing the basic ideas in the context of vector quantization in Section 2, we present the algorithm we study in Section 3. In Section 4, we introduce the basic theoretical assumptions needed in our analysis, and illustrate them considering the case where the data are sampled from a manifold. In Section 5, we present and discuss our main results, and detail the main steps in their proofs. All other proofs are deferred to the Appendix.

## 2 Vector quantization & distortion

We next introduce the problem of interest and comment on its connections with related questions.

A vector quantization (VQ) procedure is defined by a set of code vectors/centers and an associated partition of the data space. The idea is that compression can be achieved replacing all points in a given cell of the partition by the corresponding code vector. More precisely, assuming that the data space is  $\mathbb{R}^D$ , consider a set of code vectors  $c_1 \in I_1, \dots, c_k \in I_k$ , where the set of cells  $\Lambda = \{I_1, \dots, I_k\}$ , defines a partition of  $\mathbb{R}^D$ . A nonlinear projection  $P_\Lambda : \mathbb{R}^D \rightarrow \mathbb{R}^D$  can be defined by

$$P_\Lambda(x) = \sum_{j=1}^k c_j \mathbb{1}_{I_j}(x) \quad x \in \mathbb{R}^D, \quad (1)$$

where  $\mathbb{1}_I$  is the characteristic function of  $I$ , *i.e.*

$$\mathbb{1}_I(x) = \begin{cases} 1 & x \in I \\ 0 & x \notin I \end{cases}.$$

Given a set of points  $x_1, \dots, x_n$  the error (distortion) incurred by this nonlinear projection can be defined as

$$\widehat{\mathcal{E}}[P_\Lambda] = \frac{1}{n} \sum_{i=1}^n \|x_i - P_\Lambda(x_i)\|^2,$$

where  $\|\cdot\|$  is the Euclidean norm in  $\mathbb{R}^D$ . If we consider the data to be identical and independent samples of a random variable  $X$  in  $\mathbb{R}^D$ , then the following error measure can also be considered

$$\mathcal{E}[P_\Lambda] = \mathbb{E}[\|X - P_\Lambda(X)\|^2]. \quad (2)$$

The above error measure is the expected distortion associated to the quantization defined by  $P_\Lambda$ .

In the following we are interested in deriving VQ schemes with small expected distortion given a dataset  $x_1, \dots, x_n$  of  $n$  samples of  $X$ . Before describing the algorithm we propose, we add two remarks.

**Remark 1** (Comparison to supervised learning).

*Classical supervised learning is concerned with the problem of inferring a functional relationship  $f$  given a set of input-output pairs  $(x_1, y_1), \dots, (x_n, y_n)$ . A classical error measure is the least squares loss  $\|y - f(x)\|^2$  (if the outputs are vectors valued). A parallel between the above setting and supervised learning can be seen, considering the case where the input and output spaces coincide and the least squares loss would be  $\|x - f(x)\|^2$ . Clearly, in this case an optimal solution is given by the identity map, unless further constraints are imposed. Following the above remark, we can view the nonlinear projection  $P_\Lambda$  as a piece-wise constant approximation of the identity map, possibly providing a parsimonious representation.*

**Remark 2** (Vector quantization as Dictionary Learning). *A dictionary is a set of vectors (called atoms)  $a_1, \dots, a_p$  that can be used to approximately decompose each point  $x$  in the data space, i.e.  $x \approx \sum_{j=1}^p a_j \beta_j$ , with  $\beta = (\beta_1, \dots, \beta_p) \in \mathbb{R}^p$  being a coefficients vector. Given a set of points  $x_1, \dots, x_n$ , dictionary learning is the problem of estimating a dictionary, as well as a set of coefficients vectors. Vector quantization can be viewed as a form of dictionary learning where the code vectors are the atoms, and coefficients vectors are binary and have at most one non zero component [21].*

The above remarks provide different views on the problem under study. We next discuss how ideas from decision trees in supervised learning can be borrowed to define a novel VQ approach.

### 3 Multi-scale vector quantization via reconstruction trees

We next describe our approach to Multi-Scale Vector Quantization (MSVQ), based on a recursive procedure that we call reconstruction trees, since it is inspired by decision trees for function approximation. The key ingredient in the proposed approach is a family of partitions organized in a tree. The partition at the root of the tree has the largest cells, while partitions with cells of decreasing size are found in lower leaves. This *partition tree* provides a multi-scale description of the data space: the lower the leaves, the finer the scale. The idea is to use data to identify a subset of cells, and corresponding partition, providing a VQ solution (1) with low expected distortion (2). We next describe this idea in detail.

### 3.1 Partition trees and subtrees

We begin introducing the definition of a partition tree. In the following we denote by  $\mathcal{X} \subset \mathbb{R}^D$  the data space endowed with its natural Borel  $\sigma$ -algebra  $\mathcal{B}(\mathcal{X})$  and by  $\#A$  the cardinality of a set  $A$ .

**Definition 3.** A partition tree  $\mathbb{T}$  is a denumerable family  $\{\Lambda_j\}_{j \in \mathbb{N}}$  of partitions of  $\mathcal{X}$  satisfying

a)  $\Lambda_0 = \{\mathcal{X}\}$  is the root of the tree;

b) each family  $\Lambda_j = \{I\}_{I \in \Lambda_j}$  is a finite partition of  $\mathcal{X}$  of Borel subsets, i.e

$$\begin{cases} \mathcal{X} = \bigcup_{I \in \Lambda_j} I \\ I \cap J = \emptyset \quad \forall I, J \in \Lambda_j, I \neq J \\ \#\Lambda_j < +\infty \\ I \in \mathcal{B}(\mathcal{X}) \quad I \in \Lambda_j \end{cases} ;$$

c) for each  $I \in \Lambda_j$ , there exists a family  $\mathcal{C}(I) \subseteq \Lambda_{j+1}$  such that

$$\begin{cases} I = \bigcup_{J \in \mathcal{C}(I)} J \\ \#\mathcal{C}(I) \leq a \end{cases}, \quad (3)$$

where  $a \in (0, +\infty)$  is a constant depending only on  $\mathbb{T}$ .

Note that, we allow the partition tree to have arbitrary, possibly infinite, depth, needed to derive asymptotic results.

Further, notice that, since  $\#\Lambda_j \leq a^j$  the constant  $a$  characterizes how the cardinality of each partition increases at finer scales. The case  $a = 2$  corresponds to dyadic trees.

We add some further definitions. For any  $j \in \mathbb{N}$ , and  $I \in \Lambda_j$  the depth of  $I$  is  $j$  and is denoted by  $j_I$ . The cells in  $\mathcal{C}(I) \subset \Lambda_{j+1}$  are the children of  $I$ , the unique cell  $J \in \Lambda_{j-1}$  such that  $I \in \mathcal{C}(J)$  is the parent of  $I$  and is denoted by  $\mathcal{P}(I)$  (by definition  $\mathcal{P}(\mathcal{X}) = \mathcal{X}$ ). We regard  $\mathbb{T}$  as a set of nodes where each node is defined by a cell  $I$  with its parent  $\mathcal{P}(I)$  and its children  $\mathcal{C}(I)$ . The following definition will be crucial.

**Definition 4.** A (proper) subtree of  $\mathbb{T}$  is a family  $\mathcal{T} \subset \mathbb{T}$  of cells such that  $\mathcal{P}(I) \in \mathcal{T}$  for all  $I \in \mathcal{T}$  and

$$\Lambda(\mathcal{T}) = \{I \in \mathbb{T} : I \notin \mathcal{T}, \mathcal{P}(I) \in \mathcal{T}\},$$

denotes the set of outer leaves.

It is important in what follows that  $\Lambda(\mathcal{T})$  is a partition of  $\mathcal{X}$  if  $\mathcal{T}$  is finite, see Lemma 11.

### 3.2 Reconstruction trees

We next discuss a data driven procedure to derive a suitable partition and a corresponding nonlinear projection. To do this end, we need a few definitions depending on an available dataset  $x_1, \dots, x_n$ .

For each cell  $I$ , we fix an arbitrary point  $\hat{x}_I^* \in \mathcal{X}$  and define the corresponding cardinality and center of mass, respectively, as

$$n_I = \sum_{i=1}^n \mathbb{1}_I(x_i), \quad \hat{c}_I = \begin{cases} \frac{1}{n_I} \sum_{i=1}^n x_i \mathbb{1}_I(x_i) & \text{if } x \in I \text{ and } n_I \neq 0 \\ \hat{x}_I^* & \text{if } x \in I \text{ and } n_I = 0 \end{cases}. \quad (4)$$

If  $0 \in \mathcal{X}$ , a typical choice is  $\hat{x}_I^* = 0$  for all cells  $I \in \mathbb{T}$ . While  $\mathcal{E}(\hat{P}_n)$  depends on the choice of  $\hat{x}_I^*$ , our bounds hold true for all choices. We point out that it is more convenient to choose  $\hat{x}_I^* \in I$ , as this (arbitrary) choice produces an improvement of  $\mathcal{E}(\hat{P}_n)$  for free, in particular whenever  $\mathbb{E}[X \in I] > 0$  but  $n_I = 0$ .

Using this quantity we can define a local error measure for each cell  $I$ ,

$$\hat{\mathcal{E}}_I = \frac{1}{n} \sum_{x_i \in I} \|x_i - \hat{c}_I\|^2 = \frac{1}{n} \sum_{i=1}^n \|x_i - \hat{c}_I\|^2 \mathbb{1}_I(x_i),$$

as well as the potential error difference induced by considering a refinement,

$$\hat{\epsilon}_I^2 = \hat{\mathcal{E}}_I - \sum_{J \in \mathcal{C}(I)} \hat{\mathcal{E}}_J = \frac{1}{n} \sum_{J \in \mathcal{C}(I)} \|\hat{c}_J - \hat{c}_I\|^2, \quad (5)$$

where the second equality is consequence of the between-within decomposition of the variance. Following [4], we first truncate the partition tree at a given depth, depending on the size of the data set. More precisely, given  $\gamma > 0$ , we set

$$j_n = \left\lfloor \frac{\gamma \ln n}{\ln a} \right\rfloor \implies a^{j_n} \leq n^\gamma. \quad (6)$$

Deeper trees are considered as data size grows.

As a second step, we fix a threshold  $\eta > 0$  and select the cells such that  $\hat{\epsilon}_I \geq \eta$ . Since  $\hat{\epsilon}_I$  is not an decreasing function with the depth of the tree, this requires some care – see Remark 10 for an alternative construction. Indeed, we define the  $\eta$ -dependent subtree

$$\hat{\mathcal{T}}_\eta = \begin{cases} \{\mathcal{X}\} & \text{if } \hat{\epsilon}_I < \eta \quad \forall I \in \bigcup_{j \leq j_n} \Lambda_j \\ \{I \in \mathbb{T} \mid \exists j \leq j_n, J \in \Lambda_j, J \subset I, \hat{\epsilon}_J \geq \eta\} & \text{otherwise} \end{cases} \quad (7)$$

and  $\hat{\Lambda}_\eta$  is defined as outerleaves of  $\hat{\mathcal{T}}_\eta$ , *i.e.*  $\hat{\Lambda}_\eta = \Lambda(\hat{\mathcal{T}}_\eta)$ . Note that  $\hat{\mathcal{T}}_\eta$  is finite, so that by Lemma 11  $\hat{\mathcal{T}}_\eta$  is a partition of  $\mathcal{X}$  such that  $j_I \leq (\gamma \ln n) / \ln a$  for all  $I \in \hat{\Lambda}_\eta$ .

The code vectors are the centers of mass of the cells the above empirical partition, and the corresponding nonlinear projection is

$$\hat{P}_\eta = \sum_{I \in \hat{\Lambda}_\eta} \hat{c}_I \mathbb{1}_I(x) \quad \hat{\Lambda}_\eta = \Lambda(\hat{\mathcal{T}}_\eta). \quad (8)$$

We add a few comments, the above vector quantization procedure, that we call reconstruction tree, is recursive and depends on the threshold  $\eta$ . Different quantizations and corresponding distortions are achieved by different choices of  $\eta$ . Smaller values of  $\eta$  correspond to vector quantization

quantizations with smaller distortion. It is clear that the empirical distortion becomes zero for a suitably small  $\eta$  corresponding to having a single point in each cell. Understanding the behaviour of the expected distortion as function of  $\eta$  and the number of points is our main theoretical contribution. Before discussing these results we discuss the connection of the above approach to related ideas. A similar construction is given in [14], where however the thresholding criterion  $\eta$  depends on the scale, see Section 2.3 of the cited reference.

### 3.3 Comparison with related Topics

The above approach can be compared with a number of different ideas.

