

UNIVERSAL BAYES CONSISTENCY IN METRIC SPACES

BY STEVE HANNEKE¹, ARYEH KONTOROVICH^{2,*}, SIVAN SABATO^{2,†} AND ROI WEISS³

¹*Toyota Technological Institute at Chicago, steve.hanneke@gmail.com*

²*Department of Computer Science, Ben-Gurion University of the Negev, karyeh@cs.bgu.ac.il; sabatos@cs.bgu.ac.il*

³*Department of Computer Science, Ariel University, roiw@ariel.ac.il*

We extend a recently proposed 1-nearest-neighbor based multiclass learning algorithm and prove that our modification is universally strongly Bayes consistent in all metric spaces admitting *any* such learner, making it an “optimistically universal” Bayes-consistent learner. This is the first learning algorithm known to enjoy this property; by comparison, the k -NN classifier and its variants are not generally universally Bayes consistent, except under additional structural assumptions, such as an inner product, a norm, finite dimension or a Besicovitch-type property.

The metric spaces in which universal Bayes consistency is possible are the “essentially separable” ones—a notion that we define, which is more general than standard separability. The existence of metric spaces that are not essentially separable is widely believed to be independent of the ZFC axioms of set theory. We prove that essential separability exactly characterizes the existence of a universal Bayes-consistent learner for the given metric space. In particular, this yields the first impossibility result for universal Bayes consistency.

Taken together, our results completely characterize strong and weak universal Bayes consistency in metric spaces.

1. Introduction. Since their inception in the 1950s [12]—or, according to some accounts, nearly 1000 years earlier [35]—nearest-neighbor methods have provided an intuitive and reliable suite of techniques for performing classification in metric spaces. For k -NN based methods, it has been generally understood that some notion of finite dimensionality is both necessary and sufficient for the methods to be Bayes consistent under all distributions over the metric space—a property known as *universal Bayes consistency* (UBC). However, a complete characterization of the metric spaces in which *any* nearest neighbor method (or other learners, for that matter) is UBC has been so far unknown. For the problem of multiclass classification, we resolve these questions exhaustively.

To answer these questions, we study a compression-based 1-NN algorithm for multiclass classification, which was proposed in 2017 [28], and shown to be strongly UBC in all metric spaces of bounded diameter and doubling dimension. It was also shown that there exist infinite-dimensional spaces in which this algorithm is strongly Bayes consistent, while classic k -NN based methods provably are not. Left open was the full characterization of metric spaces in which this algorithm is UBC. In this work, we provide this characterization. Moreover, we prove that this algorithm is UBC in any metric space for which a UBC algorithm exists, thus resolving the above fundamental open question about nearest-neighbor methods.

Main results. We design a generalized version of the algorithm used in [28], which we call OptiNet (see Algorithm 1). The contribution of this paper is twofold: (i) We show that OptiNet is universally strongly Bayes consistent in all essentially separable metric spaces.

Received June 2019; revised October 2020.

MSC2020 subject classifications. Primary 54E70, 97K80, 62C12; secondary 03E17, 03E55.

Key words and phrases. Metric space, nearest neighbor, classification, Bayes consistency.

A formal definition of *essential separability*—our broadening of the standard notion of separability—is given in Section 3. Briefly, in an essentially separable metric space, the total mass of every probability measure is contained in some separable subspace. Whether every metric space is essentially separable is widely believed to hinge upon set-theoretic axioms that are independent of ZFC, having to do with the existence of certain measurable cardinals (we provide the relevant set-theoretic background in Section 4.1). (ii) We show that in any set-theoretic model that allows the existence of nonessentially separable metric spaces, no (strong or weak) universally Bayes-consistent learner is possible on such spaces. To our knowledge, this is the first construction of a learning setting in which universal Bayes consistency is impossible. In contrast, if one adopts a set-theoretic model in which every metric space is essentially separable, then OptiNet is always universally strongly Bayes consistent. As such, OptiNet is *optimistically* universally Bayes consistent for metric spaces, in a sense analogous to [22]: it succeeds whenever success is possible, and is the first learning algorithm known to enjoy this property. For comparison, k -NN and other existing nearest-neighbor approaches are only universally Bayes consistent under additional structural assumptions, such as an inner product, a norm, a finite dimension or a Besicovitch-type property [3, 4, 6], all of which are significantly stronger assumptions than essential separability.

Taken together, our results completely characterize strong and weak UBC in metric spaces.

Related work. Nearest-neighbor methods were initiated by Fix and Hodges in 1951 [12] and, in the celebrated k -NN formulation, have been placed on a solid theoretical foundation [7, 9, 10, 39, 42]. Following the pioneering work of [9, 39] on nearest-neighbor classification, it was shown by [10, 21, 42] that the k -NN classifier is universally strongly Bayes consistent in $(\mathbb{R}^d, \|\cdot\|_2)$. These results made extensive use of the Euclidean structure of \mathbb{R}^d , but in [38] a weak Bayes-consistency result was shown for metric spaces with a bounded diameter and finite doubling dimension, and additional distributional smoothness assumptions.

Consistency of NN-type algorithms in more general (and, in particular, infinite-dimensional) metric spaces was discussed in [1, 3, 4, 6, 13, 33]. Characterizations of Bayes consistency for the standard k -NN [6, 13] and for a generalized “moving window” classification rule [1] were given in terms of a Besicovitch-type condition (see Section 5 for a more detailed discussion). By Besicovitch’s density theorem [15], in $(\mathbb{R}^d, \|\cdot\|_2)$, and more generally in finite-dimensional normed spaces, the aforementioned condition holds for all distributions; however, in infinite-dimensional spaces this condition may be violated [36, 37]. The violation of the Besicovitch condition is not an isolated pathology—occurring, for example, in the commonly used Gaussian–Hilbert spaces [40]. Leveraging the consistency of k -NN in finite dimensions, the *filtering* technique (taking the first d coordinates in some basis representation for an appropriate d) was shown to be universally weakly consistent in [3]. However, that technique is only applicable in separable Hilbert spaces, as opposed to more general metric spaces. For compact metric spaces, the SVM algorithm can be made universally Bayes consistent by using an appropriate kernel [8].

Although the classic 1-NN classifier is well known to be inconsistent in general, in recent years a series of papers has presented various ways of learning a *regularized* 1-NN classifier, as an alternative to k -NN. Gottlieb et al. [17] showed that an approximate nearest-neighbor search can act as a regularizer, actually improving generalization performance rather than just injecting noise. This technique was extended to multiclass classification in [30]. In a follow-up work, [31] showed that applying Structural Risk Minimization (SRM) to a margin-regularized data-dependent bound very similar to that in [17] yields a strongly Bayes-consistent 1-NN classifier in doubling spaces with a bounded diameter.

Approaching the problem through the lens of sample compression, a computationally near-optimal nearest-neighbor condensing algorithm was presented in [19] and later extended to

cover semimetric spaces [18]; both were based on constructing γ -nets in spaces with a finite doubling dimension (or its semimetric analogue). As detailed in [31], margin-regularized 1-NN methods enjoy a number of statistical and computational advantages over the traditional k -NN classifier. Salient among these are explicit data-dependent generalization bounds, and considerable runtime and memory savings. Sample compression affords additional advantages, in the form of tighter generalization bounds and increased efficiency in time and space. Recently, [28] provided evidence that this technique has wider applicability than k -NN methods, by exhibiting an infinite-dimensional metric measure space where the compression-based learner is Bayes consistent, while k -NN methods provably fail.

The work of Devroye et al. [10], Theorem 21.2, has implications for 1-NN classifiers in $(\mathbb{R}^d, \|\cdot\|_2)$ that are defined based on data-dependent majority-vote partitions of the space. They showed that a *fixed* mapping from each sample size to a data-dependent partitioning rule, satisfying some regularity conditions, induces a universally strongly Bayes-consistent algorithm. This result requires the partitioning rule to have a VC dimension that grows sub-linearly in the sample size, and since this rule must be fixed in advance, the algorithm is not fully adaptive. Theorem 19.3 *ibid.* proves weak consistency for an inefficient compression-based algorithm, which selects among all the possible compression sets of a certain size, and maintains a certain rate of compression relative to the sample size. The generalizing power of sample compression was independently discovered by [10, 34], and later elaborated upon by [20, 23]. In the context of NN classification, [10] lists various condensing heuristics (which have no known performance guarantees) and leaves open the algorithmic question of how to minimize the empirical risk over all subsets of a given size.

The margin-based technique developed in [17, 30] relied on computing a minimum vertex cover. Thus, it was not possible to make it simultaneously and computationally efficient and Bayes consistent when the number of labels exceeds two, since Vertex Cover on general graphs is an NP-hard problem. Although one could resort to a 2-approximation algorithm for vertex cover, this presents an obstruction to establishing the Bayes consistency of the classifier.

In [27], an active-learning algorithm was presented, which across a broad spectrum of natural noise regimes, reduced the sample complexity roughly quadratically. Along the way, this work circumvented the computational obstacle associated with computing a minimum vertex cover on a general graph: the trick was to construct a γ -net and take the majority label (more accurately, the *plurality*—that is, the most frequent—label; we shall use the more familiar terms “majority label” and “majority vote”) in each Voronoi region. The majority is determined by actively querying each region, where the number of calls depends on the density and noise level of the region.

Paper outline. After setting down the definitions in Section 2, we describe in Section 3 the compression-based 1-NN algorithm OptiNet studied in this paper and its consistency on essentially-separable metric spaces is proved. In Section 4, we prove that no universally Bayes-consistent algorithm exists on metric spaces that are not essentially separable. We conclude with a discussion in Section 5. Appendices containing the proofs of several auxiliary lemmas are given in the Supplementary Material [24].

2. Definitions and notation. Our *instance space* is the metric probability space (\mathcal{X}, ρ, μ) , where ρ is a metric and μ is a probability measure. By definition, the Borel σ -algebra \mathcal{B} supporting μ is the smallest σ -algebra containing the open sets of ρ . For any $x \in \mathcal{X}$ and $r > 0$, denote by $B_r(x)$ the open ball of radius r around x under the metric ρ :

$$B_r(x) = \{x' \in \mathcal{X} : \rho(x, x') < r\}.$$

We consider a countable label set \mathcal{Y} . The unknown sampling distribution is a probability measure $\bar{\mu}$ over $\mathcal{X} \times \mathcal{Y}$, with marginal μ over \mathcal{X} . Denote by $(X, Y) \sim \bar{\mu}$ a pair drawn according to $\bar{\mu}$. The generalization error of a classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$ is given by

$$\text{err}(f) := \mathbb{P}_{\bar{\mu}}[Y \neq f(X)],$$

and its empirical error with respect to a labeled set $S' \subseteq \mathcal{X} \times \mathcal{Y}$ is given by

$$\widehat{\text{err}}(f, S') := \frac{1}{|S'|} \sum_{(x,y) \in S'} \mathbf{1}[y \neq f(x)].$$

The optimal Bayes risk of $\bar{\mu}$ is $R_{\bar{\mu}}^* := \inf \text{err}(f)$, where the infimum is taken over all measurable classifiers $f : \mathcal{X} \rightarrow \mathcal{Y}$. We omit the subscript $\bar{\mu}$ when there is no ambiguity and denote the optimal Bayes risk of $\bar{\mu}$ by R^* .

For a labeled sequence $S = (x_i, y_i)_{i=1}^n \in (\mathcal{X} \times \mathcal{Y})^n$ and any $x \in \mathcal{X}$, let $X_{\text{nn}}(x, S)$ be the nearest neighbor of x with respect to S and let $Y_{\text{nn}}(x, S)$ be the nearest-neighbor label of x with respect to S :

$$(X_{\text{nn}}(x, S), Y_{\text{nn}}(x, S)) := \underset{(x_i, y_i) \in S}{\text{argmin}} \rho(x, x_i),$$

where ties are broken lexicographically, that is, the smallest x_i is chosen, with respect to a fixed total ordering of the space \mathcal{X} (such an ordering can always be chosen to be measurable, see Appendix D). The 1-NN classifier induced by S is defined as $h_S(x) := Y_{\text{nn}}(x, S)$. For any $m \in \mathbb{N}$, any sequence $\mathbf{X} = \{x_1, \dots, x_m\} \in \mathcal{X}^m$ induces a *Voronoi partition* of \mathcal{X} , $\mathcal{V}(\mathbf{X}) := \{V_1(\mathbf{X}), \dots, V_m(\mathbf{X})\}$, where each Voronoi cell is

$$V_i(\mathbf{X}) := \left\{ x \in \mathcal{X} : i = \underset{1 \leq j \leq m}{\text{argmin}} \rho(x, x_j) \right\},$$

again breaking ties lexicographically. In particular, for $\mathbf{X} = \{X_i : (X_i, Y_i) \in S\}$, we have $h_S(x) = Y_i$ for all $x \in V_i(\mathbf{X})$.

