

# UNIVERSAL BAYES CONSISTENCY IN METRIC SPACES

BY STEVE HANNEKE\* AND ARYEH KONTOROVICH<sup>†</sup>  
AND SIVAN SABATO<sup>†</sup> AND ROI WEISS<sup>‡</sup>

*Toyota Technological Institute at Chicago\**  
*and Ben-Gurion University of the Negev<sup>†</sup>*  
*and Weizmann Institute of Science<sup>‡</sup>*

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

The metric spaces in which universal Bayes consistency is possible are the “essentially separable” ones — a new notion that we define, which is more general than standard separability. The existence of metric spaces that are not essentially separable is 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, these positive and negative results resolve the open problems posed in Kontorovich, Sabato, Weiss (2017).

**1. Introduction.** Since their inception, nearest-neighbor methods for classification and regression in metric spaces continue to attract the interest of theoreticians and practitioners alike. On the applied front, this seemingly naive approach remains competitive against more sophisticated methods [9, 59, 44, 12]. On the theoretical front, the most commonly investigated questions involve Bayes consistency and rates of convergence [29, 40, 21, 11]. In particular, one might seek to characterize all metric spaces in which *any* universal Bayes-consistent classification algorithm exists. One can then ask whether there is a *single* algorithm that achieves universal Bayes consistency in all metric spaces for which such an algorithm exists, and in particular, whether there exists such an algorithm which generates nearest-neighbor classification rules. For the problem of (multiclass) classification, we resolve these questions exhaustively.

*Main results.* We study the multiclass learning algorithm recently proposed by [35], which we call OptiNet (see Algorithm 1). OptiNet is based on a 1-NN rule derived from sample compression bounds, and enjoys the statistical advantages of tight, fully empirical generalization bounds, as well as the algorithmic advantages

---

*AMS 2010 subject classifications:* Primary 54E70, 97K80, 62C12; secondary 03E17, 03E55  
*Keywords and phrases:* metric space, nearest neighbor, classification, Bayes consistency

of a fast runtime and memory savings. In this work, we show that **OptiNet** is universally strongly Bayes consistent in all essentially separable metric spaces. Here, an *essentially separable* space is a new notion that we define in Section 4, which is weaker than standard separability.

We further show that the existence of metric spaces that are not essentially separable hinges upon set-theoretic axioms, independent of ZFC,<sup>1</sup> which relate to the existence of certain measurable cardinals.

Our results imply a dichotomy: If one’s model of set theory allows for metric spaces that are not essentially separable, then such spaces do not admit *any* universally Bayes-consistent learner. To our knowledge, this is the first construction of a learning setting in which universal Bayes consistency is impossible. However, 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 the sense of [30],<sup>2</sup> 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 doubling dimension, or a Besicovitch-type property [10, 6, 7], all of which are significantly stronger assumptions than essential separability.

Taken together, the results demonstrating universal Bayes consistency in essentially separable spaces and its impossibility in spaces which are not essentially separable resolve the open problems posed in [36].

*Related work.* Nearest-neighbor methods were initiated by Fix and Hodges<sup>3</sup> in 1951 [18], and, in the celebrated  $k$ -NN formulation, have been placed on a solid theoretical foundation [14, 56, 16, 60, 11]. Following the pioneering work of [14, 56] on nearest-neighbor classification, it was shown by [60, 16, 27] 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 [54] a weak Bayes-consistency result was shown for metric spaces with a bounded diameter and a bounded doubling dimension, and additional distributional smoothness assumptions. More recently, some of the classic results on the decay rates of  $k$ -NN risk were refined by [11], in an analysis that captures the interplay between the metric and the sampling distribution. The worst-case rates have an exponential dependence on the dimension (i.e., the so-called *curse of dimensionality*), and Pestov [47, 48] examines this phenomenon closely under various distributional and structural assumptions.

Consistency of NN-type algorithms in more general (and, in particular, infinite-dimensional) metric spaces was discussed in [1, 6, 7, 10, 42, 19]. In [1, 10, 19], characterizations of Bayes consistency of such algorithms (including the standard  $k$ -NN) were given in terms of a Besicovitch-type condition. By Besicovitch’s den-

<sup>1</sup>We provide below the relevant set-theoretic background.