**Decision and Regression Trees** We call the above procedure Reconstruction Tree, since its definition is formally analogous to that of decision trees for supervised learning; see for example [13], Chapter 8. In particular our construction and analysis follows closely that of tree based estimators studied in [4], in the context of least square regression. As commented in Remark 1 our problem can actually be interpreted as a special instance of vector valued regression where the regression function  $f_\rho$  is the identity function from the input space to itself (regarded as the output space): indeed, referring to the notation of [4], the conditional distribution of the output  $y$  given  $x$  is the Delta measure at  $x$ . From this point on the two formalisms overlap, in that  $\|f - f_\rho\|_{L_2}^2$  becomes  $\mathbb{E}[\|X - f(X)\|^2]$ . Since the tree based estimators  $f$  considered in [4] are piecewise constant, the expected square loss cannot vanish and its analysis is non trivial. Despite the formal similarity, the two settings do exhibit distinct features. For example, the analysis in [4] is specifically formulated for scalar functions, while our analysis is necessarily vectorial in nature. In [4] a uniform bound  $|y| < M$  is imposed, while in our setting we can assume a local bound for free; namely, if  $f$  is constant on a cell  $I \subset \mathcal{X}$  then  $\|x - f(x)\|^2 \leq \text{diam}(I)^2$  for all  $x \in I$ . The present setting finds a natural instance in the case of a probability measure supported on a smooth manifold isometrically embedded in  $\mathcal{X}$  (see Section 4), while this case is hardly addressed explicitly in the literature about least square regression. For example, the manifold case is actually discussed, in the context of classification through Decision Trees, in [28]. One last point is that the present work contains explicit quantitative results about the approximation error, see Section 5.1, while similar results are not available in the setting of [4], that aspect being typically addressed indirectly in the corresponding literature.

**Empirical risk minimization** Again in analogy to supervised learning, as in [4], one can consider the minimization problem:

$$\min_{F \in \mathcal{H}} \widehat{\mathcal{E}}[F],$$

where  $\mathcal{H}$  is the (finite-dimensional) vector space of the vector fields  $F : \mathcal{X} \rightarrow \mathbb{R}^D$ , which are piecewise constant on a given partition  $\Lambda$ . There correspond a number of independent minimization problems, one for each cell in  $\Lambda$ , so that  $\widehat{P}_\eta$  from (8) is easily shown to be a minimizer. The minimizer is not unique, since the value of  $F$  is irrelevant on cells  $I \in \Lambda$  such that  $n_I = 0$ . Similar considerations

hold for  $\min_{F \in \mathcal{H}} \mathcal{E}[F]$  as well, in which case the value of  $F$  is irrelevant whenever  $\mathbb{E}[X \in I] = 0$ . See also Section 3.2, Section 5.1 and Lemma 21.

One could consider minimization over a wider class of functions, piece-wise constant on different partitions, for example on all the partitions with a given number of cells that are induced by proper subtrees  $\mathcal{T}$  of a given partition tree  $\mathbb{T}$ . This would be a combinatorial optimization problem. The algorithm defined in (7) overcomes this issue by providing a one-parameter coarse-to-fine class of partitions, such that each refinement carries local improvements  $\hat{\epsilon}$  that are uniformly bounded. As observed in [4], such a strategy is inspired by wavelet thresholding.

**Geometric multi-resolution analysis (GMRA)** A main motivation for our work is the algorithm GMRA [1, 16, 14], which introduces the idea of learning multi-scale dictionaries by geometric approximation. The main difference between GMRA and Regression Trees is that the former represents data through a piece-wise linear approximation, while the latter through a piece-wise constant approximation. More precisely, rather than considering the center of mass of the data in each cells (4), a linear approximation is obtained by (local) Principal Component Analysis, so that the data belonging to a cell are sent to a linear subspace of suitable dimension, the latter approach being particularly natural in the case of data supported on a manifold. Another difference is in the thresholding strategy: unlike [4] and our work, in [14] the local improvement  $\hat{\epsilon}$  is scaled depending on its depth in the tree. One of our purposes is to check if these *simpler* choices affect the learning rates significantly. We provide more quantitative comparisons later in Section 5.

**Wavelets** A main motivation for GMRA is extending ideas from wavelets and multi-resolution analysis to the context of machine learning, where, given the potential high dimensionality, non-regular partitions need be considered; this point is discussed in [1, 16, 14] and references therein. Indeed partition trees generalize the classic notion of dyadic partitions. In this view, given the piece-wise constant nature of reconstruction trees, a parallel can be drawn between the latter and classical Haar wavelets.

**$k$ -Means** Our procedure being substantially a vector quantization algorithm, a comparison with the most common approach to vector quantization, namely  $k$ -means, is in order. In  $k$ -means, a set of  $k$  code vectors  $c_1, \dots, c_k$  are derived from the data and used to define corresponding partitions via the corresponding Voronoi diagram

$$V_j = \{x \in \mathbb{R}^D \mid \|x - c_j\| \leq \|x - c_i\|, \forall i = 1, \dots, k, i \neq j\}.$$

Code vectors are defined by the minimization of the following empirical objective

$$\min_{c_1, \dots, c_k} \frac{1}{n} \sum_{i=1}^n \min_{j=1, \dots, k} \|x_i - c_j\|^2.$$

This minimization problem is non convex and is typically solved by alternating minimization, a procedure referred to as Lloyd’s algorithm [15]. The inner iteration assigns each point to a center, hence a corresponding Voronoi cell. The output minimization can be easily shown to update the code vectors by computing the center of mass, *the mean*, of each Voronoi cell. In general the

algorithm is ensures to decrease or at least not increase the objective function and to converge in finite time to a local minimum. Clearly, the initialization is important, and initializations exist yielding some stronger convergence guarantees. In particular,  $k$ -means++ is a random initialization providing on average an  $\epsilon$ -approximation to the global minimum [2].

Compared to  $k$ -means, reconstruction trees restrict the search for a partition over a prescribed family defined by the partition tree. In turns, they allow a fast multi-scale exploration of the data, while  $k$ -means requires solving a new optimization problem each time  $k$  is changed. Indeed it can be shown that a solution for the  $(k-1)$ -means problem leads to a bad initialization for the  $k$ -means problem. In other words, unlike restriction trees, the partitions found by  $k$ -means at different scales (different values of  $k$ ) are generally unrelated, and cannot be seen as refinements of one another.

**Hierarchical clustering** Lastly, our coarse-to-fine approach can be compared with hierarchical clustering, in particular [with the so called Ward’s method, which proceeds in the opposite way. Indeed, this algorithm produces a coarser partition of the data starting from a finer. It starts with a Voronoi partition having all the data as centers, and at each step it merges a couple of cells that have the smallest so called between cluster inertia [30]. Interestingly this definition has an analogue in our algorithm. Our  $\widehat{\mathcal{E}}_I$  corresponds to the within cluster inertia of a cell  $I$  while  $\widehat{c}_I^2$  to the between cluster inertia (up to a factor  $1/n$ ) of cells that merge into  $I$ . Nevertheless the obtained partitions will not in general coincide, unless very specific choices are made ad hoc.

## 4 General assumptions and manifold setting

In this section, we introduce our main assumptions and then discuss a motivating example where data are sampled at random from a manifold.

We consider a statistical learning framework, in the sense that we assume the data to be random samples from an underlying probability measure. More precisely, we assume the available data to be a realization of  $n$  identical and independent random vectors  $X_1, \dots, X_n$  taking values in a bounded subset  $\mathcal{X} \subset \mathbb{R}^D$  and we denote by  $\rho$  the common law. Up to a rescaling and a translation, we assume that  $0 \in \mathcal{X}$  and

$$\text{diam}(\mathcal{X}) = \sup_{x,y \in \mathcal{X}} \|x - y\| \leq 1. \tag{9}$$

Our main assumption relates the distribution underlying the data to the partition tree to be used to derive a MSVQ via reconstruction trees. To state it, we recall the notion of essential diameter of a cell  $I$ , namely

$$\text{diam}_\rho(I) = \inf_{\substack{J \subset I \\ \rho(J)=0}} \text{diam}(I \setminus J).$$

**Assumption 5.** *There exists  $s > 0$  and  $b > 1$  such that for all  $I \in \mathbb{T}$*

$$\text{diam}_\rho(I) \leq C_1 \rho(I)^s \tag{10a}$$

$$\text{diam}_\rho(I) \leq C_2 b^{-j_I} \tag{10b}$$

where  $C_1 > 0$  and  $C_2 > 0$  are fixed constants depending only on  $\mathbb{T}$ .

To simplify the notation, we write  $c_{\mathbb{T}}$  for a constant depending only on  $s, b, C_1, C_2$  and we write  $A \lesssim B$  if there exists a constant  $c_{\mathbb{T}} > 0$  such that  $A \leq c_{\mathbb{T}}B$ .

Given the partition tree  $\mathbb{T}$ , the parameters  $s$  and  $b$  define a class  $\mathcal{P}_{b,s}(\mathbb{T})$  of probability measures  $\rho$  and for this class we are able to provide a finite sample bound on the distortion error of our estimator  $\hat{P}_\eta$ , see (14). In the context of supervised machine learning  $\mathcal{P}_{b,s}(\mathbb{T})$  is an a-priori class of distributions defining an upper learning rate, see (16a). It is an open problem to provide a lower min-max learning rate.

Clearly, (10b) is implied by the distribution-independent assumption

$$\text{diam}(I) \lesssim b^{-j_I} \quad \text{for all } I \in \mathbb{T}, \quad (11)$$

*i.e.* the diameter of the cells goes to zero exponentially with their depth. This assumption ensures that the reconstruction error goes to zero and, in supervised learning, it corresponds to the assumption that the hypotheses space is rich enough to approximate any regression function, compare with condition (A4) in [14].

Eq. (10a) is a sort of regularity condition on the shape of the cells and, if it holds true, (10b) is implied by the following condition

$$\rho(I) \lesssim c^{-j_I} \quad \text{for all } I \in \mathbb{T}, \quad (12)$$

which states that the volume of the cells goes to zero exponentially with their depth.

In [14], following ideas from [4], it is introduced a suitable model class, see Definition 5, in terms of the decay of the approximation error, compare Eq. (7) of [14] with (19) below. This important point is further discussed in Section 5.1.

In many cases the parameter  $s$  is related to the intrinsic dimension of the data. For example, if  $\mathcal{X} = [0, 1]^D$  is the unit cube and  $\rho$  is given by

$$\rho(E) = \int_E p(x) dx \quad E \in \mathcal{B}(\mathcal{X}),$$

where  $dx$  is the Lebesgue measure of  $\mathbb{R}^D$  and the density  $p$  is bounded from above and away from zero, see (13b) below, it is easily to check that the family  $\mathbb{T} = \{\Lambda_j\}$  of dyadic cubes

$$\Lambda_j = \{[2^{-j}(k_1 - 1), 2^{-j}k_1) \times \dots \times [2^{-j}(k_D - 1), 2^{-j}k_D) \mid k_1, \dots, k_D = 1, \dots, 2^j\} \quad j \in \mathbb{N}$$

is a partition tree satisfying Assumption 5 with  $s = 1/D$  and a suitable  $b > 1$ . The construction of dyadic cubes can be extended to more general settings, see [7, 9] and references therein, by providing a large class of other examples, as shown by the following result. The proof is deferred to Section 6.

**Proposition 6.** *Assume that the support  $\mathcal{M}$  of  $\rho$  is a connected submanifold of  $\mathbb{R}^D$  and the distribution  $\rho$  is given by*

$$\rho(E) = \int_{E \cap \mathcal{M}} p(x) d\rho_{\mathcal{M}}(x) \quad E \in \mathcal{B}(\mathcal{X}) \quad (13a)$$

$$0 < p_1 \leq p(x) \leq p_2 < +\infty \quad x \in \mathcal{M}, \quad (13b)$$

where  $\rho_{\mathcal{M}}$  is the Riemannian volume element of  $\mathcal{M}$ , then there exists a partition tree  $\mathbb{T}$  of  $\mathcal{X}$  satisfying Assumption 5 with  $s = 1/d$ , where  $d$  is the intrinsic dimension of  $\mathcal{M}$ .

We recall that, as a submanifold of  $\mathbb{R}^D$ ,  $\mathcal{M}$  becomes a compact Riemannian manifold with Riemannian distance  $d_{\mathcal{M}}$  and Riemannian volume element  $\rho_{\mathcal{M}}$ . We stress that the construction of the dyadic cubes only depend on  $d_{\mathcal{M}}$ . Proposition 6 has to be compared with Proposition 3 and Lemma 6 in [14].

By inspecting the proof of the above result, it is possible to show that a partition tree satisfying Assumptions 5 always exists if there are a metric  $d$  and a Borel measure  $\nu$  on  $\mathcal{M}$  such that  $(\mathcal{M}, d, \nu)$  is an Ahlfors regular metric measure [10, p. 413],  $\text{clg}\rho$  has density  $p$  with respect to  $\nu$  satisfying (13b) and the embedding of  $(\mathcal{M}, d)$  into  $(\mathbb{R}^d, \|\cdot\|)$  is a Lipschitz function.

## 5 Main result

In this section we state and discuss our main results, characterizing the expected distortion of reconstruction trees. The proofs are deferred to Section 6. Our first result is a probabilistic bound for any given threshold  $\eta$ . Recall that  $s > 0$  is defined by (10a) and  $\hat{P}_\eta$  by (7) and (8).

**Theorem 7.** *Fix  $\gamma > 0$  as in (6) and  $\eta > 0$ , for any  $0 < \sigma < s$*

$$\mathbb{P} \left[ \mathcal{E}[\hat{P}_\eta] \gtrsim \eta^{\frac{4\sigma}{2\sigma+1}} (1+t) \right] \lesssim \eta^{-\frac{2}{2\sigma+1}} \exp(-c_{\mathbb{T}} n \eta^2 t) + (n^\gamma + \eta^{-\frac{2}{2\sigma+1}}) \exp(-c_a n \eta^2) \quad t > 0, \quad (14)$$

where  $c_a = \frac{1}{128(a+1)}$  and  $c_{\mathbb{T}} > 0$  depends on the partition tree  $\mathbb{T}$ .