A 1-NN algorithm is a mapping from an i.i.d. labeled sample $S_n \sim \bar{\mu}^n$ to a labeled set $S'_n \subseteq \mathcal{X} \times \mathcal{Y}$, yielding the 1-NN classifier $h_{S'_n}$. While the classic 1-NN algorithm sets $S'_n := S_n$, the algorithm which we analyze chooses S'_n adaptively. More generally, a learning algorithm Alg is a mapping (possibly randomized) from a labeled sequence $S_n = (x_i, y_i)_{i=1}^n \in (\mathcal{X} \times \mathcal{Y})^n$ to $\text{Alg}(S_n) \in \mathcal{Y}^{\mathcal{X}}$, satisfying some natural measurability requirements spelled out in Remark 4.10 below. We say that Alg is *strongly Bayes consistent* under $\bar{\mu}$ if $\text{err}(\text{Alg}(S_n))$ converges to R^* almost surely,

$$\mathbb{P} \left[\lim_{n \rightarrow \infty} \text{err}(\text{Alg}(S_n)) = R^* \right] = 1.$$

Similarly, Alg is *weakly Bayes consistent* under $\bar{\mu}$ if $\text{err}(\text{Alg}(S_n))$ converges to R^* in expectation,

$$\lim_{n \rightarrow \infty} \mathbb{E}[\text{err}(\text{Alg}(S_n))] = R^*.$$

Obviously, the former implies the latter. We say that Alg is *universally Bayes consistent* on a metric space if Alg is Bayes consistent for every distribution supported on its Borel σ -algebra \mathcal{B} . Specializing to OptiNet , we have $\text{Alg}(S_n) = h_{S'_n}$.

For $A \subseteq \mathcal{X}$ and $\gamma > 0$, a γ -net of A is any *maximal set* $B \subseteq A$ in which all interpoint distances are at least γ . In separable metric spaces, all γ -nets are at most countable. Denote the diameter of a set $A \subseteq \mathcal{X}$ by $\text{diam}(A) \in [0, \infty]$. For a partition \mathcal{A} , $\text{diam}(\mathcal{A})$ denotes the maximum diameter $\text{diam}(A)$ among all cells $A \in \mathcal{A}$.

For $n \in \mathbb{N}$, define $[n] := \{1, \dots, n\}$. Given a labeled set $S_n = (x_i, y_i)_{i \in [n]}$, $d \in [n]$, and any $\mathbf{i} = \{i_1, \dots, i_d\} \in [n]^d$, denote the subsample of S_n indexed by \mathbf{i} by $S_n(\mathbf{i}) :=$

TABLE 1
Symbols guide

Symbol	Brief description
(\mathcal{X}, ρ, μ)	metric probability space
\mathcal{B}	Borel σ -algebra induced by ρ
$B_r(x)$	open ball of radius r around x
$\text{err}(f)$	generalization error of $f : \mathcal{X} \rightarrow \mathcal{Y}$
$\widehat{\text{err}}(f, S')$	empirical error of $f : \mathcal{X} \rightarrow \mathcal{Y}$ on S'
R^*	Bayes risk
$S_n = (X_n, Y_n)$	random sample of size n
$S_n(\mathbf{i}, \mathbf{j})$	subsample (X_j, Y_j) of S_n indexed by \mathbf{i} and \mathbf{j}
$S_n(\mathbf{i})$	subsample $S_n(\mathbf{i}, \mathbf{i})$
$S_n(\mathbf{i}, *)$	subsample (X_j, Y^*) with true majority vote labels
$X(\gamma)$	γ -net of X_n
$S_n(\gamma)$	subsample $(X(\gamma), Y(\gamma))$ with empirical majority votes
$M_n(\gamma) = 2 X(\gamma) $	size of the compression
h_S	1-NN classifier induced by the labeled set S
$\alpha_n(\gamma)$	empirical error of $h_{S_n(\gamma)}$ on S_n
$\text{UB}_\gamma(A)$	γ -envelope of $A \subseteq \mathcal{X}$
$L_\gamma(A)$	γ -missing mass of $A \subseteq \mathcal{X}$
$\mathcal{V}(X)$	Voronoi partition of \mathcal{X} induced by X
$\text{Alg}(S) = \hat{h}_S$	classifier obtained by learning algorithm Alg when given sample S

$\{(x_{i_1}, y_{i_1}), \dots, (x_{i_d}, y_{i_d})\}$. Similarly, for a vector $\mathbf{y}' = \{y'_1, \dots, y'_d\} \in \mathcal{Y}^d$, define $S_n(\mathbf{i}, \mathbf{y}') := \{(x_{i_1}, y'_{i_1}), \dots, (x_{i_d}, y'_{i_d})\}$, namely the subsample of S_n as determined by \mathbf{i} where the labels are replaced with \mathbf{y}' . Lastly, for $\mathbf{i}, \mathbf{j} \in [n]^d$, we denote $S_n(\mathbf{i}; \mathbf{j}) := \{(x_{i_1}, y_{j_1}), \dots, (x_{i_d}, y_{j_d})\}$.

We use standard order-of-magnitude notation throughout the paper; thus, for $f, g : \mathbb{N} \rightarrow [0, \infty)$ we write $f(n) \in O(g(n))$ to mean $\limsup_{n \rightarrow \infty} f(n)/g(n) < \infty$ and $f(n) \in o(g(n))$ to mean $\limsup_{n \rightarrow \infty} f(n)/g(n) = 0$. Likewise, $f(n) \in \Omega(g(n))$ means that $g(n) \in O(f(n))$. In accordance with common convention, we often use the less precise notation $f(n) = O(g(n))$, etc.

The main notation are summarized in Table 1; some are introduced in later sections.

3. Universal Bayes consistency in separable metric spaces. In this section, we describe a variant of the 1-NN majority-based compression algorithm developed in the series of papers [27–29], adapted to maintain measurability in potentially nonseparable metric spaces. We show that this variant is universally Bayes consistent in all separable metric spaces, and the extension to essential separability is immediate, as will become clear below.

Our variant, OptiNet, is formally presented in Algorithm 1. It operates as follows. The input is the sample S_n ; the set of points in the sample is denoted by $X_n = \{X_1, \dots, X_n\}$. The algorithm defines a set Γ of all scales $\gamma > 0$, which are interpoint distances in X_n , and the additional scale $\gamma = \infty$. For each scale in Γ , the algorithm constructs a γ -net of X_n ; note that any singleton in X_n is an ∞ -net. Denote the constructed γ -net by $X(\gamma) := \{X_{i_1}, \dots, X_{i_{M/2}}\}$, where $M/2 \equiv M_n(\gamma)/2 := |X(\gamma)|$ denotes its size and $\mathbf{i} \equiv \mathbf{i}(\gamma) := \{i_1, \dots, i_{M/2}\} \in [n]^{M/2}$ denotes the indices selected from S_n for this γ -net. For each γ -net, OptiNet finds the empirical majority vote labels in the Voronoi cells defined by the partition $\mathcal{V}(X(\gamma)) = \{V_1(X(\gamma)), \dots, V_{M/2}(X(\gamma))\}$; these labels are denoted by $Y'(\gamma) \in \mathcal{Y}^{M/2}$. Formally, for $i \in [M/2]$,

$$(3.1) \quad Y'_i(\gamma) := \operatorname{argmax}_{y \in \mathcal{Y}} |\{j \in [n] : X_j \in V_i(X(\gamma)), Y_j = y\}|,$$

Algorithm 1 (OptiNet) The 1-NN compression-based algorithm

Input: sample $S_n = (X_i, Y_i)_{i \in [n]}$, confidence $\delta \in (0, 1)$

Output: A 1-NN classifier

- 1: let $\Gamma := (\{\rho(X_i, X_j) : i, j \in [n]\} \cup \{\infty\}) \setminus \{0\}$
 - 2: **for** $\gamma \in \Gamma$ **do**
 - 3: let $\mathbf{X}(\gamma)$ be a γ -net of $\{X_1, \dots, X_n\}$
 - 4: let $M_n(\gamma) := 2|\mathbf{X}(\gamma)|$
 - 5: for each $i \in [M_n(\gamma)/2]$, let $Y'_i(\gamma)$ be the most frequent label in $V_i(\mathbf{X}(\gamma))$ as in (3.1)
 - 6: set $S'_n(\gamma) := (\mathbf{X}(\gamma), \mathbf{Y}'(\gamma))$
 - 7: **end for**
 - 8: Set $\alpha_n(\gamma) := \widehat{\text{err}}(h_{S'_n(\gamma)}, S_n)$
 - 9: find $\gamma_n^* \in \text{argmin}_{\gamma \in \Gamma} Q(n, \alpha_n(\gamma), M_n(\gamma), \delta)$, where Q is defined in (3.2)
 - 10: set $S'_n := S'_n(\gamma_n^*)$
 - 11: **return** $h_{S'_n}$
-

where ties are broken based on a fixed preference order on the countable set \mathcal{Y} . The result of the procedure is a labeled set $S'_n(\gamma) := S_n(\mathbf{i}(\gamma), \mathbf{Y}'(\gamma))$ for every possible scale $\gamma \in \Gamma$. The algorithm then selects one scale $\gamma^* \equiv \gamma_n^*$ from Γ , and outputs the hypothesis that it induces, $h_{S'_n(\gamma^*)}$. The choice of γ^* is based on minimizing a generalization error bound, denoted Q , which upper bounds $\text{err}(h_{S'_n(\gamma)})$ with high probability. The error bound is derived based on a compression-based analysis, as follows.

For an even integer $m \leq 2n$, we say that a specific S'_n is an (α, m) -compression of S_n if there exist $\mathbf{i}, \mathbf{j} \in [n]^{m/2}$ such that $S'_n = S_n(\mathbf{i}, \mathbf{j})$ and $\widehat{\text{err}}(h_{S'_n}, S_n) \leq \alpha$. Note that at most m examples from S_n determine $h_{S'_n}$, hence this is a compression scheme of size at most m .

The papers [28] and [29] give a consistency result for the original algorithm of [27], on metric spaces with a finite doubling dimension and a finite diameter, under the following assumptions on the generalization error bound $Q(n, \alpha, m, \delta)$:

Q1. For any $n \in \mathbb{N}$ and $\delta \in (0, 1)$, with probability at least $1 - \delta$ over $S_n \sim \bar{\mu}^n$, for all $\alpha \in [0, 1]$ and even $m \in [2n]$: If S'_n is an (α, m) -compression of S_n , then

$$\text{err}(h_{S'_n}) \leq Q(n, \alpha, m, \delta).$$

Q2. For any fixed $n \in \mathbb{N}$ and $\delta \in (0, 1)$, Q is monotonically increasing in α and in m .

Q3. There is a sequence $\{\delta_n\}_{n=1}^\infty$, $\delta_n \in (0, 1)$ such that $\sum_{n=1}^\infty \delta_n < \infty$, and for all m ,

$$\lim_{n \rightarrow \infty} \sup_{\alpha \in [0, 1]} (Q(n, \alpha, m, \delta_n) - \alpha) = 0.$$

Here, we provide a consistency result that holds for more general metric spaces. We prove that OptiNet is universally strongly Bayes consistent in all *essentially separable* metric spaces. Recall that (\mathcal{X}, ρ) is *separable* if it contains a dense countable set. A metric probability space (\mathcal{X}, ρ, μ) is *separable* if there is a measurable $\mathcal{X}' \subseteq \mathcal{X}$ with $\mu(\mathcal{X}') = 1$ such that (\mathcal{X}', ρ) is separable. We will call a metric space (\mathcal{X}, ρ) *essentially separable (ES)* if, for every probability measure μ on \mathcal{B} , the metric probability space (\mathcal{X}, ρ, μ) is separable.

To prove this stronger result, we require a slightly stronger version of property **Q3**.

Q3'. There is a sequence $\{\delta_n\}_{n=1}^\infty$, $\delta_n \in (0, 1)$ such that $\sum_{n=1}^\infty \delta_n < \infty$, and for any sequence $m_n \in o(n)$,

$$\lim_{n \rightarrow \infty} \sup_{\alpha \in [0, 1]} (Q(n, \alpha, m_n, \delta_n) - \alpha) = 0.$$

Property **Q3'** is slightly stronger than **Q3**, since it allows m to grow as $o(n)$ instead of keeping it as a constant. The compression bound used in [28] does not satisfy this property, since it includes a term of the order $m \log(n)/(n - m)$. Therefore, if $m_n = \Omega(n/\log(n))$, then $m_n = o(n)$, yet this term does not converge to zero for $n \rightarrow \infty$, thus precluding consistency of the algorithm in [28] for such cases. We provide here a tighter compression bound, which does satisfy **Q3'**.