<sup>2</sup>This terminology of optimistic learning was introduced in [30], which applies the principle in the context of removing the iid assumption, exhibiting a learning rule which is optimistically consistent for function learning, in the sense that it asymptotically learns any function under any process  $\{X_i\}_{i \in \mathbb{N}}$  for which it is possible to do so.

<sup>3</sup>Pelillo [46] traces the 1-NN method to Abū ‘Alī al-Ḥasan ibn al-Ḥasan ibn al-Haytham (Alhazen)’s book *Kitāb al-Manāẓir* (“Book of Optics”), circa 1030.

sity theorem [20], 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 a violation can occur [49, 50]. Moreover, this is not an isolated pathology, as this violation is shared by Gaussian Hilbert spaces [57]. In [1], a generalized “moving window” classification rule is used along with additional regularity conditions on the regression function. The *filtering* technique (i.e., taking the first  $d$  coordinates in some basis representation) was shown to be universally consistent in [6]. However, that technique is only applicable in Hilbert spaces, as opposed to more general metric spaces. The insight of [6] was extended to the more general Banach spaces in [7] under various regularity assumptions. For compact metric spaces, the SVM algorithm can be made universally Bayes consistent by using an appropriate kernel [13].

None of the aforementioned generalization results for proximity-based techniques are in the form of fully empirical, explicitly computable sample-dependent error bounds. Rather, they are stated in terms of the unknown Bayes-optimal risk, and some involve additional parameters quantifying the well-behavedness of the unknown distribution (see [38] for a detailed discussion). As such, these guarantees do not enable a practitioner to compute a numerical generalization error estimate for a given training sample, much less allow for a data-dependent selection of  $k$ , which must be tuned via cross-validation. The asymptotic expansions in [55, 51, 28, 52] likewise do not provide a computable finite-sample bound. The quest for such bounds was a key motivation behind the series of works [23, 39, 25, 35].

Although the classic 1-NN classifier is well-known to be inconsistent in general, in recent years a series of papers has presented variations on the theme of a *regularized* 1-NN classifier, as an alternative to  $k$ -NN. Gottlieb et al. [23] showed that 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 [39]. In a follow-up work, [38] showed that applying Structural Risk Minimization (SRM) to a margin-regularized data-dependent bound very similar to that in [23] yields a strongly Bayes consistent 1-NN classifier in doubling spaces with bounded diameter.

Approaching the problem using sample compression techniques, a computationally near-optimal nearest neighbor condensing algorithm was presented [25] and later extended to cover semimetric spaces [24]; both were based on constructing  $\gamma$ -nets in spaces with a finite doubling dimension (or its semimetric analogue). As detailed in [38], 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.

The work of Devroye et al. [16, 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. It is shown 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 [43, 16], and later elaborated upon by [26, 31]. In the context of NN classification, [16] 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 [23, 39] relied on computing a minimum vertex cover. Thus, it was not possible to make it simultaneously 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 [45], this presents an obstruction to establishing the Bayes consistency of the classifier.

In [35], 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  $\gamma$ -net and take the majority (more accurately, *plurality*) label 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. A direct precursor to the present work, [36], showed that the passive component of the active learner in [35], which was provisionally termed there KSU, is universally strongly Bayes consistent in doubling metric spaces with bounded diameter; prior to that work, there were no efficient compression-based algorithms known to be Bayes consistent. They additionally gave an example of a non-doubling probability space (that violates the Besicovitch condition) where this algorithm also succeeds, and left it as an open question whether there is *any* metric probability space where it fails to be Bayes consistent. The present paper resolves the open problems posed in [36]: We establish universal Bayes consistency in all metric spaces for which a universal Bayes-consistent learner exists, and also provide a precise characterization of which metric spaces admit the existence of universally Bayes-consistent learners.

*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. Its consistency on essentially separable metric spaces is proved in Section 4. In Section 5 we prove that no universally Bayes-consistent algorithm exists on metric spaces that are not essentially separable. We conclude with a discussion in Section 6.