As shown in Remark 19, it is possible to set  $\sigma = s$  up to an extra logarithmic factor.

Next, we show how it allows derive a choice for  $\eta$  as a function of the number of examples, and a corresponding expected distortion bound.

**Corollary 8.** *Fix  $\gamma > 1$ ,  $\beta > 0$  and set*

$$\eta_n = \sqrt{\frac{(\gamma + \beta) \ln n}{c_a n}} \quad \text{and} \quad \hat{P}_n = \hat{P}_{\eta_n} \quad n \geq 1, \quad (15)$$

where  $c_a = \frac{1}{128(a+1)}$ . Then for any  $0 < \sigma < s$

$$\mathbb{P} \left[ \mathcal{E}[\hat{P}_n] \gtrsim \left( \frac{\ln n}{n} \right)^{\frac{2\sigma}{2\sigma+1}} (1+t) \right] \lesssim \frac{1}{n^\beta} + \frac{1}{n^{\bar{c}_{\mathbb{T}} t - 1}}, \quad (16a)$$

where  $\bar{c}_{\mathbb{T}} > 0$  is a constant depending on the partition tree  $\mathbb{T}$ . Furthermore

$$\lim_{t \rightarrow +\infty} \limsup_{n \rightarrow +\infty} \sup_{\rho \in \mathcal{P}_{b,s}(\mathbb{T})} \mathbb{P} \left[ \mathcal{E}[\hat{P}_n] \gtrsim \left( \frac{\ln n}{n} \right)^{\frac{2\sigma}{2\sigma+1}} t \right] = 0, \quad (16b)$$

where  $\mathcal{P}_{b,s}(\mathbb{T})$  is the family of distributions  $\rho$  such that Assumptions 5 hold true.

If  $t$  is chosen large enough so that  $\bar{c}_{\mathbb{T}} t - 1 = \beta$ , then bound (16a) reads as

$$\mathbb{P} \left[ \mathcal{E}[\hat{P}_n] \geq c_1 \left( \frac{\ln n}{n} \right)^{\frac{2\sigma}{2\sigma+1}} \right] \leq c_2 \frac{1}{n^\beta}$$

where  $c_1$  and  $c_2$  are suitable constants depending on  $\mathbb{T}$ . This bound can be compared with Theorem 8 in [14] under the assumption that  $\mathcal{M}$  is a compact  $C^\infty$  manifold. Eq. (16a) with  $s = 1/d$  gives a convergence rate of the order  $(\ln n/n)^{\frac{2p}{2p+d}}$  for any  $p = \sigma/s \in (0, 1)$ , whereas the GMRA algorithm has a rate of the order  $(\ln n/n)^{\frac{2}{2+d}}$ , see also Proposition 3 of [14]. Hence, up to a logarithmic factor our estimator has the same convergence rate of the GMRA algorithm. However it is in order to notice that our algorithm works with a cheaper representation; indeed, given the adaptive partition  $\hat{\Lambda}_\eta$ , it only requires to compute and store the centers of mass  $\{\hat{c}_I\}_{I \in \hat{\Lambda}_\eta}$ .

In a similar setting, in [6], it is shown that the  $k$ -means algorithm with a suitable choice of  $k = k_n$  depending on  $n$ , provides has a convergence rate of the order  $(1/n)^{\frac{1}{d+1}}$  and  $k$ -flat algorithm of the order  $(1/n)^{\frac{2}{d+4}}$ .

The proof of Theorem 7 relies on splitting the error in several terms. In particular, it requires studying the stability to random sampling and the approximation properties of reconstruction trees. This latter result is relevant in the context of quantization of probability measures, hence of interest in its own right. We present this result first.

Towards this end, we need to introduce the infinite sample version of the reconstruction tree. For any cell  $I \in \mathbb{T}$ , denote the volume of the cell by

$$\rho_I = \rho(I),$$

the center of mass of the cell by

$$c_I = \begin{cases} \frac{1}{\rho_I} \int_I x d\rho(x) & \text{if } \rho_I > 0 \\ x_I^* & \text{if } \rho_I = 0, \end{cases}$$

where  $x_I^*$  is an arbitrary point in  $\mathcal{X}$ . The local expected distortion in a cell by

$$\mathcal{E}_I = \int_I \|x - c_I\|^2 d\rho(x),$$

and

$$\epsilon_I^2 = \mathcal{E}_I - \sum_{J \in \mathcal{C}(I)} \mathcal{E}_J = \sum_{J \in \mathcal{C}(I)} \rho_J \|c_J - c_I\|^2.$$

Given the threshold  $\eta > 0$ , define the subtree

$$\mathcal{T}_\eta = \begin{cases} \{\mathcal{X}\} & \text{if } \epsilon_I < \eta \quad \forall I \in \mathbb{T} \\ \{I \in \mathbb{T} \mid \exists J \in \mathbb{T} \text{ such that } J \subset I \text{ and } \epsilon_J \geq \eta\} & \text{otherwise} \end{cases}, \quad (17)$$

and let  $\Lambda_\eta = \Lambda(\mathcal{T}_\eta)$  be the corresponding outerleaves. Lemma 14 shows that  $\mathcal{T}_\eta$  is finite, so that by Lemma 11  $\Lambda_\eta$  is a partition and the corresponding nonlinear projection is

$$P_{\Lambda_\eta}(x) = \sum_{I \in \Lambda_\eta} c_I \mathbb{1}_I(x) \quad (18)$$

so that the code vectors are the centers of mass of the cells.

Comparing the definition of  $\mathcal{T}_\eta$  and  $\hat{\mathcal{T}}_\eta$ , we observe that  $\hat{\mathcal{T}}_\eta$  is truncated at the depth  $j_n$  given by (6), whereas  $\mathcal{T}_\eta$  is not truncated, but its maximal depth is bounded by Lemma 17.

Given the above definitions, we have the following result.

**Proposition 9.** *Given  $\eta > 0$ , for all  $0 < \sigma < s$*

$$\mathcal{E}(P_{\Lambda_\eta}) \lesssim \eta^{\frac{4\sigma}{2\sigma+1}}. \quad (19)$$

Note that the bound is meaningful only if  $0 < \eta < 1$ . Indeed for  $\eta \geq 1$   $\Lambda_\eta = \{\mathcal{X}\}$  and  $\mathcal{E}(P_{\Lambda_\eta}) \leq 1$ , see Remark 15.

## 5.1 Approximation Error

The quantity  $\mathcal{E}(P_{\Lambda_\eta})$  is called approximation error, by analogy with the corresponding definition in statistical learning theory, and it plays a special role in our analysis.

The problem of approximating a probability measure with a cloud of points is related to the so called optimal quantization [11]. The cost of an optimal quantizer is defined as:

$$V_{N,p}(\rho) := \inf_{S \subset \mathcal{X}, |S|=N} \mathbb{E}[d(X, S)^p],$$

where  $d(x, S) = \min_{y \in S} \|x - y\|$ . An optimal quantizer corresponds to a set  $S$  of  $N$  points attaining the infimum, with the corresponding Voronoi-Dirichlet partition of  $\mathcal{X}$ . One can interpret the approximation error  $\mathcal{E}(P_{\Lambda_\eta})$  as the quantization cost associated with the (suboptimal) quantizer given by the partition  $\Lambda_\eta$  as defined in 17 with the corresponding centers  $\{c_I\}_{I \in \Lambda_\eta}$ , and  $N := \#\Lambda_\eta$ .

This point of view is also taken through the analysis of  $k$ -means given in [6], optimal quantizers corresponding in fact to absolute minimizers of the  $k$ -means problem. Asymptotic estimates for the optimal quantization cost are available, see [6] and references therein. In the case of  $\text{supp}(\rho) = \mathcal{M}$ , being  $\mathcal{M}$  a smooth  $d$ -dimensional manifold isometrically emdedded in  $\mathcal{X}$ , they read:

$$\lim_{N \rightarrow \infty} N^{2/d} V_{N,2}(\rho) = C(\rho), \quad (20)$$

where  $C(\rho)$  is a constant depending only on  $\rho$ . We underline that the result provided by Proposition 9 is actually a non-asymptotic estimate for the quantization cost, when the quantizer is given by the outcome of our algorithm. The quantization cost is strictly higher than the optimal one, since, for instance, an optimal quantizer always corresponds to a Voronoi-Dirichlet partition [11]. Nevertheless, as observed in Section 3.3, a Voronoi quantizer is not suitable for multiscale refinements, whereas ours is. Proposition 9 does not directly compare with (20), as it depends on a different parameter quantifying the complexity of the partition, namely  $\eta$  instead of  $N$ . Though, by carefully applying (39b), in the manifold case we get:

$$\mathcal{E}(P_{\Lambda_\eta}) \lesssim \left( \frac{\log N}{N} \right)^{\frac{2}{d}}$$

so that the bound is in fact optimal up to a logarithmic factor. Furthermore, it is in order to observe that Assumption 5 together with Proposition 9 provide a more transparent understanding of the

approximation part of the analysis, as compared to what is provided in [4] and [14]. Therein, the approximation error is essentially addressed by defining the class of probability measures  $\mathcal{B}_s$  as those for which a certain approximation property holds; see Definition 5 in [14] and Definition 5 in [4], in both being the thresholding algorithm explicitly used. On the other hand Assumption 5 does not depend on the thresholding algorithm, but only on the mutual regularity of  $\rho$  and  $\mathbb{T}$ . Lastly we notice that, while for the sake of clarity none of the constants appear explicitly in our results, the proofs allow in principle to estimate them.

## 6 Proofs

In this section we collect some of the proofs of the above results. The more technical proofs are postponed to the appendix.

*Proof of Thm. 6.* We first observe that it is enough to show that there exists a partition tree  $\mathbb{T}' = \{\Lambda_j\}$  for  $\mathcal{M}$ . Indeed, by adding to each partion  $\Lambda_j$ , the cell  $I_0 = \mathcal{X} \setminus \mathcal{M}$ , we get a partion of  $\mathcal{X}$ , which satisfies Assumptions 5, since  $\text{diam}_\rho(I_0) = 0$ .

Since  $\mathcal{X}$  is bounded, then  $\mathcal{M}$  is a connected compact manifold and, hence,  $(\mathcal{M}, d_{\mathcal{M}}, \rho_{\mathcal{M}})$  is an Ahlfors regular metric measure space [10, p. 413], *i.e.*

$$d_1 r^d \leq \rho_{\mathcal{M}}(B_{\mathcal{M}}(x, r)) \leq d_2 r^d \quad r \leq \text{diam}(\mathcal{M}),$$

where  $B_{\mathcal{M}}(x, r)$  is the ball of center  $x$  and radius  $r$  with respect to the Riemannian metric  $d_{\mathcal{M}}$ . By (13b)

$$d_1 p_1 r^d \leq \rho(B_{\mathcal{M}}(x, r)) \leq d_2 p_2 r^d \quad r \leq \text{diam}(\mathcal{M}), \quad (21)$$

where  $d$  is the intrinsic dimension of  $\mathcal{M}$ . Since  $(\mathcal{M}, d_{\mathcal{M}}, \rho)$  is an Ahlfors regular metric measure, too, there exists a family of dyadic cubes, *i.e.* for each  $j \in \mathbb{Z}$  there is a family  $\Lambda_j = \{I\}$  of open subsets of  $\mathcal{M}$  such that

$$\rho(\mathcal{M} \setminus \cup_{I \in \Lambda_j} I) = 0 \quad (22a)$$

$$I \cap J = \emptyset \quad I, J \in \Lambda_j, I \neq J \quad (22b)$$

$$\text{either } I \cap J = \emptyset \text{ or } J \subset I \quad I \in \Lambda_j, J \in \Lambda_{j+\ell} \quad (22c)$$

$$I \supset B_{\mathcal{M}}(x_I, r_0 \delta^j) \quad I \in \Lambda_j \quad (22d)$$

$$I \subset B_{\mathcal{M}}(x_I, r_1 \delta^j) \quad I \in \Lambda_j \quad (22e)$$

where  $0 < r_0 < r_1$  and  $\delta \in (0, 1)$  are given constants [7, Thm. 11]. As noted in [9], it is always possible to redefine each cell  $I \in \Lambda_j$  by adding a suitable portion of its boundary in such a way that

$$\mathcal{M} = \cup_{I \in \Lambda_j} I \quad (23)$$

and (13a)–(22e) still hold true, possibly with different constants. Since  $\mathcal{M}$  is compact, there exists  $j_0 \in \mathbb{Z}$  such that  $B_{\mathcal{M}}(x_0, r_1 \delta^{j_0}) = \mathcal{M}$  for some  $x_0 \in \mathcal{M}$ . Hence, possibly redefining  $j$ ,  $r_0$  and  $r_1$ , we can assume that  $\Lambda_0 = \{\mathcal{M}\}$  and, as a consequence of (22b), (22c) and (23), the family  $\{\Lambda_j\}_{j \in \mathbb{N}}$

is a partition tree for  $\mathcal{M}$  where the bound in (3) is a consequence of the following standard volume argument. Fix  $j_0$  large enough such that for all  $j \geq j_0$ ,  $r_1 \delta^j \leq \text{diam } \mathcal{M}$ , then given  $I \in \Lambda_j$

$$\rho(I) = \sum_{J \in \mathcal{C}(I)} \rho(J) \geq \sum_{J \in \mathcal{C}(I)} \rho(B_{\mathcal{M}}(x_J, r_0 \delta^{j+1})) \geq \#\mathcal{C}(I) d_1 p_1 r_0^d \delta^{d(j+1)},$$

where the third and the fourth inequalities are consequence of (22d) and (21).