LEMMA 3.1. For $m \leq n - 2$, define

$$(3.2) \quad Q(n, \alpha, m, \delta) := \frac{n}{n - m} \alpha + \sqrt{\frac{8(\frac{n}{n-m})\alpha(m \ln(2en/m) + \ln(2n/\delta))}{n - m}} + \frac{9(m \ln(2en/m) + \ln(2n/\delta))}{n - m}.$$

For $m > n - 2$, define $Q(n, \alpha, m, \delta) := \max(1, Q(n, \alpha, n - 2, \delta))$. Then the function Q satisfies the properties **Q1**, **Q2**, **Q3'**.

The approach to obtaining property **Q3'** is inspired by refinements of compression-based generalization bounds holding for the special case of compression schemes which have a *permutation-invariant* reconstruction function [20]. While $h_{S'_n}$ cannot quite be expressed as a permutation-invariant function of a subset of the (X_i, Y_i) data points, it can be expressed as a function that is invariant to permutations of two subsets of (X_i, Y_i) points. This is used in the proof of Lemma 3.1, which is provided in Appendix A.2, to derive the tighter compression bound in (3.2). This bound is derived using Bernstein's inequality over $n - m$ random variables and applying a union bound over all $\binom{n}{m/2}^2, 1 \leq m/2 \leq n$, possible compressions.

Our main technical innovation, which allows us to dispose of the finiteness requirements on the dimension and the diameter of the metric space that were assumed in [28], is the sublinear growth of γ -nets. Another straightforward but crucial insight is to approximate functions in $L^1(\mu) := \{f : \int |f| d\mu < \infty\}$ by Lipschitz ones, rather than by continuous functions with compact support as in [28]. The latter approximation requires local compactness, which essentially amounts to a finite dimensionality condition. Our new approach does not require local compactness or finite dimensionality.

THEOREM 3.2. Let (\mathcal{X}, ρ, μ) be a separable metric probability space. Let Q be a generalization bound that satisfies Properties **Q1**, **Q2**, **Q3'**, and let δ_n be as stipulated by **Q3'**. If the input confidence δ for input size n is set to δ_n , then the 1-NN classifier $h_{S'_n(\gamma_n^*)}$ calculated by OptiNet is strongly Bayes consistent on (\mathcal{X}, ρ, μ) :

$$\mathbb{P}\left[\lim_{n \rightarrow \infty} \text{err}(h_{S'_n(\gamma_n^*)}) = R^*\right] = 1.$$

REMARK 3.3. OptiNet selects the scale γ based on a compression bound. This creates a close connection between the algorithm and the proof of consistency below. However, it is worth noting that it is possible instead to choose γ based on a hold-out validation set: for instance, using $n/2$ of the n samples to construct the predictor for each possible γ value, and then from among these values γ , one can select the γ whose predictor makes the smallest number of mistakes on the remaining $n/2$ samples. Since the analysis of [28] (see [29]), and its generalization below, show that there exists a choice of γ^* for each n such that OptiNet is Bayes consistent, this alternative technique of selecting γ based on a hold-out sample would only lose an additive $O(\sqrt{\log(n)/n})$ compared to using that γ^* , and hence would also be Bayes consistent.

REMARK 3.4. OptiNet is computationally efficient. Using a farthest first-traversal procedure such as Algorithm 1 in [32], one can construct the γ -nets simultaneously for all γ values, including their corresponding empirical errors, in $O(n^2)$ time, leading to a total runtime of $O(n^2)$.

Given a sample $S_n \sim \bar{\mu}^n$, we abbreviate the optimal empirical error $\alpha_n^* = \alpha(\gamma_n^*)$ and the optimal compression size $M_n^* = M(\gamma_n^*)$ as computed by OptiNet. As discussed above, the labeled set $S'_n(\gamma_n^*)$ computed by OptiNet is a (α_n^*, M_n^*) -compression of the sample S_n . For brevity, we denote

$$Q_n(\alpha, m) := Q(n, \alpha, m, \delta_n).$$

To prove Theorem 3.2, we first follow the standard technique, used also in [29], of decomposing the excess error over the Bayes error into two terms:

$$\begin{aligned} \text{err}(h_{S'_n(\gamma_n^*)}) - R^* &= (\text{err}(h_{S'_n(\gamma_n^*)}) - Q_n(\alpha_n^*, M_n^*)) + (Q_n(\alpha_n^*, M_n^*) - R^*) \\ &=: T_I(n) + T_{II}(n). \end{aligned}$$

We now show that each term decays to zero almost surely. For the first term, $T_I(n)$, we have, similar to [29], that Property Q1 implies that for any $n > 0$,

$$(3.3) \quad \mathbb{P}[\text{err}(h_{S'_n(\gamma_n^*)}) - Q_n(\alpha_n^*, M_n^*) > 0] \leq \delta_n.$$

Based on the Borel–Cantelli lemma and the fact that $\sum \delta_n < \infty$, we have $\limsup_{n \rightarrow \infty} T_I(n) \leq 0$ with probability 1.

The main difference from the proof in [29] is in the argument establishing $\limsup_{n \rightarrow \infty} T_{II}(n) \leq 0$ almost surely. We now show that the generalization bound $Q_n(\alpha_n^*, M_n^*)$ also approaches the Bayes error R^* , thus proving $\limsup_{n \rightarrow \infty} T_{II}(n) \leq 0$ almost surely.

We will show below that there exist $N = N(\varepsilon) > 0$, $\gamma = \gamma(\varepsilon) > 0$, and universal constants $c, C > 0$ such that $\forall n \geq N$,

$$(3.4) \quad \mathbb{P}[Q_n(\alpha_n(\gamma), M_n(\gamma)) > R^* + \varepsilon] \leq Cne^{-cne\varepsilon^2} + 1/n^2.$$

For any $\gamma > 0$ (even if $\gamma \notin \Gamma$), OptiNet finds γ_n^* such that

$$Q_n(\alpha_n^*, M_n^*) = \min_{\gamma' \in \Gamma} Q_n(\alpha_n(\gamma'), M_n(\gamma')) \leq Q_n(\alpha_n(\gamma), M_n(\gamma)).$$

The bound in (3.4) thus implies that $\forall n \geq N$,

$$(3.5) \quad \mathbb{P}[Q_n(\alpha_n^*, M_n^*) > R^* + \varepsilon] \leq Cne^{-cne\varepsilon^2} + 1/n^2.$$

By the Borel–Cantelli lemma, this implies that almost surely,

$$\limsup_{n \rightarrow \infty} T_{II}(n) = \limsup_{n \rightarrow \infty} (Q_n(\alpha_n^*, M_n^*) - R^*) \leq 0.$$

Since $\forall n, T_I(n) + T_{II}(n) \geq 0$, this implies $\lim_{n \rightarrow \infty} T_{II}(n) = 0$ almost surely, thus completing the proof of Theorem 3.2.

It remains to prove (3.4). We note that a simpler form of (3.4) is proved in [29], where they relied on the finiteness of the dimension and the diameter of the space to upper bound the compression size $M_n(\gamma)$ with probability 1. For $A \subseteq \mathcal{X}$, denote its γ -envelope by $\text{UB}_\gamma(A) := \bigcup_{x \in A} B_\gamma(x)$ and consider the γ -missing mass of S_n , defined as the following random variable:

$$(3.6) \quad L_\gamma(S_n) := \mu(\mathcal{X} \setminus \text{UB}_\gamma(S_n)).$$

We bound the left-hand side of (3.4) using a function $n \mapsto t_\gamma(n)$ of order $o(n)$, used to upper bound the compression size; t_γ will be specified below:

$$\begin{aligned}
 & \mathbb{P}[Q_n(\alpha_n(\gamma), M_n(\gamma)) > R^* + \varepsilon] \\
 (3.7) \quad & \leq \mathbb{P}\left[Q_n(\alpha_n(\gamma), M_n(\gamma)) > R^* + \varepsilon \wedge L_\gamma(S_n) \leq \frac{\varepsilon}{10} \wedge M_n(\gamma) \leq t_\gamma(n)\right] \\
 & \quad + \mathbb{P}[L_\gamma(S_n) > \varepsilon/10] + \mathbb{P}[M_n(\gamma) > t_\gamma(n)] \\
 & =: P_I + P_{II} + P_{III}.
 \end{aligned}$$

First, we bound P_I . By a union bound,

$$\begin{aligned}
 & \mathbb{P}\left[Q_n(\alpha_n(\gamma), M_n(\gamma)) > R^* + \varepsilon \wedge L_\gamma(S_n) \leq \frac{\varepsilon}{10} \wedge M_n(\gamma) \leq t_\gamma(n)\right] \\
 & \leq \sum_{d=1}^{t_\gamma(n)} \mathbb{P}\left[Q_n(\alpha_n(\gamma), M_n(\gamma)) > R^* + \varepsilon \wedge L_\gamma(S_n) \leq \frac{\varepsilon}{10} \wedge M_n(\gamma) = d\right].
 \end{aligned}$$

Thus, it suffices to bound each term in the right-hand sum separately. We do so in the following lemma.

LEMMA 3.5. *There exists a function $\varepsilon \mapsto \gamma(\varepsilon)$ for $\varepsilon > 0$, such that under the conditions of Theorem 3.2, there exists an n_0 such that for all $n \geq n_0$, and for all $d \in [t_\gamma(n)]$, letting $\gamma := \gamma(\varepsilon)$,*

$$p_d := \mathbb{P}\left[Q_n(\alpha_n(\gamma), M_n(\gamma)) > R^* + \varepsilon \wedge L_\gamma(S_n) \leq \frac{\varepsilon}{10} \wedge M_n(\gamma) = d\right] \leq e^{-\frac{n\varepsilon^2}{32}}.$$

Applying Lemma 3.5 and summing over all $1 \leq d \leq t_\gamma(n)$, we have that, for n large enough so that $t_\gamma(n) \leq n$,

$$(3.8) \quad P_I \leq \sum_{d=1}^{t_\gamma(n)} p_d \leq t_\gamma(n)e^{-\frac{n\varepsilon^2}{32}} \leq ne^{-\frac{n\varepsilon^2}{32}}.$$

Lemma 3.5 is a generalization of Lemma 10 in [29]. The main difference is that Lemma 10 holds in doubling spaces and uses the fixed map $t_\gamma(n) = 2\lceil \text{diam}(\mathcal{X})/\gamma \rceil^{\text{ddim}}$ for all $n \in \mathbb{N}$. The proof of Lemma 3.5 is the same as that of Lemma 10 in [29], except for two changes that adapt it for a general metric space. First, where Lemma 10 uses the fact that t_γ is set to a constant function, and thus $\lim_{n \rightarrow \infty} t_\gamma(n)/n = 0$, the proof of Lemma 3.5 uses instead the property that $t_\gamma(n) = o(n)$, which again leads to the same limit.

In addition, the proof of Lemma 3.5 employs a new result, Lemma 3.6 given below, instead of Lemma 8 from [29]. Lemma 8 from [29] states that for metric spaces with a finite doubling dimension and diameter, Bayes error R^* can be approached using classifiers defined by the true majority-vote labeling over fine partitions of \mathcal{X} . Here, we prove that this holds for general metric spaces. Let $\mathcal{V} = \{V_1, \dots\}$ be a countable partition of \mathcal{X} , and define the function $I_\mathcal{V} : \mathcal{X} \rightarrow \mathcal{V}$ such that $I_\mathcal{V}(x)$ is the unique $V \in \mathcal{V}$ for which $x \in V$. For any measurable set $\emptyset \neq E \subseteq \mathcal{X}$, define the true majority-vote label $y^*(E)$ by

$$(3.9) \quad y^*(E) = \underset{y \in \mathcal{V}}{\text{argmax}} \mathbb{P}(Y = y \mid X \in E),$$

where ties are broken lexicographically. Given \mathcal{V} and a measurable set $W \subseteq \mathcal{X}$, define the true majority-vote classifier $h_{\mathcal{V}, W}^* : \mathcal{X} \rightarrow \mathcal{V}$ given by

$$(3.10) \quad h_{\mathcal{V}, W}^*(x) = y^*(I_\mathcal{V}(x) \cap W).$$

The new lemma can now be stated as follows.