**2. Definitions and Notation.** Our *instance space* is the metric probability space  $(\mathcal{X}, \rho, \mu)$ , where  $\rho$  is a metric and  $\mu$  is a probability measure. By definition, the Borel  $\sigma$ -algebra  $\mathcal{B}$  supporting  $\mu$  is the smallest  $\sigma$ -algebra containing the open sets of  $\rho$ . 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  $\rho$ :

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

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  $\mu$  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

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

We omit the overline in  $\bar{\mu}$  when there is no ambiguity. 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  $\bar{\mu}$  and denote the optimal Bayes risk of  $\bar{\mu}$  by  $R^*$  when there is no ambiguity.

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 — i.e., the smallest  $x_i$  is chosen, with respect to a fixed total ordering of the space  $\mathcal{X}$ . The 1-NN classifier induced by  $S$  is denoted by  $h_S(x) := Y_{\text{nn}}(x, S)$ . The set of points in  $S$ , denoted by  $\mathbf{X} = \{X_1, \dots, X_{|S|}\} \subseteq \mathcal{X}$ , induces a *Voronoi partition* of  $\mathcal{X}$ ,  $\mathcal{V}(\mathbf{X}) := \{V_1(\mathbf{X}), \dots, V_{|S|}(\mathbf{X})\}$ , where each Voronoi cell is

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

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, **OptiNet**, 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 5.10 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  $\sigma$ -algebra  $\mathcal{B}$ . Specializing to **OptiNet**, we have  $\text{Alg}(S_n) = h_{S'_n}$ .

Symbol	Brief description
$(\mathcal{X}, \rho, \mu)$	metric probability space
$\mathcal{B}$	Borel $\sigma$ -algebra induced by $\rho$
$B_r(x)$	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
$S_n = (\mathbf{X}_n, \mathbf{Y}_n)$	random sample of size $n$
$S_n(\mathbf{i}, \mathbf{j})$	subsample $(\mathbf{X}_i, \mathbf{Y}_j)$ 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 $(\mathbf{X}_i, \mathbf{Y}^*)$ with true majority vote labels
$\mathbf{X}(\gamma)$	$\gamma$ -net of $\mathbf{X}_n$
$S_n(\gamma)$	subsample $(\mathbf{X}(\gamma), \mathbf{Y}(\gamma))$ with empirical majority votes
$M_n(\gamma) =  \mathbf{X}(\gamma) $	size of the $\gamma$ -net
$h_S$	1-NN classifier induced by the labeled set $S$
$\alpha_n(\gamma)$	empirical error of $h_{S_n(\gamma)}$ on $S_n$
$\text{UB}_\gamma(A)$	$\gamma$ -envelope of $A \subseteq \mathcal{X}$
$L_\gamma(A)$	$\gamma$ -missing mass of $A \subseteq \mathcal{X}$
$\mathcal{V}(\mathbf{X})$	Voronoi partition of $\mathcal{X}$ induced by $\mathbf{X}$

TABLE 1  
Symbols guide

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

For  $n \in \mathbb{N}$ , denote  $[n] := \{1, \dots, n\}$ . 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 sub-sample of  $S_n$  indexed by  $\mathbf{i}$  by  $S_n(\mathbf{i}) := \{(X_{i_1}, Y_{i_1}), \dots, (X_{i_d}, Y_{i_d})\}$ . Similarly, for a vector  $\mathbf{Y}' = \{Y'_1, \dots, Y'_d\} \in \mathcal{Y}^d$ , denote by  $S_n(\mathbf{i}, \mathbf{Y}') := \{(X_{i_1}, Y'_1), \dots, (X_{i_d}, Y'_d)\}$ , namely the sub-sample 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  $\limsup_{n \rightarrow \infty} f(n)/g(n) = 0$ . In accordance with common convention, we often use the less precise notation  $f(n) = O(g(n))$ , etc.

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

**3. OptiNet: 1-NN majority-based compression.** In this section we describe the 1-NN majority-based compression algorithm proposed in [35], which was provisionally titled KSU in [36], and which we name here OptiNet.

The algorithm (see Alg. 1) operates as follows. Given an input sample  $S_n$ , whose set of points is denoted by  $\mathbf{X}_n = \{X_1, \dots, X_n\}$ , OptiNet considers all possible scales  $\gamma > 0$  consisting of the interpoint distances in  $\mathbf{X}_n$ , and the additional scale  $\gamma = \infty$ .<sup>4</sup> For each such scale it constructs a  $\gamma$ -net of  $\mathbf{X}_n$  (note that any singleton in  $\mathbf{X}_n$  is an

<sup>4</sup>The analysis is streamlined by including  $\gamma = \infty$  rather than some  $\gamma > \text{diam}(\mathbf{X}_n)$ .