On the other hand, by (22e) and (21),

$$\rho(I) \leq \rho(B_{\mathcal{M}}(x_I, r_1 \delta^j)) \leq d_2 p_2 r_1^d \delta^{jd},$$

so that

$$\#\mathcal{C}(I) \leq \frac{d_2 p_2 r_1^d}{d_1 p_1 r_0^d \delta^d} = D.$$

Bound (3) holds true by setting

$$a = \max\{\max_{\substack{j < j_0 \\ I \in \Lambda_j}} \{\#\mathcal{C}(I)\}, D\}.$$

We now show that (10b) holds true. Indeed, since  $\mathcal{M}$  is Riemannian submanifold of  $\mathbb{R}^D$  it holds that

$$\|y - x\| \leq d_{\mathcal{M}}(y, x) \quad x, y \in \mathcal{M}, \quad (24)$$

see [22, Cor. 2, Prop. 21, Chapter 5]. Given  $I \in \Lambda_j$ , by (22e),

$$\text{diam}_{\rho}(I) \leq \text{diam}(I) \leq \sup_{x, y \in I} \|x - y\| \leq \sup_{x, y \in B_{\mathcal{M}}(x_I, r_1 \delta^j)} d_{\mathcal{M}}(x, y) \leq 2r_1 \delta^j,$$

so that (10b) holds true with  $b = 1/\delta > 1$  and  $C_2 = 2r_1$ . To show (10a), given  $I \in \Lambda_j$ , by (22d) and (21)

$$\rho(I) \geq \rho(B_{\mathcal{M}}(x_I, r_0 \delta^j)) \geq d_1 p_1 \delta^{jd}.$$

Hence

$$\text{diam}_{\rho}(I) \leq 2r_1 \delta^j \leq 2r_1 \left( \frac{\rho(I)}{d_1 p_1} \right)^{\frac{1}{d}},$$

so that (10a) holds true with  $C_1 = 2r_1 (d_1 p_1)^{-\frac{1}{d}}$  and  $s = 1/d$ .  $\square$

The proof of Theorem 7 borrows ideas from [4, 14] and combines a number of intermediate results given in Appendix A. For sake of clarity let  $\hat{P}_{\eta} = \hat{P}_{\Lambda_{\eta}}$ .

*Proof of Thm. 7.* Consider the following decomposition

$$\begin{aligned} x - \hat{P}_{\Lambda_{\eta}}(x) &= \left( x - P_{\Lambda(\hat{\mathcal{T}}_{\eta} \cup \mathcal{T}_{2\eta})(x)} \right) + \left( P_{\Lambda(\hat{\mathcal{T}}_{\eta} \cup \mathcal{T}_{2\eta})(x)} - P_{\Lambda(\hat{\mathcal{T}}_{\eta} \cap \mathcal{T}_{\eta/2})(x)} \right) + \\ &+ \left( P_{\Lambda(\hat{\mathcal{T}}_{\eta} \cap \mathcal{T}_{\eta/2})(x)} - \hat{P}_{\Lambda(\hat{\mathcal{T}}_{\eta} \cap \mathcal{T}_{2\eta})(x)} \right) + \left( \hat{P}_{\Lambda(\hat{\mathcal{T}}_{\eta} \cap \mathcal{T}_{2\eta})(x)} - \hat{P}_{\Lambda(\hat{\mathcal{T}}_{\eta})(x)} \right), \end{aligned}$$

which holds for all  $x \in \mathcal{X}$ . Since

$$\left\| \sum_{i=1}^4 v_i \right\|^2 \leq 4 \sum_{i=1}^4 \|v_i\|^2 \quad v_1, \dots, v_4 \in \mathbb{R}^D,$$

it holds that

$$\begin{aligned} \mathcal{E}[\widehat{P}_{\widehat{\Lambda}_\eta}] &\lesssim \underbrace{\mathcal{E}[P_{\Lambda(\widehat{\mathcal{T}}_\eta \cup \mathcal{T}_{2\eta})}]}_A + \underbrace{\int_{\mathcal{X}} \left\| P_{\Lambda(\widehat{\mathcal{T}}_\eta \cup \mathcal{T}_{2\eta})}(x) - P_{\Lambda(\widehat{\mathcal{T}}_\eta \cap \mathcal{T}_{\eta/2})}(x) \right\|^2 d\rho(x)}_B \\ &\quad + \underbrace{\int_{\mathcal{X}} \left\| P_{\Lambda(\widehat{\mathcal{T}}_\eta \cap \mathcal{T}_{\eta/2})}(x) - \widehat{P}_{\Lambda(\widehat{\mathcal{T}}_\eta \cap \mathcal{T}_{\eta/2})}(x) \right\|^2 d\rho(x)}_C + \underbrace{\int_{\mathcal{X}} \left\| \widehat{P}_{\Lambda(\widehat{\mathcal{T}}_\eta \cap \mathcal{T}_{2\eta})}(x) - \widehat{P}_{\Lambda(\widehat{\mathcal{T}}_\eta)}(x) \right\|^2 d\rho(x)}_D. \end{aligned}$$

We bound the four terms.

A) Since  $\widehat{\mathcal{T}}_\eta \cup \mathcal{T}_{2\eta} \supset \mathcal{T}_{2\eta}$ ,  $\Lambda(\widehat{\mathcal{T}}_\eta \cup \mathcal{T}_{2\eta})$  is a partition finer than  $\Lambda_{2\eta}$ , then

$$\mathcal{E}[P_{\Lambda(\widehat{\mathcal{T}}_\eta \cup \mathcal{T}_{2\eta})}] \leq \mathcal{E}[P_{\Lambda_{2\eta}}] \lesssim \eta^{\frac{4\sigma}{2\sigma+1}},$$

where the last inequality is a consequence of (19).

B) Bound (50a) implies that the term  $B$  is zero with probability greater than  $1 - p_B$ , where

$$p_B \lesssim (n^\gamma + \eta^{-\frac{2}{2\sigma+1}}) \exp(-c_a n \eta^2).$$

C) Since  $\Lambda(\widehat{\mathcal{T}}_\eta \cap \mathcal{T}_{\eta/2}) \subset \mathcal{T}_{\eta/2} \cup \Lambda_{\eta/2} = \mathcal{I}$  and  $\#\Lambda(\widehat{\mathcal{T}}_\eta \cap \mathcal{T}_{\eta/2}) \leq \#\Lambda_{\eta/2} = N$ , by (41) term C is bounded by  $t^* = \eta^{\frac{4\sigma}{2\sigma+1}} t$  with probability greater than  $1 - p_C$  with

$$p_C = 2\#\mathcal{I} \exp\left(-\frac{nt^*}{4N}\right) \lesssim \eta^{-\frac{2}{2\sigma+1}} \exp(-c_{\mathbb{T}} n \eta^2 t) \quad (25)$$

where the second inequality is a consequence of (39a) and (39b), and  $c_{\mathbb{T}} > 0$  is a suitable constant depending on the partition tree  $\mathbb{T}$ .

D) By (50b) term D is zero with probability greater than  $1 - p_D$  where

$$p_D \lesssim \eta^{-\frac{2}{2\sigma+1}} \exp(-c_a n \eta^2).$$

It follows that with probability greater than  $1 - (p_B + p_C + p_D)$

$$\mathcal{E}[\widehat{P}_{\widehat{\Lambda}_\eta}] \lesssim \underbrace{\eta^{\frac{4\sigma}{2\sigma+1}}}_A + \underbrace{\eta^{\frac{4\sigma}{2\sigma+1}} t}_C$$

*i.e.*

$$\mathbb{P} \left[ \mathcal{E}[\widehat{P}_{\widehat{\Lambda}_\eta}] \gtrsim \eta^{\frac{4\sigma}{2\sigma+1}} (1+t) \right] \lesssim \underbrace{(n^\gamma + \eta^{-\frac{2}{2\sigma+1}}) \exp(-c_a n \eta^2 t)}_{p_A + p_D} + \underbrace{\eta^{-\frac{2}{2\sigma+1}} \exp(-c_{\mathbb{T}} n \eta^2 t)}_{p_C}$$

which gives (14).  $\square$

*Proof of Cor. 8.* Since  $\eta_n^2 = \frac{(\gamma+\beta)\ln n}{c_a n}$ , then bound (14) gives (16a) since

$$\begin{aligned} (n^\gamma + \left(\frac{c_a n}{(\gamma+\beta)\ln n}\right)^{\frac{1}{2\sigma+1}}) \exp(-c_a n \eta_n^2) &\lesssim n^\gamma n^{-(\gamma+\beta)} = n^{-\beta} \\ \left(\frac{c_a n}{(\gamma+\beta)\ln n}\right)^{\frac{1}{2\sigma+1}} \exp(-c_{\mathbb{T}} n \eta_n^2 t) &\lesssim n n^{-\bar{c}_{\mathbb{T}} t} = n^{1-\bar{c}_{\mathbb{T}} t} \end{aligned}$$

where  $\bar{c}_{\mathbb{T}} = c_{\mathbb{T}} c_a^{-1} (\gamma + \beta)$ . Eq. (16b) is clear.  $\square$

*Proof of Prop. 9.* Given  $I \in \mathbb{T}$ , by (34) and (10b),

$$\sum_{J \in \mathcal{C}^{N+1}(I)} \mathcal{E}_J \lesssim a b^{-2(j_I + N+1)}$$

so that

$$\lim_{N \rightarrow +\infty} \sum_{J \in \mathcal{C}^{N+1}(I)} \mathcal{E}_J = 0$$

and, by taking the limit in (32),

$$\mathcal{E}_I = \sum_{k=0}^{+\infty} \sum_{J \in \mathcal{C}^k(I)} \epsilon_J^2.$$

Set  $\mathcal{T}_k = \mathcal{T}_{\eta/2^k}$  for all  $k \in \mathbb{N}$ , then

$$\begin{aligned} \mathcal{E}(P_{\Lambda_\eta}) &= \sum_{I \in \Lambda_\eta} \mathcal{E}_I = \sum_{I \in \Lambda_\eta} \sum_{k=0}^{+\infty} \sum_{J \in \mathcal{C}^k(I)} \epsilon_J^2 = \sum_{J \notin \mathcal{T}_\eta} \epsilon_J^2 = \sum_{k=0}^{+\infty} \sum_{J \in \mathcal{T}_{k+1} \setminus \mathcal{T}_k} \epsilon_J^2 \leq \sum_{k=0}^{+\infty} \#\mathcal{T}_{k+1} \left(\frac{\eta}{2^k}\right)^2 \\ &\lesssim \sum_{k=0}^{+\infty} \frac{\eta^2}{2^{2k}} \left(\frac{\eta}{2^{k+1}}\right)^{-\frac{2}{2\sigma+1}} = \eta^{\frac{4\sigma}{2\sigma+1}} \sum_{k=0}^{+\infty} 4^{\frac{k+1}{2\sigma+1} - k} \lesssim \eta^{\frac{4\sigma}{2\sigma+1}}, \end{aligned}$$

where the first inequality is a consequence of the fact that  $\epsilon_J < (\frac{\eta}{2^k})^2$  if  $J \notin \mathcal{T}_k$ , the second inequality follows from (39a), whereas the last inequality holds since the series  $\sum_{k=0}^{+\infty} 4^{-\frac{2\sigma k - 1}{2\sigma+1}}$  converges.  $\square$

## 7 Conclusions

In this paper, we proposed and analyzed a multiscale vector quantization approach inspired by ideas in [14]. We provided non asymptotic error bounds on the corresponding distortion/reconstruction error combining geometric and probabilistic tools. The analysis is developed in a general setting with manifold supported data as a special case.

A number of research directions remain to be explored. Following again [14], it would be interesting to understand the role of noise and the impact of data driven partitioning. Further, following the parallel with wavelets it would be interesting to use results from reproducing kernel Hilbert spaces to develop Gabor like geometric wavelets and analyse the properties of the corresponding quantization schemes.

## Acknowledgments

L. R. acknowledges the financial support of the AFOSR projects FA9550-17-1-0390 and BAA-AFRL-AFOSR-2016-0007 (European Office of Aerospace Research and Development), and the EU H2020-MSCA-RISE project NoMADS - DLV-777826.

## A Further Proofs

This section is devoted to show all the results cited in the proof in Section 6. We first show some preliminary results.

### A.1 Technical results

We introduce some further notations about partition trees.

With slight abuse of notation, we regard  $\mathbb{T} = \cup_{j \in \mathbb{N}} \cup_{I \in \Lambda_j} I$  as the (disjoint) union of the cells in each partition  $\Lambda_j$ . If a cell does not split, *i.e.*  $\mathcal{C}(I) = I$ , we regard  $I \in \Lambda_j$  and  $\mathcal{C}(I) \in \Lambda_{j+1}$  as different cells.