LEMMA 3.6. *Let $\bar{\mu}$ be a probability measure on $\mathcal{X} \times \mathcal{Y}$, where \mathcal{X} is a metric probability space. For any $\nu > 0$, there exists a diameter $\beta = \beta(\nu) > 0$ such that for any countable measurable partition $\mathcal{V} = \{V_1, \dots\}$ of \mathcal{X} and any measurable set $W \subseteq \mathcal{X}$ satisfying:*

- (i) $\mu(\mathcal{X} \setminus W) \leq \nu$
- (ii) $\text{diam}(\mathcal{V} \cap W) \leq \beta$,

the true majority-vote classifier $h_{\mathcal{V},W}^$ defined in (3.10) satisfies*

$$\text{err}(h_{\mathcal{V},W}^*) \leq R^* + 5\nu.$$

The proof of Lemma 3.6 is identical to that of Lemma 8 from [29], except for the following change: in the proof of Lemma 8 from [29], they use their Lemma 7, which states that on doubling spaces, the set of continuous functions with compact support is dense in $L^1(\mu)$. To remove the requirement of compact support, which restricts the type of spaces for which this lemma holds, we use instead a stronger approximation result, which states that Lipschitz functions are dense in $L^1(\mu)$ for any metric probability space. For completeness, we include a proof of this fact in Lemma A.1 in the Supplementary Material [24], where a complete proof of Lemma 3.6 is also given.

Having established Lemma 3.6, this completes the necessary generalizations to obtain Lemma 3.5, whose proof is given in Appendix A.4 for completeness. This proves the bound on P_I claimed in (3.8).

We now turn to constructing the function t_γ , which bounds the compression size (i.e., twice the γ -net size) with high probability, and bounding P_{II} and P_{III} .

LEMMA 3.7. *Let (\mathcal{X}, ρ, μ) be a separable metric probability space. For $S_n \sim \mu^n$, let $X(\gamma)$ be any γ -net of S_n . Then, for any $\gamma > 0$, there exists a function $t_\gamma : \mathbb{N} \rightarrow \mathbb{R}_+$ in $o(n)$ such that*

$$(3.11) \quad \mathbb{P}\left[\sup_{\gamma\text{-nets } X(\gamma)} 2|X(\gamma)| \geq t_\gamma(n) \right] \leq 1/n^2.$$

This result can be compared to the case of finite-dimensional and finite-diameter metric spaces, in which one can set $t_\gamma(n) := 2 \lceil \frac{\text{diam}(\mathcal{X})}{\gamma} \rceil^{\text{ddim}}$ for all $n \in \mathbb{N}$, where ddim is the (finite) doubling dimension and $\text{diam}(\mathcal{X})$ is the diameter of the space, and get that $\mathbb{P}[M_n(\gamma) \geq t_\gamma(n)] = 0$. The proof of Lemma 3.7 is provided in Appendix A of the Supplementary Material [24].

This lemma implies that $P_{III} \leq 1/n^2$, while a bound on P_{II} , which bounds the γ -missing-mass $L_\gamma(S_n)$, is furnished by the following lemma, whose proof is given in Appendix A of the Supplementary Material [24].

LEMMA 3.8. *Let (\mathcal{X}, ρ, μ) be a separable metric probability space, $\gamma > 0$ be fixed, and the γ -missing mass L_γ defined as in (3.6). Then there exists a function $u_\gamma : \mathbb{N} \rightarrow \mathbb{R}_+$ in $o(1)$, such that for $S_n \sim \mu^n$ and all $t > 0$,*

$$(3.12) \quad \mathbb{P}[L_\gamma(S_n) \geq u_\gamma(n) + t] \leq \exp(-nt^2).$$

Taking n sufficiently large so that $u_\gamma(n)$, as furnished by Lemma 3.8, satisfies $u_\gamma(n) \leq \varepsilon/20$, and invoking Lemma 3.8 with $t = \varepsilon/20$, we have

$$(3.13) \quad P_{II} = \mathbb{P}[L_\gamma(S_n) > \varepsilon/10] \leq e^{-\frac{n\varepsilon^2}{400}}.$$

Plugging (3.8), (3.13) and $P_{III} \leq 1/n^2$ into (3.7), we get that (3.4) holds, which completes the proof of Theorem 3.2.

4. Essential separability is necessary for universal Bayes consistency. Recall that a metric space (\mathcal{X}, ρ) is *essentially separable* (ES) if for every probability measure μ on the Borel σ -algebra \mathcal{B} , the metric probability space (\mathcal{X}, ρ, μ) is separable; namely, there is an $\mathcal{X}' \subseteq \mathcal{X}$ with $\mu(\mathcal{X}') = 1$ such that (\mathcal{X}', ρ) is separable. In Theorem 3.2, we established that OptiNet is indeed universally Bayes consistent (UBC) for all such spaces. As such, essential separability of a metric space is sufficient for the existence of a UBC learning rule in that space. In this section, we show that essential separability is also necessary for such a rule to exist.

The metric spaces one typically encounters in statistics and machine learning are all ES, as reflected by Dudley’s remark that “for practical purposes, a probability measure defined on the Borel sets of a metric space is always concentrated in some separable subspace” [11]. The question of whether non-ES metric spaces exist at all turns out to be rather subtle. It is widely believed that the existence of non-ES spaces is *independent of the ZFC axioms of set theory* (see Section 4.1 for further details). In other words, it is believed that assuming that ZFC is consistent, its axioms neither necessitate nor preclude the existence of non-ES metric spaces.

The main contribution of this section is to show that in any non-ES metric space (if one exists), no learning rule is UBC.

THEOREM 4.1. *Let (\mathcal{X}, ρ) be a non-ES metric space equipped with the Borel σ -algebra \mathcal{B} . Then no (weak or strong) UBC algorithm exists on (\mathcal{X}, ρ) .*

Combining this result with Theorem 3.2, the following result is immediate, revealing that OptiNet is *optimistically* UBC (adopting the terminology of [22]), in the sense that the only required assumption on (\mathcal{X}, ρ) is that UBC learning is *possible*.

COROLLARY 4.2. *OptiNet is UBC in every metric space for which there exists a UBC learning rule.*

REMARK 4.3. Theorem 4.1 is somewhat unusual, in that it identifies a setting in which no universal Bayes-consistent procedure exists. To our knowledge, this is the first such impossibility result. Also unusual, for a statistics paper, is the appearance of esoteric set theory. See [2] for another recent result discussing a setting in which learnability is independent of ZFC.

In the next section, we provide necessary preliminaries. Theorem 4.1 is proved in Section 4.2.

4.1. Preliminaries. We collect necessary definitions and known results about non-ES metric spaces. In particular, we connect the existence of non-ES metric spaces with the existence of *real-valued measurable cardinals* (Definition 4.5 below). A thorough treatment of the latter, including most of the material in this subsection, can be found in [26]; a more gentle introduction to the subject can be found in [25]. Throughout the following presentation, we work under the standard Zermelo–Fraenkel set theory together with the Axiom of Choice, commonly abbreviated as ZFC.

Cardinals. We denote the cardinality of a set A by $|A|$. The first infinite (countable) cardinal is denoted by $\aleph_0 = |\omega|$, where ω is the set of all finite cardinals. In particular, \aleph_0 is the cardinality of the set of natural numbers, $\aleph_0 = |\mathbb{N}|$. We write $[A]^n$ to denote the family of all subsets of A of size $n \in \mathbb{N}$, and $[A]^{<\omega} := \bigcup_{n \in \mathbb{N}} [A]^n$ is the family of all finite subsets of A .

The smallest uncountable cardinal is denoted by \aleph_1 . The cardinality of the real numbers, also known as the *continuum*, is $\mathfrak{c} = |\mathbb{R}|$. It is well known that $\mathfrak{c} = 2^{\aleph_0} \geq \aleph_1 > \aleph_0$. The continuum hypothesis states that $\mathfrak{c} = \aleph_1$. It is known that its truth value is independent of ZFC, so that either the continuum hypothesis or its negation can be added as an axiom to ZFC set theory while maintaining its consistency status. In the following, we do not include the continuum hypothesis (or its negation) in our set theory; thus, our discussion includes models of ZFC in which $\mathfrak{c} > \aleph_1$.

Nontrivial probability measures. Let $(\mathcal{X}, \mathcal{B})$ be a measurable space. Recall that a probability measure on \mathcal{X} , henceforth called a *measure*, is a function $\mu : \mathcal{B} \rightarrow [0, 1]$ satisfying:

- (i) $\mu(\emptyset) = 0$ and $\mu(\mathcal{X}) = 1$;
- (ii) if $A, B \in \mathcal{B}$ and $A \subseteq B$ then $\mu(A) \leq \mu(B)$;
- (iii) if $\{A_i\}_{i=1}^\infty \subseteq \mathcal{B}$ are pairwise disjoint then $\mu(\bigcup_{i=1}^\infty A_i) = \sum_{i=1}^\infty \mu(A_i)$.

A measure μ is *nontrivial* if it vanishes on singletons: $\mu(\{x\}) = 0, \forall x \in \mathcal{X}$. Nontrivial measures play a key role in establishing the impossibility of UBC in non-ES spaces and we will be concerned mainly with such measures.

Another important property of a measure is its additivity. For a cardinal κ , a measure μ is κ -*additive* if for any $\beta < \kappa$ and any pairwise disjoint measurable family $\{A_\alpha \in \mathcal{B} : \alpha < \beta\}$,

$$(4.1) \quad \mu\left(\bigcup_{\alpha < \beta} A_\alpha\right) = \sum_{\alpha < \beta} \mu(A_\alpha) := \sup_{B \in [\beta]^{<\omega}} \sum_{\alpha \in B} \mu(A_\alpha).$$

By definition, any measure is \aleph_1 -additive, commonly known as σ -*additive*. The following lemma states the main property of nontrivial and κ -additive measures that will be used here. Its proof follows directly from the definitions.

LEMMA 4.4. *Let μ be a nontrivial and κ -additive measure on \mathcal{B} . Then any set $A \in \mathcal{B}$ with $|A| < \kappa$ has $\mu(A) = 0$.*

Non-ES metric spaces and real-valued measurable cardinals. Before giving the formal definition of real-valued measurable cardinals and establishing their relation to general non-ES metric spaces, let us first illustrate the main ideas which will be presented below, using a simple example of an uncountable discrete metric space. Consider the metric space $([0, 1], \rho_{\text{dis}})$, where ρ_{dis} is the *discrete metric*, defined as

$$(4.2) \quad \rho_{\text{dis}}(x, x') := \mathbf{1}[x \neq x'], \quad x, x' \in \mathcal{X}.$$

The Borel σ -algebra on $([0, 1], \rho_{\text{dis}})$ is all of $2^{[0,1]}$; thus, all subsets of $[0, 1]$ are measurable. This metric space is clearly nonseparable; the interesting question is whether it is ES. In other words, does there exist a measure on $([0, 1], \rho_{\text{dis}})$ that does not have a separable support?

Note that any nontrivial measure on $([0, 1], \rho_{\text{dis}})$ suffices to prove that it is non-ES. Indeed, for any such measure, Lemma 4.4, together with the fact that all measures are σ -additive, implies that any set of positive measure must have an uncountable cardinality. But any such set is clearly nonseparable, due to the discrete nature of the metric space. Therefore, if a nontrivial measure exists on $([0, 1], \rho_{\text{dis}})$, then it is non-ES. Conversely, if there are no nontrivial measures on $([0, 1], 2^{[0,1]})$, then the discrete metric space admits only trivial measures with countable support. So in this case $([0, 1], \rho_{\text{dis}})$ is ES. Thus, the question of whether $([0, 1], \rho_{\text{dis}})$ is non-ES is equivalent to the question of whether a nontrivial measure exists on this space. It is known that the Lebesgue measure, defined on the Borel σ -algebra generated by open sets on $([0, 1], |\cdot|)$, cannot be extended to the measurable space $([0, 1], 2^{[0,1]})$ while

simultaneously being translation invariant [15]. However, currently, other nontrivial measures on $([0, 1], 2^{[0,1]})$ are not ruled out in ZFC.

More generally, given a cardinal κ , let \mathcal{X} be some set of that cardinality, and consider the measurable space $\mathfrak{X}_\kappa := (\mathcal{X}, 2^\mathcal{X})$. As above, such a space is induced, for example, by the discrete metric ρ_{dis} . Moreover, whether $(\mathcal{X}, \rho_{\text{dis}})$ is ES depends only on the cardinality κ , and is closely related to the existence of nontrivial measures on \mathfrak{X}_κ , similarly to the example of $([0, 1], \rho_{\text{dis}})$ above. To characterize the cardinalities for which \mathfrak{X}_κ is ES, we use the known concept of *real-valued measurable cardinals* (RVMC).