---

**Algorithm 1** OptiNet: 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  $\gamma$ -net of  $\{X_1, \dots, X_n\}$   
 4:   let  $M_n(\gamma) := |\mathbf{X}(\gamma)|$   
 5:   for each  $i \in [M_n(\gamma)]$ , let  $Y'_i$  be the majority label in  $V_i(\mathbf{X}(\gamma))$  as defined 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  $\gamma_n^* \in \arg\min_{\gamma \in \Gamma} Q(n, \alpha_n(\gamma), 2M_n(\gamma), \delta)$ , where  $Q$  is defined in (3.2)  
 10: set  $S'_n := S'_n(\gamma_n^*)$   
 11: **return**  $h_{S'_n}$

---

$\infty$ -net). Denote this  $\gamma$ -net by  $\mathbf{X}(\gamma) := \{X_{i_1}, \dots, X_{i_M}\}$ , where  $M \equiv M_n(\gamma)$  denotes its size and  $\mathbf{i} \equiv \mathbf{i}(\gamma) := \{i_1, \dots, i_M\} \in [n]^M$  denotes the indices selected from  $S_n$  for this  $\gamma$ -net. For every such  $\gamma$ -net, the algorithm chooses the labels  $\mathbf{Y}'(\gamma) \in \mathcal{Y}^M$ , which are the empirical majority-vote labels in the respective Voronoi cells in the partition  $\mathcal{V}(\mathbf{X}(\gamma)) = \{V_1, \dots, V_M\}$ . Formally, for  $i \in [M]$ ,

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

where ties are broken lexicographically (the “smallest”  $y \in \mathcal{Y}$  is chosen, according to some fixed ordering). This procedure creates a labeled set  $S'_n(\gamma) := S_n(\mathbf{i}(\gamma), \mathbf{Y}'(\gamma))$  for every candidate  $\gamma \in \{\rho(X_i, X_j) : i, j \in [n]\} \cup \{\infty\} \setminus \{0\}$ . The algorithm then selects a single  $\gamma$ , denoted  $\gamma^* \equiv \gamma_n^*$ , and outputs  $h_{S'_n(\gamma^*)}$ . The scale  $\gamma^*$  is selected so as to minimize a generalization error bound, which upper bounds  $\text{err}(S'_n(\gamma))$  with high probability. This error bound, denoted by  $Q$  in the algorithm, can be derived using a compression-based analysis, as described below.

We say that a specific  $S'_n$  is an  $(\alpha, 2m)$ -compression of  $S_n$  if there exist  $\mathbf{i}, \mathbf{j} \in [n]^m$  such that  $S'_n = S_n(\mathbf{i}, \mathbf{j})$  and  $\widehat{\text{err}}(h_{S'_n}, S_n) \leq \alpha$ . We say that a mapping from samples to classifiers  $S'_n \mapsto \Phi[S'_n]$  is permutation-invariant if  $\Phi[S'_n]$  is invariant under permutations of  $S'_n$ . Clearly, the subsample  $S'_n(\gamma)$  computed by Alg. 1 is an  $(\alpha_n(\gamma), 2M_n(\gamma))$ -compression of  $S_n$  for any  $\gamma > 0$ , and the 1-NN mapping  $S'_n \mapsto h_{S'_n}$  is permutation-invariant. The following generalization bound is a straightforward adaptation of the compression bound in [24, Theorem 8] to the permutation-invariant case; our proof requires the sharpened variant.