Given a cell  $I \in \mathbb{T}$ , for any  $k \in \mathbb{N}$  we set

$$\mathcal{C}^{k+1}(I) = \mathcal{C}(\mathcal{C}^k(I)) \quad \mathcal{P}^{k+1}(I) = \mathcal{P}(\mathcal{P}^k(I)),$$

where  $\mathcal{C}^0(I) = \mathcal{P}^0(I) = \{I\}$  and, clearly,

$$\mathcal{P}^k(I) = \{\mathcal{X}\} \quad k \geq j_I.$$

Furthermore, for any pair  $I, J \in \mathbb{T}$ ,  $I \neq J$  one and only one of the following alternative possibilities holds true

$$I \cap J = \emptyset \quad \text{or} \quad J \in \mathcal{C}^k(I) \quad \text{or} \quad I \in \mathcal{C}^k(J) \quad (26)$$

for some  $k \geq 1$ .

If  $\{\mathcal{T}_t\}_{t \in T}$  is an arbitrary family of subtrees, clearly the intersection  $\cap_{t \in T} \mathcal{T}_t$  and the union  $\cup_t \mathcal{T}_t$  are the *smallest* and the *largest* subtrees in the family.

Given a subset  $S \subset \mathbb{T}$ , we set

$$\mathcal{T}(S) = \bigcap \{\mathcal{T} \mid \mathcal{T} \text{ is a subtree and } S \subset \mathcal{T}\},$$

which is the smallest subtree containing all the cells in  $S$ , and

$$\mathcal{T}(S) = \bigcup_{I \in S} \{\mathcal{P}^k(I) \mid k = 0, \dots, j_I\} = \bigcup_{I \in S} \mathcal{T}(I) \quad (27a)$$

$$\#\mathcal{T}(S) \leq 1 + \sum_{I \in S} j_I. \quad (27b)$$

The following remark provides an alternative definition of  $\hat{\mathcal{T}}_\eta$  and a similar procedure can be applied to  $\mathcal{T}_\eta$

**Remark 10.** Given  $\eta > 0$  and  $j_n \in \mathbb{N}$ , set

$$\hat{S}_\eta = \{I \in \mathbb{T} \mid j_I \geq j_n, \hat{\epsilon}_I \geq \eta\}$$

which is a finite subset of  $\mathbb{T}$ . It is easy to check that  $\mathcal{T}(\hat{S}_\eta) = \hat{T}_\eta$ .

The following lemma is quite obvious.

**Lemma 11.** If  $\mathcal{T}$  is finite subtree, the family of outer leaves

$$\Lambda(\mathcal{T}) = \{I \in \mathbb{T} : I \notin \mathcal{T}, \mathcal{P}(I) \in \mathcal{T}\}$$

is a partition with

$$\#\Lambda(\mathcal{T}) \leq (a-1)\#\mathcal{T} + 1 \leq a\#\mathcal{T}. \quad (28)$$

*Proof.* For any  $j \in \mathbb{N}$ , we fix an arbitrary order among the cells in  $\Lambda_j$ , i.e.

$$\Lambda_j = \{I_{j,1}, \dots, I_{j,N_j}\} \quad N_j = \#\Lambda_j.$$

Given two different cells  $I = I_{j,k}$  and  $J = I_{j',k'}$  in  $\mathbb{T}$ , we say that  $I$  is older than  $J$  either if  $j < j'$  or  $j = j'$  and  $k < k'$ .

By definition  $\mathcal{T}$  contains the parents of all its elements and, hence, the root  $\mathcal{X}$ . Define, by induction, the family of subtrees

$$\mathcal{T}_1 = \{\mathcal{X}\} \quad \mathcal{T}_{n+1} = \mathcal{T}_n \cup \{I\},$$

where  $I$  is the oldest cell in  $\mathcal{T} \setminus \mathcal{T}_n$ . Note that, by construction,  $I \in \Lambda(\mathcal{T}_n)$ . Since  $\mathcal{T}$  is finite by assumption,  $\mathcal{T}_N = \mathcal{T}$  with  $N = \#\mathcal{T}$ .

We now prove by induction on  $n = 1, \dots, N$  that  $\Lambda(\mathcal{T}_n)$  is a partition. If  $n = 1$

$$\Lambda(\mathcal{T}_1) = \{I \in \mathbb{T} \mid I \neq \mathcal{X}, \mathcal{P}(I) = \mathcal{X}\} = \mathcal{C}(\mathcal{X}) = \Lambda_1$$

which is a partition by assumption and, by (3), it satisfies (28). Assume that  $\Lambda(\mathcal{T}_n)$  is partition satisfying (28) with  $\mathcal{T} = \mathcal{T}_n$ . By construction,  $\mathcal{T}_{n+1} = \mathcal{T}_n \cup I$  with  $I \in \Lambda(\mathcal{T}_n)$ , then

$$\Lambda(\mathcal{T}_{n+1}) = \left( \bigcup_{J \in \mathcal{C}(I)} \{J\} \right) \cup (\Lambda(\mathcal{T}_n) \setminus \{I\})$$

which is a partition since  $I = \cup_{J \in \mathcal{C}(I)} J$  and

$$\#\Lambda(\mathcal{T}_{n+1}) \leq a + \#\Lambda(\mathcal{T}_n) - 1 \leq a + (a-1)\#\mathcal{T}_n = (a-1)(\#\mathcal{T}_n + 1) + 1 = (a-1)(\#\mathcal{T}_{n+1}) + 1,$$

so that (28) holds true with  $\mathcal{T} = \mathcal{T}_{n+1}$ .  $\square$

We observe that, given a cell  $I$ , by definition of  $\text{diam}_\rho(I)$ , it exists a measurable subset  $J \subset I$  such that  $\rho(J) = 0$  and  $\text{diam}(I \setminus J) = \text{diam}_\rho(I)$ . Furthermore,

$$\mathbb{P}[\|X_i - X_j\| > \text{diam}_\rho(I) \mid X_i, X_j \in I] = 0 \quad i, j = 1, \dots, n.$$

The following simple lemmas will be useful.

**Lemma 12.** *Given a cell  $I \in \mathbb{T}$  with  $\rho(I) > 0$*

$$\|x - c_I\| \leq \text{diam}_\rho(I) \quad \rho - \text{almost all } x \in I \quad (29a)$$

$$\|c_J - c_K\| \leq \text{diam}_\rho(I) \quad J, K \in \mathcal{C}(I) \quad (29b)$$

$$\|\widehat{c}_I - c_I\| \leq \begin{cases} \text{diam}_\rho(I) & n_I \neq 0 \\ \text{diam}(\mathcal{X}) & n_I = 0 \end{cases} \quad \text{almost surely} \quad (29c)$$

*Proof.* By definition of essential diameter, there exists  $I_0 \subset I$  such that  $\text{diam}(I_0) = \text{diam}_\rho(I)$  and  $\rho(I \setminus I_0) = 0$ . Let  $C$  the closed convex hull of  $I_0$ . It is known that  $\text{diam}(C) = \text{diam}_\rho(I)$  and, by convexity theorem, see [27, Thm. 5.7.35],

$$c_I = \frac{1}{\rho(I_0)} \int_{I_0} x d\rho(x) \in C,$$

so that (29a) is clear. Since  $J, K \subset I$ , Eq. (29b) is a consequence of the fact that  $c_J, c_K \in C$ . If  $n_I = 0$ ,  $\widehat{c}_I = \widehat{x}_I^* \in \mathcal{X}$  so that (29c) is clear. If  $n_I \neq 0$ , almost surely  $\widehat{c}_I \in C$  so that

$$\|\widehat{c}_I - c_I\| \leq \text{diam}(C) = \text{diam}(I_0) = \text{diam}_\rho(I).$$

□

Given a cell  $I \in \mathbb{T}$ , the within-between decomposition of the variance

$$\mathcal{E}_I = \sum_{J \in \mathcal{C}(I)} \mathcal{E}_J + \sum_{J \in \mathcal{C}(I)} \rho_I \|c_J - c_I\|^2, \quad (30)$$

implies

$$\epsilon_I^2 = \sum_{J \in \mathcal{C}(I)} \rho_I \|c_J - c_I\|^2 \quad (31)$$

As a consequence we have the following decomposition.

**Lemma 13.** *Given  $I \in \mathbb{T}$ , for all  $N \in \mathbb{N}$*

$$\mathcal{E}_I = \sum_{k=0}^N \sum_{J \in \mathcal{C}^k(I)} \epsilon_J^2 + \sum_{J \in \mathcal{C}^{N+1}(I)} \mathcal{E}_J. \quad (32)$$

*Proof.* The claim is clear for  $N = 0$ . Assume that it holds true for  $N$ . Then, for any  $J \in \mathcal{C}^{N+1}(I)$

$$\mathcal{E}_J = \epsilon_J^2 + \sum_{J' \in \mathcal{C}(J)} \mathcal{E}_{J'},$$

hence

$$\mathcal{E}_I = \sum_{k=0}^N \sum_{J \in \mathcal{C}^k(I)} \epsilon_J^2 + \sum_{J \in \mathcal{C}^{N+1}(I)} \left( \epsilon_J^2 + \sum_{J' \in \mathcal{C}(J)} \mathcal{E}_{J'} \right),$$

by observing that a cell  $J' \in \mathcal{C}(J)$  for some  $J \in \mathcal{C}^{N+1}(I)$  if and only if  $J' \in \mathcal{C}^{N+2}(I)$ , so that (32) holds true for  $N + 1$ . □

We now show that the set  $\mathcal{T}_\eta$  defined by (17) is a finite subtree.

**Lemma 14.** *The family  $\mathcal{T}_\eta$  is a finite subtree of  $\mathbb{T}$ .*

*Proof.* If  $\mathcal{T}_\eta = \{\mathcal{X}\}$ , there is nothing to prove. Otherwise, if  $I \in \mathcal{T}_\eta$ , then by definition there exists  $J \in \mathbb{T}$  such that  $J \subset I$  and  $\epsilon_J \geq \eta$ . Since  $\mathcal{P}(I) \supset I \supset J$ , then  $P(I) \in \mathcal{T}_\eta$ , so that  $\mathcal{T}_\eta$  is a subtree.

We now show that  $\mathcal{T}_\eta$  is finite. From (32) with  $I = \mathcal{X}$  and (34) we get that for all  $N \in \mathbb{N}$

$$\sum_{k=0}^N \sum_{J \in \mathcal{C}^k(\mathcal{X})} \epsilon_J^2 \leq \mathcal{E}_\mathcal{X} \leq 1.$$

Then the series  $\sum_{k=0}^{+\infty} \sum_{J \in \mathcal{C}^k(\mathcal{X})} \epsilon_J^2 = \sum_{I \in \mathbb{T}} \epsilon_I^2$  is sommable. Hence, the set

$$S_\eta = \{I \in \mathbb{T} : \epsilon_I \geq \eta, \}$$

is finite. Furthermore, by construction

$$\mathcal{T}_\eta = \mathcal{T}(S_\eta) = \bigcup_{I \in S_\eta} \mathcal{T}(I). \quad (33)$$

Bound (27b) implies that  $\mathcal{T}_\eta$  is finite. □

**Remark 15.** *By (9),*

$$\epsilon_I^2 \leq \mathcal{E}_I \leq \text{diam}_\rho(I)^2 \rho_I \leq \text{diam}(\mathcal{X})^2 \leq 1, \quad (34)$$

and  $\epsilon_I < \mathcal{E}_I \leq 1$  provided that  $\epsilon_I > 0$ . Furthermore, by (34)  $\mathcal{T}_\eta = \Lambda_\eta = \{\mathcal{X}\}$  for all  $\eta \geq 1$  and  $\mathcal{E}(P_{\Lambda_\eta}) = \mathcal{E}_\mathcal{X} \leq 1$ . By the same argument  $\hat{\mathcal{T}}_\eta = \hat{\Lambda}_\eta = \{\mathcal{X}\}$ . Hence, it is enough to consider the case  $0 < \eta < 1$ .

## A.2 A-term: Approximation error

In this section, we bound the approximation error, which is based in an estimation of the number of cells  $I$  such that  $\epsilon_I$  is big enough.