DEFINITION 4.5. A cardinal κ is *real-valued measurable* if there exists a nontrivial and κ -additive measure on \mathfrak{X}_κ . Any such measure is called a *witnessing measure* for \mathfrak{X}_κ .

Clearly, any RVMC must be uncountable. We denote by κ_{min} the smallest RVMC; this cardinal exists if some RVMC exists, by the well ordering of the cardinals. The following theorem from [5] characterizes ES metric spaces in terms of κ_{min} . Recall that a set $D \subseteq \mathcal{X}$ is *discrete* if for any $x \in D$ there exists some $r_x > 0$ such that $B_{r_x}(x) \cap D = \{x\}$.

THEOREM 4.6 ([5], Appendix III, Theorem 2). *Let κ_{min} be the smallest real-valued measurable cardinal (if one exists). Then a metric space (\mathcal{X}, ρ) is ES if and only if every discrete $D \subseteq \mathcal{X}$ has $|D| < \kappa_{\text{min}}$.*

REMARK 4.7. Theorem 4.6 is stated in [5] only for the case $\kappa_{\text{min}} \leq \mathfrak{c}$. However, one can readily verify that the proof extends essentially verbatim (by replacing “atomless” with “nontrivial”) to the case $\kappa_{\text{min}} > \mathfrak{c}$ as well.

It follows that whether a given metric space (\mathcal{X}, ρ) is ES or not depends on whether any RVMC exists and, if one exists, on the cardinality of the smallest such cardinal, κ_{min} . Assuming that ZFC is consistent (which cannot be proved in ZFC, by Gödel’s second incompleteness theorem), it is well known that one cannot prove in ZFC the existence of real-valued measurable cardinals (RVMC). While it is possible that one can prove in ZFC that RVMCs do not exist, no such proof has been discovered yet. However, quoting Fremlin [14], “at present, almost no-one is seriously searching for a proof in ZFC that real-valued measurable cardinals do not exist.” In fact, currently, the vast majority of set-theoreticians believe that the existence of RVMC is *independent* of ZFC, that is, assuming that ZFC is consistent, the existence of an RVMC can neither be proven nor disproven from the axioms of ZFC.

In particular, if one adds to ZFC the axiom that no RVMC exists, then *any* metric space is ES. Alternatively, under some additional properties that are beyond the scope of this paper, one can take κ_{min} to be of cardinality that is arbitrarily large. For more details, see [26], Section 12.

The above relations (and their connection to the results to follow) are further discussed in Section 5 and illustrated in Figure 1 therein.

REMARK 4.8. It is worth mentioning that if one adds to ZFC the continuum hypothesis, $\mathfrak{c} = \aleph_1$, which is well known to be independent of ZFC, then if a RVMC exists, then it must hold that $\kappa_{\text{min}} > \mathfrak{c}$ [26]. In this case, all metric spaces of cardinality $\leq \mathfrak{c}$ are ES. In particular, the metric space $([0, 1], \rho_{\text{dis}})$ discussed at the beginning of this section admits only trivial measures with a countable support.

4.2. *UBC is impossible in non-ES metric spaces.* In this section, we prove the following theorem, which readily implies Theorem 4.1.

THEOREM 4.9. *Let (\mathcal{X}, ρ) be a non-ES metric space and let Alg be any (possibly random) learning algorithm mapping samples $S \in (\mathcal{X} \times \{0, 1\})^{<\omega}$ to classifiers $\text{Alg}(S) \in \{0, 1\}^{\mathcal{X}}$. Then, there exist a measure $\bar{\mu}$ on $\mathcal{X} \times \{0, 1\}$ (w.r.t. the Borel sets induced by ρ), a measurable classifier $h^* : \mathcal{X} \rightarrow \{0, 1\}$, and an $\varepsilon > 0$ such that, for $n \in \mathbb{N}$ and $S_n \sim \bar{\mu}^n$,*

$$(4.3) \quad \limsup_{n \rightarrow \infty} \mathbb{E}[\text{err}_{\bar{\mu}}(\text{Alg}(S_n))] \geq \text{err}_{\bar{\mu}}(h^*) + \varepsilon = R_{\bar{\mu}}^* + \varepsilon,$$

where $R_{\bar{\mu}}^*$ is the optimal Bayes error. In particular, no weak or strong UBC algorithm exists for (\mathcal{X}, ρ) .

For notational simplicity, in the following we denote $\hat{h}_S := \text{Alg}(S)$ (not to be confused with the 1-NN classifier h_S which we used in previous sections).

REMARK 4.10 (Measurability of Alg). To be strictly clear about definitions here, note that we require that the learning algorithm be *measurable*, in the sense that for every $\bar{\mu}$ and n , for $S \sim \bar{\mu}^n$, \hat{h}_S is a $\mathcal{B}(L^1(\mu))$ -measurable random variable, where μ is the marginal of $\bar{\mu}$ on \mathcal{X} , and $\mathcal{B}(L^1(\mu))$ is the Borel σ -algebra on the set of all measurable functions $\mathcal{X} \rightarrow \{0, 1\}$, induced by the $L^1(\mu)$ pseudo-metric. This is a basic criterion, without which the expected risk of \hat{h}_S is not well defined (among other pathologies).

For deterministic algorithms \hat{h} , to satisfy the above criterion, it suffices that the function $(s, x) \mapsto \hat{h}_s(x)$ on $(\mathcal{X} \times \{0, 1\})^n \times \mathcal{X}$ is a measurable $\{0, 1\}$ -valued random variable, under the product σ -algebra on $(\mathcal{X} \times \{0, 1\})^n \times \mathcal{X}$. To see this, note that for any such function, for $X \sim \mu$ independent of $S \sim \bar{\mu}^n$, for any measurable function $f : \mathcal{X} \rightarrow \{0, 1\}$, we have that $|\hat{h}_S(X) - f(X)|$ is a measurable random variable; hence the variable $\mathbb{E}[|\hat{h}_S(X) - f(X)| \mid S]$ is well defined and measurable. Therefore, for any $\varepsilon > 0$, the event that $\mathbb{E}[|\hat{h}_S(X) - f(X)| \mid S] \leq \varepsilon$ is measurable. Thus, the inverse images of balls in the $L^1(\mu)$ pseudo-metric are measurable sets, and since these balls generate $\mathcal{B}(L^1(\mu))$, this implies \hat{h}_S is a $\mathcal{B}(L^1(\mu))$ -measurable random variable.

In particular, we note that OptiNet satisfies this measurability criterion, since calculating its prediction $\hat{h}_s(x)$ involves only simple operations based on the metric ρ (which are measurable, since by definition, ρ induces the topology generating the Borel σ -algebra), and other basic measurability-preserving operations such as argmin for a finite number of indices indexing measurable quantities. Thus, our requirements of \hat{h}_S in Theorem 4.9 are satisfied by OptiNet.

REMARK 4.11. In [6], Section 2.1, the authors define the metric space (\mathcal{X}, ρ) , where $\mathcal{X} = [0, 1]$ and

$$\rho(x, x') = \mathbf{1}[x \neq x'] \cdot (1 + \mathbf{1}[xx' \neq 0])$$

and endow it with the distribution μ , which places a mass of $1/2$ on $x = 0$ and spreads the rest of the mass “uniformly” on $(0, 1]$. The deterministic labeling $h^*(x) = \mathbf{1}[x > 0]$ is imposed. The authors observe that the optimal Bayes risk is $R^* = 0$ while the (classical) 1-NN classifier achieves an asymptotic expected risk of $1/2$ —in contradistinction to the standard result that in finite-dimensional spaces 1-NN is Bayes consistent in the realizable case. The authors then use this example to argue that “[separability] is required even in finite dimension.” We find the example somewhat incomplete, because care is not taken to ensure that (\mathcal{X}, ρ, μ) is a *metric probability space*, that is, that the σ -algebra supporting μ is generated by the open sets

of ρ . Indeed, the Borel σ -algebra generated by ρ is the discrete one, $\mathcal{B} = 2^{[0,1]}$. Endowing the latter with a “uniform” measure implicitly assumes that the Lebesgue measure on the *standard* Borel σ -algebra can be extended to all subsets of $[0, 1]$ —a statement known to be equivalent to \mathfrak{c} being larger than or equal to a real-valued measurable cardinal [26]. So the above metric probability space is assumed to be non-ES, as in Theorem 4.9. Another objection is that, under any reasonable notion of *dimension*, the metric space (\mathcal{X}, ρ) would be considered \mathfrak{c} -dimensional rather than finite-dimensional.

It is worth mentioning that by Remark 4.8, if one accepts, say, the continuum hypothesis, then the above metric space becomes ES and admits only trivial measures with a countable support (so the standard k -NN, and many other algorithms, are in fact UBC in this space).

To prove Theorem 4.9, we first note that (4.3) indeed implies that no weak or strong UBC algorithm exists for (\mathcal{X}, ρ) by an application of [5], Theorem 5.4. To establish (4.3), note that since (\mathcal{X}, ρ) is non-ES, Theorem 4.6 implies that there exists a discrete set $D \subseteq \mathcal{X}$ with $|D| = \kappa_{\min}$, where κ_{\min} is the smallest real-valued measurable cardinal (see Section 4.1). Let $\mathfrak{X}_{|D|} := (D, 2^D)$. By Lemma B.1 in the Supplementary Material [24], $2^D \subseteq \mathcal{B}$. Hence, it suffices to construct the required adversarial measure on $D \times \{0, 1\}$. That being the case, from now on we set without loss of generality $\mathcal{X} := D$ and $\mathcal{B} := 2^D$.

Below, we split the argument for the construction of the required adversarial measure on $\mathcal{X} \times \{0, 1\}$ into two cases:

Case (I): $\kappa_{\min} \leq \mathfrak{c}$ and Case (II): $\kappa_{\min} > \mathfrak{c}$.

This is manifested by what is known as Ulam’s dichotomy. This dichotomy dictates the nature of nontrivial measures in the two cases. To formally state the dichotomy, we first need some additional definitions.

Let μ be a measure on \mathcal{X} . A set $A \subseteq \mathcal{X}$ is an *atom* of μ if $\mu(A) > 0$ and for every measurable $B \subseteq A$ either $\mu(B) = 0$ or $\mu(B) = \mu(A)$. A measure μ is *atomless* if it has no atoms. So in an atomless measure, for any $A \in \mathcal{B}$ with $\mu(A) > 0$ there exists a $B \subset A$ with $0 < \mu(B) < \mu(A)$. Conversely, μ is *purely atomic* if every $A \in \mathcal{B}$ with $\mu(A) > 0$ contains an atom.

Clearly, in a countable space all measures are trivial and purely atomic. However, in uncountable spaces matters are more subtle. While any atomless measure is nontrivial, one might expect that conversely a nontrivial measure cannot contain an atom. However, this is not necessarily the case. In particular, when $\kappa_{\min} > \mathfrak{c}$, measures on \mathcal{X} that are simultaneously nontrivial and purely atomic exist.

Formally, let κ be an RVMC. Recall that a witnessing measure for \mathfrak{X}_κ is a nontrivial and κ -additive measure on \mathfrak{X}_κ , namely, a measure defined over all subsets of \mathcal{X} and that vanishes on any set of cardinality $< \kappa$. We say that κ is *two-valued measurable* if there is a $\{0, 1\}$ -valued witnessing measure on \mathfrak{X}_κ , where a measure is $\{0, 1\}$ -valued (or *two-valued*) if $\mu(A) \in \{0, 1\}$ for all $A \in \mathcal{B}$. Clearly, a two-valued measure is purely atomic and satisfies that, for any countable partition $\{P_i\}_{i \in \mathbb{N}} \subseteq \mathcal{B}$ of \mathcal{X} , there exists one and only one $j \in \mathbb{N}$ such that $\mu(P_j) = 1$. We say that κ is *atomlessly measurable* if there is an atomless witnessing measure on \mathfrak{X}_κ . In 1930, Ulam established the following dichotomy (see [15], Section 543).

THEOREM 4.12 (Ulam’s dichotomy [41]). *Let κ be a real-valued measurable cardinal. Then:*

- (i) *if $\kappa \leq \mathfrak{c}$ then κ is atomlessly measurable and every witnessing measure on \mathfrak{X}_κ is atomless;*
- (ii) *if $\kappa > \mathfrak{c}$ then κ is two-valued measurable and every witnessing measure on \mathfrak{X}_κ is purely atomic.*

In other words, if κ is atomlessly measurable then $\kappa \leq \mathfrak{c}$, while if κ is two-valued measurable then $\kappa > \mathfrak{c}$.