**PROPOSITION 3.1.** *Let  $S_n \sim \bar{\mu}^n$  and  $\delta > 0$ . With probability  $1 - \delta$ , for any  $\alpha \in [0, 1]$  and any  $m \geq 1$ , if  $S'_n$  is a  $(\alpha, m)$ -compression of  $S_n$ , and the mapping  $S'_n \mapsto \Phi[S'_n]$  is permutation-invariant, then*

$$\begin{aligned}
 \text{err}(\Phi[S'_n]) &\leq \frac{\alpha n}{n - m} + O\left(\frac{m \log \frac{n}{m} + \log \frac{1}{\delta}}{n - m} + \sqrt{\frac{\frac{\alpha n m}{n - m} \log \frac{n}{m} + \log \frac{1}{\delta}}{n - m}}\right) \\
 &=: Q(n, \alpha, m, \delta).
 \end{aligned} \quad (3.2)$$

As presented in Alg. 1, **OptiNet** has a naive runtime complexity of  $O(n^4)$ , since  $O(n^2)$  values of  $\gamma$  are considered and a  $\gamma$ -net is constructed for each one in time  $O(n^2)$  (see [25, Algorithm 1]). Improved runtimes can be obtained in doubling spaces, e.g., using the methods in [41, 23]. In this work we focus on the Bayes consistency of **OptiNet**, rather than optimizing its computational complexity. Our Bayes consistency results below hold for **OptiNet**, whenever the generalization bound  $Q(n, \alpha, m, \delta)$  satisfies the following properties:

- Q1.** For any  $n \in \mathbb{N}$  and  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$  over the  $S_n \sim \bar{\mu}^n$ , for all  $\alpha \in [0, 1]$  and  $m \in [n]$ : If  $S'_n$  is an  $(\alpha, 2m)$ -compression of  $S_n$ , then  $\text{err}(h_{S'_n}) \leq Q(n, \alpha, 2m, \delta)$ .
- Q2.**  $Q$  is monotonically increasing in  $\alpha$  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 any sequence  $m_n \in o(n)$ ,

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

Properties **Q1** and **Q2** are equivalent to properties assumed in [35]. The compression bound in (3.2) clearly satisfies these properties. In contrast, property **Q3** is slightly stronger than the corresponding property in [35], since it allows  $m$  to grow as  $o(n)$ . The latter property holds due to the tighter compression bound in (3.2). Indeed, **Q3** is satisfied by (3.2) via any convergent series  $\sum_{n=1}^\infty \delta_n < \infty$  such that  $\delta_n = e^{-o(n)}$ ; note that the decay of  $\delta_n$  cannot be too rapid.

**REMARK 3.2.** We presented the algorithm **OptiNet** with  $\gamma$  selected based on a compression bound. This choice of presentation creates a close connection between the algorithm and the proof of consistency below. However, it is worth noting that we could instead have chosen  $\gamma$  based on a hold-out validation set: for instance, using  $n/2$  of the  $n$  samples to construct the predictor for each possible  $\gamma$  value, and then from among these values  $\gamma$ , we can select the  $\gamma$  whose predictor makes the smallest number of mistakes on the remaining  $n/2$  samples. Since the analysis below shows that there exists a choice of  $\gamma^*$  for each  $n$  such that **OptiNet** is Bayes consistent, this alternative technique of selecting  $\gamma$  based on a hold-out sample would only lose an additive  $O\left(\sqrt{\log(n)/n}\right)$  compared to using that  $\gamma^*$ , and hence would also be Bayes consistent.  $\blacktriangleleft$

**4. OptiNet is universally Bayes consistent in all essentially separable metric spaces.** In this section, we prove that **OptiNet** is universally strongly Bayes consistent in all essentially separable metric spaces. Recall that  $(\mathcal{X}, \rho)$  is *separable* if it contains a dense countable set. A metric probability space  $(\mathcal{X}, \rho, \mu)$  is separable if there is an  $\mathcal{X}' \subseteq \mathcal{X}$  with  $\mu(\mathcal{X}') = 1$  such that  $(\mathcal{X}', \rho)$  is separable. We will call a metric space  $(\mathcal{X}, \rho)$  *essentially separable* (ES) if for every probability measure  $\mu$  on  $\mathcal{B}$ , the metric probability space  $(\mathcal{X}, \rho, \mu)$  is separable.