**Lemma 16.** *Given a partition  $\Lambda \subset \mathbb{T}$ , given  $0 < \eta < 1$ ,*

$$\#\{I \in \Lambda : \epsilon_I \geq \eta\} \leq \#\{I \in \Lambda : \mathcal{E}_I \geq \eta^2\} \lesssim \eta^{-\frac{2}{2s+1}}. \quad (35)$$

*Proof.* First inequality in (35) is a direct consequence of (34). By (10a) there exists  $C > 0$  such that for all  $I \in \mathbb{T}$

$$\text{diam}_\rho(I) \leq C \rho_I^s \quad I \in \Lambda,$$

then, by (34)

$$\mathcal{E}_I \leq C^2 \rho_I^{2s+1}.$$

Set  $\Lambda_+ = \{I \in \Lambda : \mathcal{E}_I \geq \eta^2\}$  and  $N_+ = \#\Lambda_+$ . Fix  $q \geq 1$ , clearly

$$\eta^{\frac{2}{q}} N_+ \leq \sum_{I \in \Lambda_+} \mathcal{E}_I^{\frac{1}{q}} \leq C^{\frac{2}{q}} \sum_{I \in \Lambda_+} \rho_I^{\frac{2s+1}{q}}. \quad (36)$$

The Holder inequality with  $1/p + 1/q = 1$  gives

$$\sum_{I \in \Lambda_+} \rho_I^{\frac{2s+1}{q}} \leq \left( \sum_{I \in \Lambda_+} \rho_I^{\frac{p(2s+1)}{q}} \right)^{\frac{1}{p}} \left( \sum_{I \in \Lambda_+} 1^q \right)^{\frac{1}{q}} \leq \left( \sum_{I \in \Lambda} \rho_I \right)^{\frac{2s+1}{2s+2}} N_+^{\frac{1}{2s+2}} = N_+^{\frac{1}{2s+2}},$$

where the last inequality follows by choosing  $\frac{p(2s+1)}{q} = 1$ , *i.e.*  $q = 2s + 2$ . By replacing in (36)

$$N_+^{1 - \frac{1}{2s+2}} \leq C^{\frac{1}{s+1}} \eta^{-\frac{2}{2s+2}}$$

and, since  $1 - \frac{1}{2s+2} = \frac{2s+1}{2s+2}$ , we get (35).  $\square$

We first observe that the proof of the above lemma only depends on Assumption (10a) and the constant in the inequality (35) only depends on the constant in (10a), denoted by  $C$  in the proof. Furthermore, without Assumption (10a) we always have the following bound

$$\#\{I \in \Lambda : \mathcal{E}_I \geq \eta^2\} \leq \#\{I \in \Lambda : \mathcal{E}_I \geq \eta^2\} \leq \eta^{-2},$$

where the first inequality is consequence of (34) and the last bounds is due to (30) with  $I = \mathcal{X}$

$$1 \geq \mathcal{E}_{\mathcal{X}} \geq \sum_{I \in \Lambda} \mathcal{E}_I \geq \sum_{I \in \Lambda, \mathcal{E}_I \geq \eta^2} \mathcal{E}_I = \eta^2 \#\{I \in \Lambda : \mathcal{E}_I \geq \eta^2\}.$$

We recall that  $\Lambda_\eta = \Lambda(T_\eta)$  is the family of the corresponding outer leaves, which is a partition of  $\mathcal{X}$  by Lemma 11. In order to bound the cardinality of  $\Lambda_\eta$  we need an auxiliary lemma based on Assumption (10b).

**Lemma 17.** *Given  $\eta > 0$ , set*

$$j_\eta = \sup\{j_I \in \mathbb{N} \mid I \in \mathbb{T} \text{ and } \epsilon_I \geq \eta\}, \quad (37)$$

*then*

$$j_\eta \lesssim \ln\left(\frac{2}{\eta}\right). \quad (38)$$

If  $\{j_I \in \mathbb{N} \mid I \in \mathbb{T} \text{ and } \epsilon_I \geq \eta\} = \emptyset$  we set  $j_\eta = 0$ .

*Proof.* If  $\mathcal{T}_\eta = \{\mathcal{X}\}$  or  $\{j_I \in \mathbb{N} \mid I \in \mathbb{T} \text{ and } \epsilon_I \geq \eta\} = \emptyset$ , then  $j_\eta = 0$ , so that the claim is evident. If  $\mathcal{T}_\eta \neq \{\mathcal{X}\}$ , then  $0 < \eta < 1$ . Take  $I \in \mathbb{T}$  such that  $\epsilon_I \geq \eta$ . By (34) and (10b)

$$\eta^2 \leq \epsilon_I^2 \leq \rho_I \text{diam}_\rho(I)^2 \leq C_2 b^{-2j_I}.$$

Hence

$$j_I \leq \frac{1}{\ln b} \ln\left(\frac{1}{\eta}\right) + \frac{1}{2 \ln b} \ln C_2 \leq E \ln\left(\frac{2}{\eta}\right),$$

where  $E = \max\{1, \frac{\ln C_2}{2 \ln 2}\} / \ln b$ .  $\square$

**Proposition 18.** *Given  $\eta > 0$ , for all  $\sigma < s$ ,*

$$\#\mathcal{T}_\eta \lesssim \eta^{-\frac{2}{2s+1}} \ln\left(\frac{2}{\eta}\right) \lesssim \eta^{-\frac{2}{2\sigma+1}} \quad (39a)$$

$$\#\Lambda_\eta \lesssim \eta^{-\frac{2}{2s+1}} \ln\left(\frac{2}{\eta}\right) \lesssim \eta^{-\frac{2}{2\sigma+1}}, \quad (39b)$$

where the constants in  $\lesssim$  also depend on  $\sigma$ .

*Proof.* As observed in Remark 15, we can assume that  $0 < \eta < 1$ . Let

$$\Upsilon = \{I \in \mathcal{T}_\eta \mid \epsilon_I \geq \eta \text{ and } \epsilon_J < \eta \forall J \in \mathcal{C}^k(I), k \geq 1\}.$$

By (26) the elements of  $\Upsilon$  are disjoint. Let  $\Lambda \subset \mathbb{T}$  be a partition such that  $\Upsilon \subset \Lambda$ . Hence, by (35)

$$\#\Upsilon \leq \#\{I \in \Lambda \mid \epsilon_I \geq \eta\} \lesssim \eta^{-\frac{2}{2s+1}}. \quad (40)$$

We claim that

$$\mathcal{T}_\eta = \bigcup_{J \in \Upsilon} \mathcal{T}(J) = \mathcal{T}(\Upsilon).$$

By construction  $\mathcal{T}_\eta \supset \bigcup_{J \in \Upsilon} \mathcal{T}(J)$ . To prove the opposite inclusion, fix  $I \in \mathcal{T}_\eta$ , then there exists  $J_1 \in \mathbb{T}$  such that  $J_1 \subset I$  and  $\epsilon_{J_1} \geq \eta$ . If  $J_1 \in \Upsilon$ , then  $I \in \mathcal{T}(J_1)$ . Otherwise, there exists  $J_2 \subset J_1 \subset I$  and  $\epsilon_{J_2} \geq \eta$ . If  $J_2 \in \Upsilon$ ,  $I \in \mathcal{T}(J_2)$ . Otherwise, because of  $\mathcal{T}_\eta$  is finite, we can repeat the procedure until we get  $J_k \in \Upsilon$  such that  $J_k \subset I$ , then  $I \in \mathcal{T}(J_k)$  and the claim is proven. By (27b)

$$\#\mathcal{T}_\eta = \#\mathcal{T}(\Upsilon) \leq 1 + \sum_{J \in \Upsilon} j_J \leq 1 + \#\Upsilon j_\eta \lesssim \eta^{-\frac{2}{2s+1}} \ln\left(\frac{2}{\eta}\right),$$

since, by definition,  $j_J \leq j_\eta$  and the last inequality is due to (40) and (38). This shows the first inequality in (39a). Since  $\sigma < s$ , for some  $\delta > 0$

$$\eta^{-\frac{2}{2s+1}} \ln\left(\frac{2}{\eta}\right) = \eta^{-\frac{2}{2\sigma+1}} \eta^\delta \ln\left(\frac{2}{\eta}\right) \leq C \eta^{-\frac{2}{2\sigma+1}}$$

where  $C = \sup_{0 < \eta \leq 1} \eta^\delta \ln\left(\frac{2}{\eta}\right)$ , which is finite, since  $\lim_{\eta \rightarrow 0} \eta^\delta \ln\left(\frac{2}{\eta}\right) = 0$ . Bound (39b) is a direct consequence of (28).  $\square$

**Remark 19.** *In the following, for sake of clarity we bound the logarithmic dependence on  $\eta$  by considering  $\sigma < s$ . However our results can be extended to  $\sigma = s$  by adding a logarithmic factor, as in (39a) and (39b).*

### A.3 C-term: sample error

The following result bounds the sample error for a given partition.

**Proposition 20.** *Fix a data-independent subset  $\mathcal{I} \subset \mathbb{T}$  of cells and  $N > 0$ . Given a partition  $\widehat{\Lambda} \subset \mathcal{I}$  (possibly depending on the data) such that  $\#\widehat{\Lambda} \leq N$ , for any  $t > 0$ ,*

$$\mathbb{P} \left[ \int_{\mathcal{X}} \left\| P_{\widehat{\Lambda}}(x) - \widehat{P}_{\widehat{\Lambda}}(x) \right\|^2 d\rho(x) > t \right] \leq 2 \#\mathcal{I} \exp \left( -\frac{nt}{8N} \right). \quad (41)$$

*Proof.* Consider the following event

$$\Omega = \bigcup_{I \in \mathcal{I}} \{\sqrt{\rho_I} \|\widehat{c}_I - c_I\| > t\},$$

which is well-defined since  $\mathcal{I}$  does not depend on the data  $X_1, \dots, X_n$ . By union bound

$$\mathbb{P}[\Omega] \leq \#\mathcal{I} \sup_{\substack{I \in \mathcal{I} \\ \rho_I > 0}} \mathbb{P}[\|\widehat{c}_I - c_I\| > \frac{t}{\sqrt{\rho_I}}]. \quad (42)$$

Fix  $I \in \mathcal{I}$  with  $\rho_I > 0$ . The tower property with respect to the binomial random variable  $n_I$  gives

$$\mathbb{P}[\|\widehat{c}_I - c_I\| > t] = \sum_{k=0}^n \binom{n}{k} \rho_I^k (1 - \rho_I)^{n-k} \mathbb{P}[\|\widehat{c}_I - c_I\| > t \mid n_I = k].$$

Conditionally to the event  $\{n_I = k\}$  with  $k > 0$ , up to a permutation of the indexes, we can assume that  $X_1, \dots, X_k \in I$  and  $X_{k+1}, \dots, X_n \notin I$ . Furthermore,

$$\widehat{c}_I - c_I = \frac{1}{k} \sum_{i=1}^k (X_i - c_I) = \frac{1}{k} \sum_{i=1}^k \xi_i$$

where  $\xi_1, \dots, \xi_k$  are independent zero mean random vectors bounded by  $M = \text{diam}_\rho(I)$  almost surely by (29a). Hence, by (56)

$$\mathbb{P}[\|\widehat{c}_I - c_I\| > t \mid n_I = k] \leq 2 \exp\left(-\frac{kt^2}{4 \text{diam}_\rho(I)^2}\right),$$

which trivially holds true also if  $k = 0$ . Hence,

$$\begin{aligned} \mathbb{P}[\|\widehat{c}_I - c_I\| > t] &\leq 2 \sum_{k=0}^n \binom{n}{k} \left(\rho_I \exp\left(-\frac{t^2}{4 \text{diam}_\rho(I)^2}\right)\right)^k (1 - \rho_I)^{n-k} = 2 \left(1 - \rho_I \left(1 - \exp\left(-\frac{t^2}{4 \text{diam}_\rho(I)^2}\right)\right)\right)^n \\ &\leq 2 \exp\left(-n\rho_I \left(1 - \exp\left(-\frac{t^2}{4 \text{diam}_\rho(I)^2}\right)\right)\right) \end{aligned}$$

where in the fourth line we use the bound  $(1 - \tau)^n \leq \exp(-n\tau)$  with  $0 \leq \tau \leq 1$ . Since

$$1 - \exp(-\tau) \geq \frac{\tau}{2} \quad \text{for all } \tau \leq 1,$$

then, for all  $t \leq \text{diam}_\rho(I)$ ,

$$\mathbb{P}[\|\widehat{c}_I - c_I\| > t] \leq 2 \exp\left(-n\rho_I \frac{t^2}{8 \text{diam}_\rho(I)^2}\right) \leq 2 \exp\left(-n\rho_I \frac{t^2}{8}\right). \quad (43)$$

If  $\text{diam}_\rho(I) < t \leq \text{diam } \mathcal{X} \leq 1$ , by (29c), clearly  $\mathbb{P}[\|\widehat{c}_I - c_I\| > t \mid n_I > 0] = 0$ , so that

$$\mathbb{P}[\|\widehat{c}_I - c_I\| > t] = \mathbb{P}[\|\widehat{x}_I^* - c_I\| > t \mid n_I = 0] \mathbb{P}[n_I = 0] \leq (1 - \rho_I)^n \leq \exp(-n\rho_I) \leq 2 \exp\left(-n\rho_I \frac{t^2}{8}\right),$$

where the last bound holds true for any  $t \leq 2\sqrt{2}$ . Finally, if  $t > \text{diam } \mathcal{X}$ , as above

$$\mathbb{P}[\|\widehat{c}_I - c_I\| > t] = \mathbb{P}[\|\widehat{x}_I^* - c_I\| > t \mid n_I = 0] \mathbb{P}[n_I = 0] = 0$$

since  $\widehat{x}_I^*, c_I \in \mathcal{X}$ , compare with (29c). It follows that (43) holds true for all  $t > 0$ . Summarizing, from (42) we get

$$\mathbb{P}[\Omega] \leq 2 \#\mathcal{I} \exp\left(-\frac{nt^2}{8}\right).$$

Since  $\widehat{\Lambda} \subset \mathcal{I}$ , on the complement of  $\Omega$ ,

$$\int_{\mathcal{X}} \left\| P_{\widehat{\Lambda}}(x) - \widehat{P}_{\widehat{\Lambda}}(x) \right\|^2 d\rho(x) = \sum_{\substack{I \in \widehat{\Lambda} \\ \rho_I > 0}} \rho_I \|\widehat{c}_I - c_I\|^2 \leq Nt^2,$$

and bound (41) follows by replacing  $t$  with  $\sqrt{t/N}$ .  $\square$

**Remark 21.** *By inspecting the proof, it is possible to check that the assumption  $\widehat{x}_I^* \in \mathcal{X}$  is needed only in this proposition and it can be replaced by the condition that  $\inf_{x \in \mathcal{X}} \|\widehat{x}_I^* - x\| \leq 1$ , so that  $\|\widehat{x}_I^* - c_I\| \leq 2$ .*

#### A.4 $B$ and $D$ terms: Stability of $P_{\Lambda}$ with respect to the partition.

The following result is well-known.