We now proceed to construct the adversarial measures on $\mathcal{X} \times \{0, 1\}$ by considering the two cases (I) and (II) above separately.

(I) *The case $\kappa_{\min} \leq \mathfrak{c}$.* By Ulam’s dichotomy in Theorem 4.12, $|\mathcal{X}| = \kappa_{\min}$ is atomlessly measurable, so there exists an atomless witnessing measure μ on \mathcal{B} . Fix such a μ and define the induced set-difference pseudometric

$$\Delta(A, B) = \mu(\{A \cup B\} \setminus \{A \cap B\}), \quad A, B \in \mathcal{B}.$$

Define the metric space (\mathcal{U}, Δ) , where $\mathcal{U} \subseteq \mathcal{B}$ is the quotient σ -algebra under the equivalence relation $A \sim B \Leftrightarrow \Delta(A, B) = 0$. The measure μ induces the corresponding functional $\tilde{\mu} : \mathcal{U} \rightarrow [0, 1]$ which agrees with μ on the equivalence classes. The following is proved in Appendix B of the Supplementary Material [24] by an application of Gitik–Shelah theorem [16].

LEMMA 4.13. *Let \mathcal{X} be a set of an atomlessly-measurable cardinality κ and let μ be a witnessing measure on \mathfrak{X}_κ . Let (\mathcal{U}, Δ) be as above. Then there exist $\varepsilon > 0$ and $\mathcal{H}_\varepsilon \subseteq \mathcal{U}$ of cardinality $|\mathcal{H}_\varepsilon| = \kappa$ that is ε -separated:*

$$\Delta(U, V) \geq \varepsilon \quad \forall U, V \in \mathcal{H}_\varepsilon, U \neq V.$$

By Lemma 4.13, there exist an $\varepsilon > 0$ and a set $\mathcal{H}_\varepsilon \subseteq \{0, 1\}^\mathcal{X}$ such that $\forall g, h \in \mathcal{H}_\varepsilon$ with $g \neq h$, $\mu(\{x : g(x) \neq h(x)\}) \geq \varepsilon$ and, furthermore, \mathcal{H}_ε has cardinality κ_{\min} : that is, the same cardinality as \mathcal{X} . Since κ_{\min} is atomlessly-measurable, there exists an atomless witnessing measure π on $(\mathcal{H}_\varepsilon, 2^{\mathcal{H}_\varepsilon})$. We will construct the distribution $\bar{\mu}$ using a random construction, by fixing the marginal μ on \mathcal{X} and setting $\bar{\mu}$ to agree with the classifier h^* , which is π -distributed, independently of the input to the algorithm. This process is described formally below.

First, we introduce a relaxed objective for the learning algorithm Alg. Recall that given a labeled sample $S_n \in (\mathcal{X} \times \{0, 1\})^n$, Alg outputs a classifier $\hat{h}_{S_n} \in \{0, 1\}^\mathcal{X}$. For any sequence $S'_x := \{x'_1, x'_2, \dots\} \in \mathcal{X}$, $n \in \mathbb{N}$, and $\mathbf{y}' := (y'_1, \dots, y'_n) \in \{0, 1\}^n$, denote $\hat{h}_{S'_x, \mathbf{y}'} := \hat{h}_{\{(x'_1, y'_1), \dots, (x'_n, y'_n)\}}$ and let $H_{S'_x} := \{\hat{h}_{S'_x, \mathbf{y}'} : n \in \mathbb{N}, \mathbf{y}' \in \{0, 1\}^n\}$. This set may be random if the learning algorithm is randomized. Then note that, for any fixed $\bar{\mu}$, denoting by $S := \{(x_1, y_1), (x_2, y_2), \dots\}$ a countably-infinite sequence of independent $\bar{\mu}$ -distributed random variables, and further denoting $S_n := \{(x_1, y_1), \dots, (x_n, y_n)\}$ and $S_x := \{x_1, x_2, \dots\}$, we have

$$\inf_n \mathbb{E}[\text{err}_{\bar{\mu}}(\hat{h}_{S_n})] \geq \mathbb{E}\left[\inf_{h \in H_{S_x}} \text{err}_{\bar{\mu}}(h)\right].$$

Now take $S_x = \{x_1, x_2, \dots\}$ to be an i.i.d. μ -distributed sequence, and let $h^* \sim \pi$ independently of S_x . Let $\bar{\mu}$ have marginal μ over \mathcal{X} and define $\bar{\mu}$ such that $\bar{\mu}(\{(x, h^*(x)) : x \in \mathcal{X}\}) = 1$; that is, $\bar{\mu}$ is an h^* -dependent random measure. Note that $\text{err}_{\bar{\mu}}(h^*) = 0$ (a.s.), and hence also that any h has $\text{err}_{\bar{\mu}}(h) = \mu(\{x' : h(x') \neq h^*(x')\})$ (a.s.). Furthermore, by the assumed measurability of the learning algorithm, for each \mathbf{y} we have that $\hat{h}_{S_x, \mathbf{y}}$ is a $\mathcal{B}(L^1(\mu))$ -measurable random variable, and h^* is also $\mathcal{B}(L^1(\mu))$ -measurable (its distribution is π , which is defined on this σ -algebra). Therefore, $\mu(\{x' : \hat{h}_{S_x, \mathbf{y}}(x') \neq h^*(x')\})$ is a measurable random variable, equal (a.s.) to $\text{err}_{\bar{\mu}}(\hat{h}_{S_x, \mathbf{y}})$.

In particular, this implies that $\mathbb{E}[\inf_{h \in H_{S_x}} \text{err}_{\bar{\mu}}(h)]$ is well defined, and by the law of total expectation,

$$\begin{aligned} \mathbb{E}\left[\inf_{h \in H_{S_x}} \text{err}_{\bar{\mu}}(h)\right] &= \mathbb{E}\left[\mathbb{E}\left[\inf_{h \in H_{S_x}} \text{err}_{\bar{\mu}}(h) \mid H_{S_x}\right]\right] \\ &\geq \mathbb{E}\left[(\varepsilon/2)\mathbb{P}\left(\inf_{h \in H_{S_x}} \text{err}_{\bar{\mu}}(h) > \varepsilon/2 \mid H_{S_x}\right)\right] \\ &= \mathbb{E}\left[(\varepsilon/2)\pi\left(h' \in \mathcal{H}_\varepsilon : \inf_{h \in H_{S_x}} \mu(\{x' : h(x') \neq h'(x')\}) > \varepsilon/2\right)\right]. \end{aligned}$$

Then note that each element of H_{S_x} can be $(\varepsilon/2)$ -close to at most one element of \mathcal{H}_ε , and since H_{S_x} is a countable set, this implies that, given H_{S_x} , the set $H_{S_x}^\varepsilon = \{h' \in \mathcal{H}_\varepsilon : \inf_{h \in H_{S_x}} \mu(\{x' : h(x') \neq h'(x')\}) \leq \varepsilon/2\}$ is countable.

But since π vanishes on singletons, we have $\pi(H_{S_x}^\varepsilon) = 0$. Thus, given H_{S_x} ,

$$\pi\left(h' \in \mathcal{H}_\varepsilon : \inf_{h \in H_{S_x}} \mu(\{x' : h(x') \neq h'(x')\}) > \varepsilon/2\right) = 1,$$

so that altogether we have

$$\mathbb{E}\left[\inf_{h \in H_{S_x}} \text{err}_{\bar{\mu}}(h)\right] \geq \varepsilon/2.$$

In particular, this also implies there exist fixed choices of h^* for which (4.3) holds. This completes the proof for the case (I).

(II) *The case $\kappa_{\min} > c$.* By Ulam’s dichotomy in Theorem 4.12, $|\mathcal{X}| = \kappa_{\min}$ is two-valued measurable, so there exists a two-valued witnessing measure μ on \mathcal{B} . As the following lemma shows, μ can be taken to further satisfy a key homogeneity property. The lemma is proved in Appendix C of the Supplementary Material [24], where it is shown to follow by combining Theorems 10.20, 10.22 in [26] and Ulam’s theorem 4.12.

LEMMA 4.14. *Let \mathcal{X} be of a two-valued measurable cardinality κ and let $\mathfrak{X}_\kappa = (\mathcal{X}, 2^\mathcal{X})$. Then there is a witnessing measure μ on \mathfrak{X}_κ such that for any function $f : [\mathcal{X}]^{<\omega} \rightarrow \mathbb{R}$, there exists a $U \subseteq \mathcal{X}$ with $\mu(U) = 1$ such that U is homogeneous for f , that is, for every $n \in \mathbb{N}$, there exists a $C_n \in \mathbb{R}$ such that $f(W) = C_n$ for all $W \in [U]^n$.*

Let μ be a two-valued witnessing measure on $\mathcal{B} = 2^\mathcal{X}$ as furnished by Lemma 4.14. For a label $y \in \{0, 1\}$ and any two-valued witnessing measure ϕ , let $\bar{\phi}_y$ be the measure over $\mathcal{X} \times \{0, 1\}$ with ϕ as its marginal over \mathcal{X} and

$$\bar{\phi}_y(Y = y \mid X = x) = 1 \quad \forall x \in \mathcal{X}.$$

For μ as above, and any other two-valued witnessing measure ϕ , define

$$(4.4) \quad \lambda_\phi := \frac{2}{3}\bar{\phi}_1 + \frac{1}{3}\bar{\mu}_0.$$

We will show that there exists a two-valued witnessing measure $\nu := \nu(\mu, \text{Alg}) \neq \mu$ such that Alg cannot be Bayes-consistent on both λ_μ and λ_ν . To this end, we will use the following properties of the mixture λ_ϕ , proved in Appendix C of the Supplementary Material [24].

LEMMA 4.15. *Let $\nu \neq \mu$ be any two distinct two-valued measures on $(\mathcal{X}, \mathcal{B})$ and let λ_ϕ with $\phi \in \{\mu, \nu\}$ be as in (4.4).*

(i) Any Bayes-optimal classifier h^* on λ_μ achieves the optimal Bayes-error $\text{err}_{\lambda_\mu}(h^*) = \frac{1}{3}$ if and only if $\mathbb{E}_{X \sim \mu}[h^*(X)] = 1$.

(ii) Any Bayes-optimal classifier h^* on λ_ν achieves the optimal Bayes-error $\text{err}_{\lambda_\nu}(h^*) = 0$ if and only if $\mathbb{E}_{X \sim \mu}[h^*(X)] = 0$ and $\mathbb{E}_{X \sim \nu}[h^*(X)] = 1$.

Let $\text{Alg} : (\mathcal{X} \times \{0, 1\})^{<\omega} \rightarrow 2^{\mathcal{X}}$ be any (possibly randomized) learning algorithm, and recall that \hat{h}_S denotes the classifier output for data set S ; for S and X independent samples from Borel measures on \mathcal{X} , we suppose that $\hat{h}_S(X)$ is a measurable random variable (by definition of learning algorithm; see Remark 4.10). Let $\nu \neq \mu$ be a two-valued witnessing measure to be chosen below. Consider the quantity

$$Z_n^\phi := \mathbb{E}_{S_n \sim (\lambda_\phi)^n} [\mathbb{E}_{X \sim \mu} [\hat{h}_{S_n}(X)]], \quad \phi \in \{\mu, \nu\}.$$

In the case of a randomized Alg, also add an innermost expectation over the independent randomness of Alg in the above expression. By Lemma 4.15, for Alg to be Bayes consistent on both λ_ν and λ_μ we must have

$$(4.5) \quad Z_n^\phi \xrightarrow[n \rightarrow \infty]{} \delta_{\mu, \phi}, \quad \phi \in \{\mu, \nu\},$$

where $\delta_{\mu, \phi}$ is the Kronecker delta. So to prove the claim it suffices to show that we can choose $\nu := \nu(\mu, \text{Alg})$ such that (4.5) does not hold.