A Bayes consistency result was obtained in [36] for spaces with finite doubling dimension and diameter. Our main technical innovation, which allowed us to circumvent both finiteness assumptions of [36], is the sublinear growth of  $\gamma$ -nets established in Lemma 4.1:

LEMMA 4.1. *Let  $(\mathcal{X}, \rho, \mu)$  be a separable metric probability space. For  $S_n \sim \mu^n$ , let  $\mathbf{X}(\gamma)$  be any  $\gamma$ -net of  $S_n$ , and  $M_n(\gamma) := |\mathbf{X}(\gamma)|$ . Then, for any  $\gamma > 0$ , there exists a function  $t_\gamma(n) : \mathbb{N} \rightarrow \mathbb{R}_+$  such that  $t_\gamma(n) \in o(n)$  and*

$$\mathbb{P} \left[ M_n(\gamma) \geq t_\gamma(n) \right] \leq 1/n^2. \quad (4.1)$$

The proof of Lemma 4.1 is provided in Appendix A. Another straightforward but crucial insight is to approximate  $L_1(\mu)$  functions by Lipschitz ones (Lemma A.1) rather than by continuous functions with compact support as in [36]. The latter approximation requires local compactness, which essentially amounts to a finite dimensionality condition.

THEOREM 4.2. *Let  $(\mathcal{X}, \rho, \mu)$  be a separable metric probability space. Let  $Q$  be a generalization bound that satisfies Properties **Q1**, **Q2**, **Q3**, and let  $\delta_n$  be as stipulated by **Q3**. If the input confidence  $\delta$  for input size  $n$  is set to  $\delta_n$ , then the 1-NN classifier  $h_{S'_n(\gamma_n^*)}$  calculated by OptiNet is strongly Bayes consistent on  $(\mathcal{X}, \rho, \mu)$ :*

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

Given a sample  $S_n \sim \mu^n$ , 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 shown in Section 3, the labeled set  $S'_n(\gamma_n^*)$  computed by Alg. 1 is a permutation-invariant  $(\alpha_n^*, 2M_n^*)$ -compression of the sample  $S_n$ . For brevity we denote

$$Q_n(\alpha, m) := Q(n, \alpha, 2m, \delta_n).$$

To prove Theorem 4.2 we decompose the excess error over the Bayes into two terms:

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

and show that each term decays to zero almost surely.

For the first term,  $T_{\text{I}}(n)$ , from Property **Q1** of generalization bound  $Q$ , we have that for any  $n > 0$ ,

$$\mathbb{P} \left[ \text{err}(h_{S'_n(\gamma_n^*)}) - Q_n(\alpha_n^*, M_n^*) > 0 \right] \leq \delta_n. \quad (4.2)$$

Since  $\sum \delta_n < \infty$ , the Borel-Cantelli lemma implies  $\limsup_{n \rightarrow \infty} T_{\text{I}}(n) \leq 0$  with probability 1.

The main challenge is to prove that  $\limsup_{n \rightarrow \infty} T_{\text{II}}(n) \leq 0$  almost surely. We begin by showing that the Bayes error  $R^*$  can be approached using classifiers defined by the true majority-vote labeling over fine partitions of  $\mathcal{X}$ . Formally, 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

$$y^*(E) = \underset{y \in \mathcal{Y}}{\text{argmax}} \mathbb{P}(Y = y | X \in E), \quad (4.3)$$

where ties are broken lexicographically. To ensure that  $y^*$  is always well-defined, when  $E = \emptyset$  we arbitrarily define it to be the lexicographically first  $y \in \mathcal{Y}$ . Given  $\mathcal{V}$  and a measurable set  $W \subseteq \mathcal{X}$ , consider the true majority-vote classifier  $h_{\mathcal{V},W}^* : \mathcal{X} \rightarrow \mathcal{Y}$  given by

$$h_{\mathcal{V},W}^*(x) = y^*(I_{\mathcal{V}}(x) \cap W). \quad (4.4)$$

Note that if  $x \notin W$ , this classifier assigns a label to  $x$  based on the true majority-vote in a set that does not contain  $x$ . To bound the error of  $h_{\mathcal{V},W}^*$  for any conditional distribution of labels, we use the fact that for any metric probability space, Lipschitz functions are dense in  $L_1(\mu)$  (Lemma A.1, see Appendix).