**Lemma 22.** *Given a cell  $I \in \mathbb{T}$  with  $\rho_I > 0$ , for all  $t > 0$*

$$\mathbb{P}\left[\left|\frac{n_I}{n} - \rho_I\right| \geq \rho_I t\right] \leq 2 \exp\left(-\frac{n\rho_I t^2}{2(1+t/3)}\right) \leq 2 \exp\left(-\frac{n\rho_I t^2}{M_I}\right), \quad (44)$$

where  $M_I = 2/3 \max\{4, (1 + 2\rho_I)/\rho_I\}$ .

*Proof.* We apply the Bernstein inequality, see [5, Cor. 2.11], to the family of independent random variables  $\mathbb{1}_I(X_1), \dots, \mathbb{1}_I(X_n)$ , which satisfy

$$\begin{aligned} \mathbb{E}[\mathbb{1}_I(X_i)] &= \rho_I & i = 1, \dots, n \\ \sum_{i=1}^n \mathbb{E}[\mathbb{1}_I(X_i)^2] &= n\rho_I \\ \sum_{i=1}^n \mathbb{E}[\mathbb{1}_I(X_i)^m] &= n\rho_I \leq \frac{n\rho_I}{2} m! \left(\frac{1}{3}\right)^{m-2} \quad m \in \mathbb{N}, m \geq 3, \end{aligned}$$

then,

$$\mathbb{P}[|n_I - n\rho_I| \geq n\rho_I t] \leq 2 \exp\left(-\frac{(n\rho_I t)^2}{2(n\rho_I + n\rho_I t/3)}\right) = 2 \exp\left(-\frac{n\rho_I t^2}{2(1+t/3)}\right).$$

Observing that

$$|n_I - n\rho_I| \leq n \max\{\rho_I, 1 - \rho_I\},$$

if  $t > \max\{1, 1/\rho_I - 1\} = t^*$ , then  $\mathbb{P}[|n_I - n\rho_I| \geq n\rho_I t] = 0$ , whereas if  $t \leq \max\{1, 1/\rho_I - 1\} = t^*$ , it holds that

$$2(1 + t/3) \leq 2(1 + t^*/3) = M_I,$$

so that the second bound in (44) is clear.  $\square$

The following lemma provides a concentration inequality of  $\sqrt{\frac{n_I}{n}}$ .

**Lemma 23.** *Given a cell  $I \in \mathbb{T}$  with  $\rho_I > 0$ , for all  $t > 0$*

$$\mathbb{P}\left[\left|\sqrt{\frac{n_I}{n}} - \sqrt{\rho_I}\right| \geq t\right] \leq 2 \exp\left(-\frac{nt^2}{2}\right). \quad (45)$$

*Proof.* We apply Proposition 27 with  $\mathcal{Y} = \{0, 1\}$

$$\xi_i = \mathbb{1}_I(X_i) \quad f(y_1, \dots, y_n) = \frac{1}{n} \sum_{i=1}^n y_i,$$

where  $f$  is clearly bounded, and

$$f(\xi_1, \dots, \xi_n) = \frac{n_I}{n} \quad \mathbb{E}[f(\xi_1, \dots, \xi_n)] = \rho_I.$$

Given  $k = 1, \dots, n$ , it holds that

$$V_k(\xi_1, \dots, \xi_n) = \frac{1}{n} \sup_{y \in \mathcal{Y}} (\mathbb{1}_I(X_i) - y) = \frac{1}{n} \mathbb{1}_I(X_i) \leq \frac{1}{n},$$

so that  $\alpha = 1/n$ . Furthermore

$$\sum_k V_k^2(\xi_1, \dots, \xi_n) = \frac{1}{n^2} \sum_k \mathbb{1}_I(X_i) = \frac{1}{n} f(\xi_1, \dots, \xi_n),$$

then  $\beta = 1/n$  and Eq. (58) implies (45).  $\square$

The following lemma shows that, given a cell  $I \in \mathbb{T}$ ,  $\widehat{\epsilon}_I$  concentrates around  $\epsilon_I$ .

**Lemma 24.** *Given  $I \in \mathbb{T}$ , for all  $t > 0$*

$$\mathbb{P}[|\widehat{\epsilon}_I - \epsilon_I| > t] \leq 2\ell \exp\left(-n \frac{t^2}{64\ell \text{diam}_\rho(I)^2}\right). \quad (46)$$

where  $\ell = 1 + \#\mathcal{C}(I)$ .

*Proof.* Fix  $I \in \mathbb{T}$ . If  $\rho_J = 0$  for some  $J \in \mathcal{C}(I)$ , then almost surely  $X_i \notin J$  for all  $i = 1, \dots, n$  and, hence,  $n_J = 0$ , so that both  $\widehat{\epsilon}_I$  and  $\epsilon_I$  do not depend on the children  $J$ . Hence, without loss of generality, we can assume that  $\rho_J > 0$  for all  $J \in \mathcal{C}(I)$ .

Let  $\ell = \#\mathcal{C}(I) + 1$ . Set  $L^2(\mathcal{C}(I)) = \mathbb{R}^{\ell-1}$ , regarded as Euclidean vector space whose norm is denoted by  $\|\cdot\|_2$ . Define  $v, \widehat{v}, \widehat{w} \in L^2(\mathcal{C}(I))$

$$v(J) = \sqrt{\rho_J} \|c_J - c_I\| \quad \widehat{v}(J) = \sqrt{\frac{n_J}{n}} \|\widehat{c}_J - \widehat{c}_I\| \quad \widehat{w}(J) = \sqrt{\rho_J} \|\widehat{c}_J - \widehat{c}_I\|.$$

Then

$$|\widehat{\epsilon}_I - \epsilon_I| = \left| \|\widehat{v}\|_2 - \|v\|_2 \right| \leq \|\widehat{v} - v\|_2 \leq \|\widehat{w} - v\|_2 + \|\widehat{v} - \widehat{w}\|_2. \quad (47)$$

We now bound the first term. Set  $\mathcal{I} = \mathcal{C}(I) + \{I\}$  and  $\#\mathcal{I} = \ell$ . Then

$$\begin{aligned} \|\widehat{w} - v\|_2^2 &= \sum_{J \in \mathcal{C}(I)} \rho_J \left| \|\widehat{c}_J - \widehat{c}_I\| - \|c_J - c_I\| \right|^2 \leq 2 \sum_{J \in \mathcal{C}(I)} \rho_J \left( \|\widehat{c}_J - c_J\|^2 + \|\widehat{c}_I - c_I\|^2 \right) \\ &\leq 2 \sum_{J \in \mathcal{I}} \rho_J \|\widehat{c}_J - c_J\|^2 \leq 2\ell \max_{J \in \mathcal{I}} \rho_J \|\widehat{c}_J - c_J\|^2. \end{aligned}$$

Setting

$$\Omega_I = \bigcup_{J \in \mathcal{I}} \left\{ \|\widehat{c}_J - c_J\| > \frac{t}{\sqrt{2\ell\rho_J}} \right\},$$

bound (43) gives

$$\mathbb{P}[\Omega_I] \leq 2\ell \exp\left(-\frac{nt^2}{16\ell \operatorname{diam}_\rho(I)^2}\right),$$

so that

$$\mathbb{P}[\|\widehat{w} - v\|_2 > t] \leq 2\ell \exp\left(-\frac{nt^2}{16\ell \operatorname{diam}_\rho(I)^2}\right). \quad (48)$$

We now bound the second term in (47). Set

$$\Omega'_I = \bigcup_{J \in \mathcal{C}(I)} \left\{ \left| \sqrt{\frac{n_J}{n}} - \sqrt{\rho_J} \right| > \frac{t}{\sqrt{\ell \operatorname{diam}_\rho(I)}} \right\},$$

bound (44) gives

$$\mathbb{P}[\Omega'_I] \leq 2\ell \exp\left(-\frac{nt^2}{2\ell \operatorname{diam}_\rho(I)^2}\right),$$

On the complement on  $\Omega'_I$ ,

$$\begin{aligned} \|\widehat{w} - \widehat{v}\|_2^2 &= \sum_{J \in \mathcal{C}(I)} \|\widehat{c}_J - \widehat{c}_I\|^2 \left( \sqrt{\frac{n_J}{n}} - \sqrt{\rho_J} \right)^2 \\ &\leq \ell \operatorname{diam}_\rho(I)^2 \sup_{J \in \mathcal{C}(I)} \left| \sqrt{\frac{n_J}{n}} - \sqrt{\rho_J} \right|^2 \leq t^2. \end{aligned}$$

Hence

$$\mathbb{P}[\|\widehat{w} - \widehat{v}\|_2 > t] \leq 2\ell \exp\left(-\frac{nt^2}{2\ell \operatorname{diam}_\rho(I)^2}\right). \quad (49)$$

Inequality (47) with bounds (48) and (49) implies (46).  $\square$

The next proposition shows that  $P_\Lambda$  is stable under suitable small perturbations of the partition  $\Lambda$ .

**Proposition 25.** For any  $\eta > 0$

$$\mathbb{P} \left[ \left( \int_{\mathcal{X}} \left\| P_{\Lambda(\hat{\mathcal{T}}_\eta \cup \mathcal{T}_{2\eta})}(x) - P_{\Lambda(\hat{\mathcal{T}}_\eta \cap \mathcal{T}_{\eta/2})}(x) \right\|^2 d\rho(x) \right) > 0 \right] \lesssim (n^\gamma + \eta^{-\frac{2}{2\sigma+1}}) \exp(-c_a n \eta^2) \quad (50a)$$

and

$$\mathbb{P} \left[ \left( \int_{\mathcal{X}} \left\| P_{\Lambda(\hat{\mathcal{T}}_\eta \cup \mathcal{T}_{2\eta})}(x) - P_{\Lambda(\hat{\mathcal{T}}_\eta)}(x) \right\|^2 d\rho(x) \right) > 0 \right] \lesssim \eta^{-\frac{2}{2\sigma+1}} \exp(-c_a n \eta^2), \quad (50b)$$

where

$$c_a = \frac{1}{128(a+1)}. \quad (50c)$$

*Proof.* Recalling the definition of  $j_n$  given by (6), set  $\mathcal{T}_n^* = \bigcup_{j \leq j_n} \Lambda_j$ , which is a subtree with

$$\#\mathcal{T}_n^* \leq \sum_{j=0}^{j_n} a^j = \frac{a^{j_n+1} - 1}{a - 1} \lesssim n^\gamma, \quad (51)$$

and, by construction,  $\widehat{\mathcal{T}}_\eta \subset \mathcal{T}_n^*$ . The probability of the event in the left hand side of (50a) is clearly bounded by the probability of the event

$$\{\widehat{\mathcal{T}}_\eta \cap \mathcal{T}_{\eta/2} \subsetneq \widehat{\mathcal{T}}_\eta \cup \mathcal{T}_{2\eta}\} = \bigcup_{I \in \mathbb{T}} \{I \in \widehat{\mathcal{T}}_\eta \wedge I \notin \mathcal{T}_{\eta/2}\} \cup \{I \in \mathcal{T}_{2\eta} \wedge I \notin \widehat{\mathcal{T}}_\eta\}.$$

About the first term, we observe that if  $I \in \widehat{\mathcal{T}}_\eta \subset \mathcal{T}_n^*$ , then there exist  $k \geq 0$  and  $J \in \mathcal{C}^k(I) \cap \mathcal{T}_n^*$  such that  $\widehat{\epsilon}_J \geq \eta$  and, since  $I \notin \mathcal{T}_{\eta/2}$  and  $J \in \mathcal{C}^k(I)$ , then  $\epsilon_J < \frac{\eta}{2}$ , so that

$$\bigcup_{I \in \mathcal{T}_n^*} \{I \in \widehat{\mathcal{T}}_\eta \wedge I \notin \mathcal{T}_{\eta/2}\} \subset \bigcup_{J \in \mathcal{T}_n^*} \{\widehat{\epsilon}_J \geq \eta \wedge \epsilon_J < \frac{\eta}{2}\} \subset \bigcup_{J \in \mathcal{T}_n^*} \{|\widehat{\epsilon}_J - \epsilon_J| > \frac{\eta}{2}\}.$$

By union bound and (46) with  $t = \eta/2$ ,  $\text{diam}_\rho(I) \leq 1$  and  $\ell \leq a + 1$  give

$$\begin{aligned} \mathbb{P} \left[ \bigcup_{I \in \mathcal{T}_n^*} \{\widehat{\epsilon}_I \geq \eta \wedge \epsilon_I < \frac{\eta}{2}\} \right] &\lesssim \#\mathcal{T}_n^* \exp(-c_a n \eta^2) \\ &\lesssim n^\gamma \exp(-c_a n \eta^2), \end{aligned} \quad (52)$$

where the second inequality is a consequence of (51) and  $c_a$  is given by (50c).