Given a labeled sample $S_n = (\mathbf{X}_n, \mathbf{Y}_n) \sim (\lambda_\phi)^n$ with $\phi \in \{\mu, \nu\}$, let $n_1 := n_1(\mathbf{Y}_n) = \sum_{i=1}^n Y_i$ and $n_0 := n - n_1$ be the random number of samples in S_n with labels 1 and 0, respectively, and let $\mathbf{X}_n^0 \in \mathcal{X}^{n_0}$ and $\mathbf{X}_n^1 \in \mathcal{X}^{n_1}$ be the corresponding instances in \mathbf{X}_n . For notational simplicity, we write $\mathbf{X}_n = (\mathbf{X}_n^0, \mathbf{X}_n^1)$ where it is understood that the embedding of \mathbf{X}_n^0 and \mathbf{X}_n^1 in \mathbf{X}_n is in accordance with \mathbf{Y}_n . Note that $\mathbf{Y}_n \sim (\text{Bernoulli}(\frac{2}{3}))^n$ irrespectively of μ and ϕ . In addition, given \mathbf{Y}_n we have that \mathbf{X}_n^0 and \mathbf{X}_n^1 are independent and $\mathbf{X}_n^0 | \mathbf{Y}_n \sim \mu^{n_0}$ and $\mathbf{X}_n^1 | \mathbf{Y}_n \sim \phi^{n_1}$. We decompose

$$(4.6) \quad \begin{aligned} Z_n^\phi &= \mathbb{E}_{S_n \sim (\lambda_\phi)^n} [\mathbb{E}_{X \sim \mu} [\hat{h}_{S_n}(X)]] \\ &= \mathbb{E}_{\mathbf{Y}_n} \mathbb{E}_{\mathbf{X}_n | \mathbf{Y}_n} [\mathbb{E}_{X \sim \mu} [\hat{h}_{(\mathbf{X}_n, \mathbf{Y}_n)}(X)]] \\ &= \mathbb{E}_{\mathbf{Y}_n} \mathbb{E}_{\mathbf{X}_n^1 \sim \phi^{n_1}} \mathbb{E}_{\mathbf{X}_n^0 \sim \mu^{n_0}} [\mathbb{E}_{X \sim \mu} [\hat{h}_{((\mathbf{X}_n^0, \mathbf{X}_n^1), \mathbf{Y}_n)}(X)]] \end{aligned}$$

Toward applying Lemma 4.14, we first need to translate our reasoning about a random vector $\mathbf{X} = (X_1, \dots, X_k) \sim \phi^k$ with $k \in \mathbb{N}$ into reasoning about the random set of its distinct elements, $W_{\mathbf{X}} := \bigcup_{i=1}^k \{X_i\}$. Since ϕ vanishes on singletons, all instances in \mathbf{X} are distinct with probability one,

$$(4.7) \quad \mathbb{P}_{\mathbf{X} \sim \phi^k} [|W_{\mathbf{X}}| = k] = 1.$$

Fixing an ordering on \mathcal{X} , for any finite set $W = \{w_1, \dots, w_k\} \in [\mathcal{X}]^k$, denote by $\Pi(W)$ the distribution over vectors $\mathbf{X}' = (w_{\pi(1)}, \dots, w_{\pi(k)}) \in W^k$ as induced by a random permutation π of the instances in W . Then, by (4.7) and the fact that ϕ^k is a product measure, we have that for any measurable function $f : \mathcal{X}^k \rightarrow [0, 1]$ the following symmetrization holds:

$$(4.8) \quad \mathbb{E}_{\mathbf{X} \sim \phi^k} [f(\mathbf{X})] = \mathbb{E}_{\mathbf{X} \sim \phi^k} [\mathbb{E}_{\mathbf{X}' \sim \Pi(W_{\mathbf{X}})} [f(\mathbf{X}')] \mid |W_{\mathbf{X}}| = k].$$

For every $\mathbf{Y}_n \in \{0, 1\}^n$ define $F_{\mathbf{Y}_n} : [\mathcal{X}]^{n_1} \rightarrow \mathbb{R}$ by

$$F_{\mathbf{Y}_n}(W) = \mathbb{E}_{\mathbf{X}^1 \sim \Pi(W)} \mathbb{E}_{\mathbf{X}^0 \sim \mu^{n_0}} [\mathbb{E}_{X \sim \mu} [\hat{h}_{((\mathbf{X}^0, \mathbf{X}^1), \mathbf{Y}_n)}(X)]], \quad W \in [\mathcal{X}]^{n_1}.$$

In the case of randomized Alg, we also include an innermost conditional expectation over the value of $\hat{h}_{((\mathbf{X}_n^0, \mathbf{X}_n^1), \mathbf{Y}_n)}(X)$ given $\mathbf{X}_n^0, \mathbf{X}_n^1, \mathbf{Y}_n, X$. Putting this in (4.6) while using (4.7) and (4.8),

$$Z_n^\phi = \mathbb{E}_{\mathbf{Y}_n} \mathbb{E}_{\mathbf{X} \sim \phi^{n_1}} [F_{\mathbf{Y}_n}(W_{\mathbf{X}}) \mid |W_{\mathbf{X}}| = n_1].$$

By the choice of μ , Lemma 4.14 implies there exist $C_{\mathbf{Y}_n} \in \mathbb{R}$ and $U_{\mathbf{Y}_n} \subseteq \mathcal{X}$ with $\mu(U_{\mathbf{Y}_n}) = 1$ such that $U_{\mathbf{Y}_n}$ is homogeneous for $F_{\mathbf{Y}_n}$, namely, $F_{\mathbf{Y}_n}(W) = C_{\mathbf{Y}_n}, \forall W \in [U_{\mathbf{Y}_n}]^{n_1}$. Let

$$(4.9) \quad U = \bigcap_{n \in \mathbb{N}} \bigcap_{\mathbf{Y}_n \in \{0,1\}^n} U_{\mathbf{Y}_n}.$$

Then U is simultaneously homogeneous for all $\{F_{\mathbf{Y}_n}\}$,

$$(4.10) \quad F_{\mathbf{Y}_n}(W) = C_{\mathbf{Y}_n} \quad \forall n \in \mathbb{N}, \forall \mathbf{Y}_n \in \{0,1\}^n, \forall W \in [U]^{n_1}.$$

In addition, by Lemma C.1 in the Supplementary Material [24], $\mu(U) = 1$.

We are now in position to choose $\nu := \nu(\mu, \text{Alg})$. By Lemma 4.4, we may split U in (4.9) into two disjoint sets B and $U \setminus B$ such that $|B| = |U \setminus B| = |\mathcal{X}|$. Since μ is two-valued, we may assume without loss of generality that $\mu(B) = 0$ (so $\mu(U \setminus B) = 1$). Since $|B|$ is a two-valued measurable cardinal, there exists a two-valued witnessing measure ν' on $(B, 2^B)$ with $\nu'(B) = 1$. Extend ν' to a measure ν over all \mathcal{B} by $\nu(A) = \nu'(A \cap B)$, $\forall A \subseteq \mathcal{X}$. Then $\nu \neq \mu$ and $\nu(U) = \mu(U) = 1$. By the last equality, for $\phi \in \{\mu, \nu\}$ and $\forall k \in \mathbb{N}$, $\Pr_{\mathbf{X} \sim \phi^k} [W_{\mathbf{X}} \in [U]^k \mid |W_{\mathbf{X}}| = k] = 1$. So, for $\phi \in \{\mu, \nu\}$,

$$\begin{aligned} Z_n^\phi &= \mathbb{E}_{\mathbf{Y}_n} \mathbb{E}_{\mathbf{X} \sim \phi^{n_1}} [F_{\mathbf{Y}_n}(W_{\mathbf{X}}) \mid |W_{\mathbf{X}}| = n_1] \\ &= \mathbb{E}_{\mathbf{Y}_n} \mathbb{E}_{\mathbf{X} \sim \phi^{n_1}} [F_{\mathbf{Y}_n}(W_{\mathbf{X}}) \mid |W_{\mathbf{X}}| = n_1 \wedge W_{\mathbf{X}} \in [U]^{n_1}] \\ &= \mathbb{E}_{\mathbf{Y}_n} \mathbb{E}_{\mathbf{X} \sim \phi^{n_1}} [C_{\mathbf{Y}_n} \mid |W_{\mathbf{X}}| = n_1 \wedge W_{\mathbf{X}} \in [U]^{n_1}] \\ &= \mathbb{E}_{\mathbf{Y}_n} [C_{\mathbf{Y}_n}], \end{aligned}$$

where we used (4.10) and the fact that $C_{\mathbf{Y}_n}$ does not depend on \mathbf{X} . Since $\mathbb{E}_{\mathbf{Y}_n} [C_{\mathbf{Y}_n}]$ is independent of ϕ , we conclude that $Z_n^\mu = Z_n^\nu$ for all $n \in \mathbb{N}$. However by (4.5), for Alg to be Bayes consistent on λ_μ and λ_ν we must have $Z_n^\mu \xrightarrow{n \rightarrow \infty} 1$ and $Z_n^\nu \xrightarrow{n \rightarrow \infty} 0$. Thus Alg cannot be Bayes consistent on both λ_μ and λ_ν . In particular, (4.3) holds with $\varepsilon = 1/4$.

5. Discussion. We have exhibited a computationally efficient multiclass learning algorithm, OptiNet, that is universally strongly Bayes consistent (UBC) in all essentially separable (ES) metric spaces. In contrast, we showed that in non-ES spaces, no algorithm can be UBC. As such, OptiNet is optimistically universal (in the terminology of [22])—it is universally Bayes consistent in all metric spaces that admit such a learner. We note that in this work, we do not study the rates of decay of the excess risk, leaving this challenging open problem for future study.

By definition, any separable metric space is ES. As discussed in Section 1, consistency of NN-type algorithms in general separable metric spaces was studied in [1, 3, 4, 6, 13, 33]. In particular, in [1, 6, 13], a characterization of the metric spaces in which an algorithm is universally Bayes consistent was given for several such algorithms, in terms of Besicovitch-type conditions. As a notable example, it is shown in [6] that for any separable metric space \mathcal{X} , a sufficient condition for the k -NN algorithm (with an appropriate choice of the number of neighbors k) to be Bayes consistent for a distribution $\bar{\mu}$ over $\mathcal{X} \times \{0, 1\}$ is that for all $\varepsilon > 0$,

$$(5.1) \quad \lim_{r \rightarrow 0^+} \mathbb{P} \left\{ \frac{1}{\mu(B_r(X))} \int_{B_r(X)} |\eta(z) - \eta(X)| d\mu(z) > \varepsilon \right\} = 0,$$

where μ is the marginal of $\bar{\mu}$ over \mathcal{X} and $\eta(x) := \mathbb{P}(Y = 1 \mid X = x)$. It is also shown in [6] that in the realizable case, where $\eta(x) \in \{0, 1\}$ for all $x \in X$, a violation of (5.1) implies that k -NN is inconsistent. Say that a metric space satisfies the *universal Besicovitch condition* if (5.1) holds for all measures $\bar{\mu}$ over the Borel σ -algebra. By Besicovitch’s density

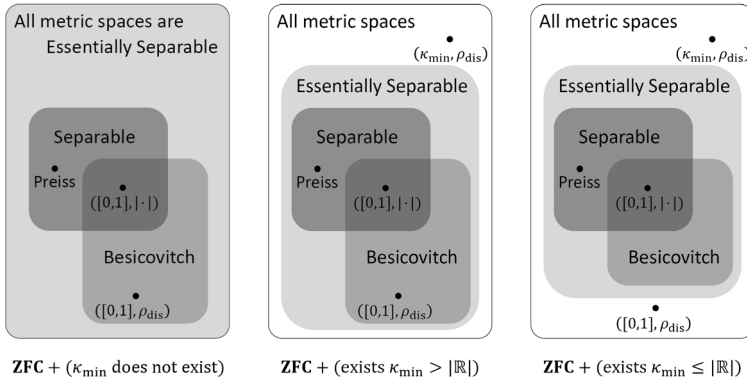


FIG. 1. The classes of metric spaces discussed in this paper and their inclusion relationships in the three cases where: no RVMC exists (left); minimal RVMC is $\kappa_{\min} > \mathfrak{c}$ (middle); and minimal RVMC is $\kappa_{\min} \leq \mathfrak{c}$ (right). All three cases are believed to be valid extensions of ZFC. The metric space $(\kappa_{\min}, \rho_{\text{dis}})$ corresponds to a discrete one of cardinality κ_{\min} ; it is not ES but any discrete metric space of cardinality $< \kappa_{\min}$ is ES. The shaded area named “Besicovitch” and the specific metric space “Preiss” are as discussed in the text.

theorem [15], Section 472, the metric space $(\mathbb{R}^d, \|\cdot\|_2)$ —and more generally, any finite-dimensional normed space—satisfies this condition, so k -NN is UBC on such spaces. In contrast, in infinite-dimensional separable spaces, such as ℓ_2 , a violation of (5.1) can occur [36, 37, 40]. One such example is the separable metric probability space studied in [28], building upon a construction of Preiss [36]. While the k -NN algorithm is provably not UBC in this space, OptiNet is. As far as we know, OptiNet is the first algorithm known to be UBC (weakly or strongly) in any separable metric space.