By a similar argument

$$\{I \in \mathcal{T}_{2\eta} \wedge I \notin \widehat{\mathcal{T}}_\eta\} \subset \bigcup_{J \in \mathcal{T}_{2\eta}} \{|\widehat{\epsilon}_J - \epsilon_J| > \eta\}.$$

By union bound and (46) with  $t = \eta$  and  $\text{diam}_\rho(I) \leq 1$  give

$$\mathbb{P} \left[ \bigcup_{J \in \mathcal{T}_{2\eta}} \{\epsilon_J \geq 2\eta \wedge \widehat{\epsilon}_J < \eta\} \right] \lesssim \#\mathcal{T}_{2\eta} \exp(-c_a n \eta^2) \lesssim \eta^{-\frac{2}{2\sigma+1}} \exp(-c_a n \eta^2), \quad (53)$$

where the second inequality is a consequence of (39a).

By (53) and (52), we get (50a). The proof of (50b) can be deduced reasoning as for (53).  $\square$

## B Auxialiry results

We recall the following probabilistic inequality based on a result of [24, 25], see also [31, Thm. 3.3.4] and [23] for concentration inequalities for Hilbert-space-valued random variables.

**Proposition 26.** *Let  $\xi_1, \dots, \xi_n$  be a family of independent zero-mean random variables taking values in a real separable Hilbert space and satisfying*

$$\mathbb{E}[\|\xi_i\|^m] \leq \frac{1}{2} m! \Sigma^2 M^{m-2} \quad \forall m \geq 2, \quad (54)$$

where  $\Sigma$  and  $M$  are two positive constants. Then, for all  $n \in \mathbb{N}$  and  $t > 0$

$$\mathbb{P} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \right\| \geq t \right] \leq 2 \exp \left( - \frac{nt^2}{\Sigma^2 + Mt + \Sigma \sqrt{\Sigma^2 + 2Mt}} \right) = 2 \exp \left( -n \frac{\Sigma^2}{M^2} g \left( \frac{Mt}{\Sigma^2} \right) \right) \quad (55)$$

where  $g(t) = \frac{t^2}{1+t+\sqrt{1+2t}}$ . In particular, if  $\xi_i$  are bounded by  $M$  with probability 1, then

$$\mathbb{P} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \right\| \geq t \right] \leq 2 \exp \left( - \frac{nt^2}{4M^2} \right). \quad (56)$$

*Proof.* Bound(55) is given in [24] with a wrong factor, see [25]. To show(56), note that (54) is satisfied with  $\Sigma = M$ . Furthermore, for all  $t \leq 1$ ,  $g(t) \geq t^2/4$ , so that if  $t \leq M$ ,

$$\mathbb{P} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \right\| \geq t \right] \leq 2 \exp \left( - \frac{nt^2}{4M^2} \right)$$

whereas, if  $t > M$ , (56) is trivially satisfied.  $\square$

The following concentration inequality is based on [18] and we adapt the proof of Theorem 10 in [19]. We introduce the following notation. Given a family  $\xi_1, \dots, \xi_n$  of independent random variables taking value in some measurable space  $\mathcal{Y}$  and a measurable positive bounded function  $f : \mathcal{Y}^n \rightarrow \mathbb{R}$ , for any  $k = 1, \dots, n$  set

$$\begin{aligned} V_k &= f(\xi_1, \dots, \xi_n) - \inf_{y \in \mathcal{Y}} f(\xi_1, \dots, \xi_{k-1}, y, \xi_{k+1}, \dots, \xi_n) \\ &= \sup_{y \in \mathcal{Y}} \left( f(\xi_1, \dots, \xi_n) - f(\xi_1, \dots, \xi_{k-1}, y, \xi_{k+1}, \dots, \xi_n) \right). \end{aligned}$$

**Proposition 27.** *With the above notation, if there exist two constants  $\alpha, \beta > 0$  such that*

$$\max_{k=1, \dots, n} V_k \leq \alpha \quad (57a)$$

$$\sum_{k=1}^n V_k^2 \leq \beta f(\xi_1, \dots, \xi_n) \quad (57b)$$

then, for any  $t > 0$

$$\mathbb{P} \left[ \left| \sqrt{f(\xi_1, \dots, \xi_n)} - \sqrt{\mathbb{E}[f(\xi_1, \dots, \xi_n)]} \right| > t \right] \leq 2 \exp \left( - \frac{t^2}{2 \max\{\alpha, \beta\}} \right). \quad (58)$$

*Proof.* Set  $Z_k = V_k/\alpha$  and  $Z = f(\xi_1, \dots, \xi_n)/\alpha$ . By construction

$$\begin{aligned} Z_k &\leq 1 \quad k = 1, \dots, n \\ \sum_{k=1}^n Z_k^2 &\leq \frac{\beta}{\alpha} Z. \end{aligned}$$

Let  $\gamma = \max\{\beta/\alpha, 1\}$ . Theorem 13 of [18] gives that

$$\mathbb{P} [\mathbb{E}[Z] - Z > t] \leq \exp\left(-\frac{t^2}{2\gamma\mathbb{E}[Z]}\right) \quad (59)$$

$$\mathbb{P} [Z - \mathbb{E}[Z] > t] \leq \exp\left(-\frac{t^2}{2\gamma\mathbb{E}[Z] + \gamma t}\right). \quad (60)$$

By replacing  $t$  with  $2t\sqrt{\mathbb{E}[Z]}$  in (59)

$$\begin{aligned} \exp\left(-\frac{2t^2}{\gamma}\right) &\geq \mathbb{P} \left[ \mathbb{E}[Z] - 2t\sqrt{\mathbb{E}[Z]} + t^2 > Z + t^2 \right] = \mathbb{P} \left[ \left| \sqrt{\mathbb{E}[Z]} - t \right| > \sqrt{Z + t^2} \right] \\ &\geq \mathbb{P} \left[ \sqrt{\mathbb{E}[Z]} - \sqrt{Z} > 2t \right] \end{aligned}$$

since

$$\sqrt{\mathbb{E}[Z]} - t \leq \left| \sqrt{\mathbb{E}[Z]} - t \right| \leq \sqrt{Z + t^2} \leq \sqrt{Z} + t$$

provided that  $\left| \sqrt{\mathbb{E}[Z]} - t \right| \leq \sqrt{Z + t^2}$ . Hence

$$\mathbb{P} \left[ \sqrt{\mathbb{E}[Z]} - \sqrt{Z} > t \right] \leq \exp\left(-\frac{t^2}{2\gamma}\right). \quad (61)$$

Setting  $\frac{t^2}{2\mathbb{E}[Z]+t} = 2\tau^2$ , bound (60) gives

$$\begin{aligned} \exp\left(-\frac{2\tau^2}{\gamma}\right) &\geq \mathbb{P} \left[ Z - \mathbb{E}[Z] > \tau^2 + \sqrt{\tau^4 + 4\tau^2\mathbb{E}[Z]} \right] \geq \mathbb{P} \left[ Z - \mathbb{E}[Z] > 4\tau^2 + 4\tau\sqrt{\mathbb{E}[Z]} \right] \\ &= \mathbb{P} \left[ Z > \left( \sqrt{\mathbb{E}[Z]} + 2\tau \right)^2 \right] = \mathbb{P} \left[ \sqrt{Z} - \sqrt{\mathbb{E}[Z]} > 2\tau \right], \end{aligned}$$

so that, setting  $\tau = t/2$ ,

$$\mathbb{P} \left[ \sqrt{Z} - \sqrt{\mathbb{E}[Z]} > t \right] \leq \exp\left(-\frac{t^2}{2\gamma}\right). \quad (62)$$

Bounds (61) and (62) imply that

$$\mathbb{P} \left[ \left| \sqrt{Z} - \sqrt{\mathbb{E}[Z]} \right| > t \right] \leq 2 \exp\left(-\frac{t^2}{2\gamma}\right),$$

and, by replacing  $t$  with  $t/\sqrt{\alpha}$ ,

$$\mathbb{P} \left[ \left| \sqrt{f(\xi_1, \dots, \xi_n)} - \sqrt{\mathbb{E}[f(\xi_1, \dots, \xi_n)]} \right| > t \right] \leq 2 \exp\left(-\frac{t^2}{2\alpha\gamma}\right).$$

where  $\alpha\gamma = \alpha \max\{\beta/\alpha, 1\} = \max\{\beta, \alpha\}$ . □

## References

- [1] William K Allard, Guangliang Chen, and Mauro Maggioni. Multi-scale geometric methods for data sets II: Geometric multi-resolution analysis. *Appl. Comput. Harmon. Anal.*, 32(3):435–462, 2012.
- [2] David Arthur and Sergei Vassilvitskii. k-means++: The advantages of careful seeding. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 1027–1035. Society for Industrial and Applied Mathematics, 2007.
- [3] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps and spectral techniques for embedding and clustering. In *Advances in Neural Information Processing Systems*, volume 14, pages 585–591, 2001.
- [4] Peter Binev, Albert Cohen, Wolfgang Dahmen, Ronald DeVore, and Vladimir Temlyakov. Universal algorithms for learning theory part I: piecewise constant functions. *J. Mach. Learn. Res.*, 6(Sep):1297–1321, 2005.
- [5] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities*. Oxford University Press, Oxford, 2013.
- [6] Guillermo Canas, Tomaso Poggio, and Lorenzo Rosasco. Learning manifolds with k-means and k-flats. In *Advances in Neural Information Processing Systems*, pages 2465–2473, 2012.
- [7] Michael Christ. A  $T(b)$  theorem with remarks on analytic capacity and the Cauchy integral. *Colloq. Math.*, 60/61(2):601–628, 1990.
- [8] David L. Donoho and Carrie Grimes. Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data. *Proc. Natl. Acad. Sci. USA*, 100(10):5591–5596, 2003.
- [9] Giacomo Gigante and Paul Leopardi. Diameter bounded equal measure partitions of Ahlfors regular metric measure spaces. *Discrete Comput. Geom.*, 57(2):419–430, 2017.
- [10] Misha Gromov. *Metric structures for Riemannian and non-Riemannian spaces*. Birkhäuser Boston, Inc., Boston, MA, english edition, 2007.
- [11] Peter M Gruber. Optimum quantization and its applications. *Adv. Math.*, 186(2):456–497, 2004.
- [12] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer New York Inc., 2001.
- [13] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An introduction to statistical learning*, volume 112. Springer, 2013.
- [14] Wenjing Liao and Mauro Maggioni. Adaptive geometric multiscale approximations for intrinsically low-dimensional data. *J. Mach. Learn. Res.*, 20(98):1–83, 2019.
- [15] Stuart Lloyd. Least squares quantization in PCM. *IEEE Trans. Inform. Theory*, 28(2):129–137, 1982.

- [16] Mauro Maggioni, Stanislav Minsker, and Nate Strawn. Dictionary learning and non-asymptotic bounds for geometric multi-resolution analysis. *PAMM. Proc. Appl. Math. Mech.*, 14(1):1013–1016, 2014.
- [17] Stéphane Mallat. *A Wavelet Tour of Signal Processing, Third Edition: The Sparse Way*. Academic Press, Inc., 3rd edition, 2008.
- [18] Andreas Maurer. Concentration inequalities for functions of independent variables. *Random Structures Algorithms*, 29(2):121–138, 2006.
- [19] Andreas Maurer and Massimiliano Pontil. Empirical Bernstein bounds and sample-variance penalization. In *COLT*, 2009.
- [20] Boaz Nadler, Stéphane Lafon, Ronald R Coifman, and Ioannis G Kevrekidis. Diffusion maps, spectral clustering and reaction coordinates of dynamical systems. *Appl. Comput. Harmon. Anal.*, 21(1):113–127, 2006.
- [21] Bruno Olshausen and David Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(607):6583, 1996.
- [22] Peter Petersen. *Riemannian geometry*. Springer, Cham, third edition, 2016.
- [23] I. F. Pinelis and A. I. Sakhanenko. Remarks on inequalities for probabilities of large deviations. *Theory Probab. Appl.*, 30(1):143–148, 1985.
- [24] Iosif Pinelis. Optimum bounds for the distributions of martingales in Banach spaces. *Ann. Probab.*, 22(4):1679–1706, 1994.
- [25] Iosif Pinelis. Correction: “Optimum bounds for the distributions of martingales in Banach spaces” [*Ann. Probab.* **22** (1994), no. 4, 1679–1706; MR1331198 (96b:60010)]. *Ann. Probab.*, 27(4):2119, 1999.
- [26] Bernhard Schölkopf, Alexander J. Smola, and Klaus-Robert Müller. Advances in kernel methods. In Bernhard Schölkopf, Christopher J. C. Burges, and Alexander J. Smola, editors, *International conference on artificial neural networks*, chapter Kernel Principal Component Analysis, pages 327–352. MIT Press, Cambridge, MA, USA, 1999.
- [27] Laurent Schwartz. *Cours d’analyse. 3*. Hermann, Paris, second edition, 1993.
- [28] Clayton Scott, Robert D Nowak, et al. Minimax-optimal classification with dyadic decision trees. *IEEE Trans. Inform. Theory*, 52(4):1335–1353, 2006.
- [29] Josh Tenenbaum, Vin De Silva, and John Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.
- [30] Joe H Ward Jr. Hierarchical grouping to optimize an objective function. *J. Amer. Statist. Assoc.*, 58(301):236–244, 1963.
- [31] V. Yurinsky. *Sums and Gaussian vectors*, volume 1617. Springer-Verlag, Berlin, 1995.