As discussed in Section 4.1, the essential separability of nonseparable metric spaces is believed to depend on set-theoretic axioms that are independent of ZFC, and in particular on the cardinality of the minimal RVMC, κ_{\min} : a metric space is non-ES if and only if it contains a discrete subset of cardinality κ_{\min} . Figure 1 gives a pictorial illustration of the possible relationships between the following types of metric spaces: separable, (uniform) Besicovitch, ES, and all spaces, depending on the set-theoretic model. If one adopts a model in which no RVMC exist, then any discrete subspace of a metric space admits only trivial, purely-atomic measures. In this case, abbreviated as ZFC + (κ_{\min} does not exist) in the left panel of Figure 1, all metric spaces are ES, and OptiNet is UBC on any metric space. Alternatively, if one adopts a set-theoretic model in which an RVMC exists, then discrete subspaces of \mathcal{X} of cardinality $\geq \kappa_{\min}$ admit also nontrivial measures. As shown in Section 4.2, such measures exclude the possibility of a UBC algorithm. The nature of the nontrivial measures, being atomless or purely atomic, depends on whether $\kappa_{\min} > \mathfrak{c}$ or $\kappa_{\min} \leq \mathfrak{c}$, which are illustrated on the middle and right panels of Figure 1, respectively.

Lastly, we note that our argument for the impossibility of UBC in non-ES metric spaces is based solely on the real-valued measurability of the cardinality of discrete subspaces of \mathcal{X} . This raises a natural question: Assuming no cardinal is real-valued measurable, are there any topological spaces (which by the results above must be nonmetric) in which no UBC algorithm exists?

To summarize, in this work we provided the first multiclass learning algorithm that is universally Bayes consistent in any metric space where such an algorithm exists. Moreover, we provided a characterization of these metric spaces. The study of learnability in general spaces is fundamental, and provides many open questions for future research.

Acknowledgments. We thank Vladimir Pestov for sharing with us his proof of the existence of a measurable total order. We also thank Robert Furber, Iosif Pinelis, Menachem Kojman, and Roberto Colomboni for helpful discussions.

Funding. Aryeh Kontorovich was supported in part by the Israel Science Foundation (Grant No. 755/15), Paypal and IBM. Sivan Sabato was supported in part by the Israel Science Foundation Grant No. 555/15.

SUPPLEMENTARY MATERIAL

Supplementary material (DOI: [10.1214/20-AOS2029SUPP](https://doi.org/10.1214/20-AOS2029SUPP); .pdf). In the Supplementary Material file, we provide the proofs of Lemmas 3.6, 3.7, 3.8, 4.13, 4.14, 4.15, as well as the proofs of the auxiliary Lemmas A.1, B.1, C.1, C.2, D.1.

REFERENCES

- [1] ABRAHAM, C., BIAU, G. and CADRE, B. (2006). On the kernel rule for function classification. *Ann. Inst. Statist. Math.* **58** 619–633. [MR2327897 https://doi.org/10.1007/s10463-006-0032-1](https://doi.org/10.1007/s10463-006-0032-1)
- [2] BEN-DAVID, S., HRUBES, P., MORAN, S., SHPILKA, A. and YEHUDAYOFF, A. (2019). Learnability can be undecidable. *Nat. Mach. Intell.* **1** 44–48. <https://doi.org/10.1038/s42256-018-0002-3>
- [3] BIAU, G., BUNEA, F. and WEGKAMP, M. H. (2005). Functional classification in Hilbert spaces. *IEEE Trans. Inf. Theory* **51** 2163–2172. [MR2235289 https://doi.org/10.1109/TIT.2005.847705](https://doi.org/10.1109/TIT.2005.847705)
- [4] BIAU, G., CÉROU, F. and GUYADER, A. (2010). Rates of convergence of the functional k -nearest neighbor estimate. *IEEE Trans. Inf. Theory* **56** 2034–2040. [MR2654492 https://doi.org/10.1109/TIT.2010.2040857](https://doi.org/10.1109/TIT.2010.2040857)
- [5] BILLINGSLEY, P. (1968). *Convergence of Probability Measures*, 2nd ed. *Wiley Series in Probability and Statistics: Probability and Statistics*. Wiley, New York. A Wiley-Interscience Publication. [MR1700749 https://doi.org/10.1002/9780470316962](https://doi.org/10.1002/9780470316962)
- [6] CÉROU, F. and GUYADER, A. (2006). Nearest neighbor classification in infinite dimension. *ESAIM Probab. Stat.* **10** 340–355. [MR2247925 https://doi.org/10.1051/ps:2006014](https://doi.org/10.1051/ps:2006014)
- [7] CHAUDHURI, K. and DASGUPTA, S. (2014). Rates of convergence for nearest neighbor classification. In *Advances in Neural Information Processing Systems* 3437–3445.
- [8] CHRISTMANN, A. and STEINWART, I. (2010). Universal kernels on non-standard input spaces. In *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems* 2010 406–414.
- [9] COVER, T. M. and HART, P. E. (1967). Nearest neighbor pattern classification. *IEEE Trans. Inf. Theory* **13** 21–27.
- [10] DEVROYE, L., GYÖRFI, L. and LUGOSI, G. (1996). *A Probabilistic Theory of Pattern Recognition. Applications of Mathematics (New York)* **31**. Springer, New York. [MR1383093 https://doi.org/10.1007/978-1-4612-0711-5](https://doi.org/10.1007/978-1-4612-0711-5)
- [11] DUDLEY, R. M. (1999). *Uniform Central Limit Theorems. Cambridge Studies in Advanced Mathematics* **63**. Cambridge Univ. Press, Cambridge. [MR1720712 https://doi.org/10.1017/CBO9780511665622](https://doi.org/10.1017/CBO9780511665622)
- [12] FIX, E. and HODGES J. L., J. (1989). Discriminatory analysis. Nonparametric discrimination: Consistency properties. *Int. Stat. Rev.* **57** 238–247.
- [13] FORZANI, L., FRAIMAN, R. and LLOP, P. (2012). Consistent nonparametric regression for functional data under the Stone–Besicovitch conditions. *IEEE Trans. Inf. Theory* **58** 6697–6708. [MR2991802 https://doi.org/10.1109/TIT.2012.2209628](https://doi.org/10.1109/TIT.2012.2209628)
- [14] FREMLIN, D. H. (1993). Real-valued-measurable cardinals. In *Set Theory of the Reals (Ramat Gan, 1991)*. *Israel Math. Conf. Proc.* **6** 151–304. Bar-Ilan Univ., Ramat Gan. [MR1234282](https://doi.org/10.1007/978-1-4612-0711-5)
- [15] FREMLIN, D. H. (2000). *Measure Theory* **1–5**. Torres Fremlin, Colchester.
- [16] GITIK, M. and SHELAH, S. (1989). Forcings with ideals and simple forcing notions. *Israel J. Math.* **68** 129–160. [MR1035887 https://doi.org/10.1007/BF02772658](https://doi.org/10.1007/BF02772658)
- [17] GOTTLIEB, L.-A., KONTOROVICH, A. and KRAUTHGAMER, R. (2014). Efficient classification for metric data. *IEEE Trans. Inf. Theory* **60** 5750–5759. [MR3252418 https://doi.org/10.1109/TIT.2014.2339840](https://doi.org/10.1109/TIT.2014.2339840)
- [18] GOTTLIEB, L.-A., KONTOROVICH, A. and NISNEVITCH, P. (2017). Nearly optimal classification for semi-metrics. *J. Mach. Learn. Res.* **18** Paper No. 37, 22. [MR3646632](https://doi.org/10.1007/978-1-4939-9831-7_37)
- [19] GOTTLIEB, L.-A., KONTOROVICH, A. and NISNEVITCH, P. (2018). Near-optimal sample compression for nearest neighbors. *IEEE Trans. Inf. Theory* **64** 4120–4128. [MR3809730 https://doi.org/10.1109/TIT.2018.2822267](https://doi.org/10.1109/TIT.2018.2822267)
- [20] GRAEPEL, T., HERBRICH, R. and SHAWE-TAYLOR, J. (2005). PAC-Bayesian compression bounds on the prediction error of learning algorithms for classification. *Mach. Learn.* **59** 55–76.

- [21] GYÖRFI, L., KOHLER, M., KRZYŻAK, A. and WALK, H. (2002). *A Distribution-Free Theory of Nonparametric Regression*. *Springer Series in Statistics*. Springer, New York. MR1920390 <https://doi.org/10.1007/b97848>
- [22] HANNEKE, S. (2017). Learning whenever learning is possible: Universal learning under general stochastic processes. Preprint. Available at [arXiv:1706.01418](https://arxiv.org/abs/1706.01418).
- [23] HANNEKE, S. and KONTOROVICH, A. (2019). A sharp lower bound for agnostic learning with sample compression schemes. In *Algorithmic Learning Theory 2019. Proc. Mach. Learn. Res. (PMLR)* **98** 489–505. MR3932856
- [24] HANNEKE, S., KONTOROVICH, A., SABATO, S. and WEISS, R. (2021). Supplement to “Universal Bayes consistency in metric spaces.” <https://doi.org/10.1214/20-AOS2029SUPP>
- [25] HRBACEK, K. and JECH, T. (1999). *Introduction to Set Theory*, 3rd ed. *Monographs and Textbooks in Pure and Applied Mathematics* **220**. Dekker, New York. MR1697766
- [26] JECH, T. (2003). *Set Theory*. *Springer Monographs in Mathematics*. Springer, Berlin. The third millennium edition, revised and expanded. MR1940513
- [27] KONTOROVICH, A., SABATO, S. and URNER, R. (2017). Active nearest-neighbor learning in metric spaces. *J. Mach. Learn. Res.* **18** Paper No. 195, 38. MR3827083
- [28] KONTOROVICH, A., SABATO, S. and WEISS, R. (2017). Nearest-neighbor sample compression: Efficiency, consistency, infinite dimensions. In *Advances in Neural Information Processing Systems* 1573–1583.
- [29] KONTOROVICH, A., SABATO, S. and WEISS, R. (2017). Nearest-neighbor sample compression: Efficiency, consistency, infinite dimensions. Preprint. Available at [arXiv:1705.08184](https://arxiv.org/abs/1705.08184).
- [30] KONTOROVICH, A. and WEISS, R. (2014). Maximum margin multiclass nearest neighbors. In *International Conference on Machine Learning (ICML 2014)*.
- [31] KONTOROVICH, A. and WEISS, R. (2014). A Bayes consistent 1-NN classifier. In *Artificial Intelligence and Statistics (AISTATS 2015)*.
- [32] KPOTUFE, S. and VERMA, N. (2017). Time-accuracy tradeoffs in kernel prediction: Controlling prediction quality. *J. Mach. Learn. Res.* **18** Paper No. 44, 29. MR3655309
- [33] KULKARNI, S. R. and POSNER, S. E. (1995). Rates of convergence of nearest neighbor estimation under arbitrary sampling. *IEEE Trans. Inf. Theory* **41** 1028–1039. MR1366756 <https://doi.org/10.1109/18.391248>
- [34] LITTLESTONE, N. and WARMUTH, M. K. (1986). Relating data compression and learnability. Unpublished.
- [35] PELILLO, M. (2014). Alhazen and the nearest neighbor rule. *Pattern Recogn. Lett.* **38** 34–37. <https://doi.org/10.1016/j.patrec.2013.10.022>
- [36] PREISS, D. (1979). Invalid Vitali theorems. In *Abstracts. 7th Winter School on Abstract Analysis* 58–60.
- [37] PREISS, D. (1981). Gaussian measures and the density theorem. *Comment. Math. Univ. Carolin.* **22** 181–193. MR0609946
- [38] SHALEV-SHWARTZ, S. and BEN-DAVID, S. (2014). *Understanding Machine Learning: From Theory to Algorithms*. Cambridge Univ. Press, Cambridge.
- [39] STONE, C. J. (1977). Consistent nonparametric regression. *Ann. Statist.* **5** 595–645. MR0443204
- [40] TIŠER, J. (2003). Vitali covering theorem in Hilbert space. *Trans. Amer. Math. Soc.* **355** 3277–3289. MR1974687 <https://doi.org/10.1090/S0002-9947-03-03296-3>
- [41] ULAM, S. M. (1930). *Zur Masstheorie in der Allgemeinen Mengenlehre*. Uniwersytet, seminarjum matematyczne.
- [42] ZHAO, L. C. (1987). Exponential bounds of mean error for the nearest neighbor estimates of regression functions. *J. Multivariate Anal.* **21** 168–178. MR0877849 [https://doi.org/10.1016/0047-259X\(87\)90105-9](https://doi.org/10.1016/0047-259X(87)90105-9)