We have the following uniform approximation bound for the error of classifiers in the form of (4.4), essentially extending the approximation analysis done in the proof of [16, Theorem 21.2], which holds for the special case  $|\mathcal{Y}| = 2$  and  $\mathcal{X} = \mathbb{R}^d$ , to the more general multiclass problem in metric probability spaces.<sup>5</sup>

LEMMA 4.3. *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 (4.4) satisfies*

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

PROOF. Let  $\eta_y : \mathcal{X} \rightarrow [0, 1]$  be the conditional probability function for label  $y \in \mathcal{Y}$ ,

$$\eta_y(x) = \mathbb{P}(Y = y | X = x),$$

which is measurable by [53, Corollary B.22]. Define  $\tilde{\eta}_y : \mathcal{X} \rightarrow [0, 1]$  as  $\eta_y$ 's conditional expectation function with respect to  $(\mathcal{V}, W)$ : For  $x$  such that  $I_{\mathcal{V}}(x) \cap W \neq \emptyset$ ,

$$\tilde{\eta}_y(x) = \mathbb{P}(Y = y | X \in I_{\mathcal{V}}(x) \cap W) = \frac{\int_{I_{\mathcal{V}}(x) \cap W} \eta_y(z) d\mu(z)}{\mu(I_{\mathcal{V}}(x) \cap W)}.$$

For other  $x$ , define  $\tilde{\eta}_y(x) = \mathbf{1}[y \text{ is lexicographically first}]$ . Note that  $(\tilde{\eta}_y)_{y \in \mathcal{Y}}$  are piecewise constant on the cells of the restricted partition  $\mathcal{V} \cap W$ . By definition, the Bayes classifier  $h^*$  and the true majority-vote classifier  $h_{\mathcal{V},W}^*$  satisfy

$$\begin{aligned} h^*(x) &= \operatorname{argmax}_{y \in \mathcal{Y}} \eta_y(x), \\ h_{\mathcal{V},W}^*(x) &= \operatorname{argmax}_{y \in \mathcal{Y}} \tilde{\eta}_y(x). \end{aligned}$$

---

<sup>5</sup>An erroneous claim that would have contradicted Lemma 4.3 was made in [36, Thm. 5]. See [37] for a corrected version of [36].

It follows that

$$\begin{aligned}
 & \mathbb{P}(h_{\mathcal{Y},W}^*(X) \neq Y \mid X = x) - \mathbb{P}(h^*(X) \neq Y \mid X = x) \\
 &= \eta_{h^*(x)}(x) - \eta_{h_{\mathcal{Y},W}^*(x)}(x) \\
 &= \max_{y \in \mathcal{Y}} \eta_y(x) - \max_{y \in \mathcal{Y}} \tilde{\eta}_y(x) \\
 &\leq \max_{y \in \mathcal{Y}} |\eta_y(x) - \tilde{\eta}_y(x)|.
 \end{aligned}$$

By condition (i) in the lemma statement,  $\mu(\mathcal{X} \setminus W) \leq \nu$ . Thus,

$$\begin{aligned}
 \text{err}(h_{\mathcal{Y},W}^*) - R^* &= \mathbb{P}(h_{\mathcal{Y},W}^*(X) \neq Y) - \mathbb{P}(h^*(X) \neq Y) \\
 &\leq \mu(\mathcal{X} \setminus W) + \int_W \max_{y \in \mathcal{Y}} |\eta_y(x) - \tilde{\eta}_y(x)| \, d\mu(x) \\
 &\leq \nu + \sum_{y \in \mathcal{Y}} \int_W |\eta_y(x) - \tilde{\eta}_y(x)| \, d\mu(x).
 \end{aligned}$$

Let  $\mathcal{Y}_\nu \subseteq \mathcal{Y}$  be a finite set of labels such that  $\mathbb{P}[Y \in \mathcal{Y}_\nu] \geq 1 - \nu$ . Then

$$\text{err}(h_{\mathcal{Y},W}^*) - R^* \leq 2\nu + \sum_{y \in \mathcal{Y}_\nu} \int_W |\eta_y(x) - \tilde{\eta}_y(x)| \, d\mu(x). \quad (4.5)$$

To bound the integrals in (4.5), we approximate  $(\eta_y)_{y \in \mathcal{Y}}$  with functions from the dense set of Lipschitz functions, applying Lemma A.1. Since  $\eta_y \in L_1(\mu)$  for all  $y \in \mathcal{Y}_\nu$  and  $|\mathcal{Y}_\nu| < \infty$ , Lemma A.1 implies that there are  $|\mathcal{Y}_\nu|$  Lipschitz functions  $(r_y)_{y \in \mathcal{Y}_\nu}$  such that

$$\max_{y \in \mathcal{Y}_\nu} \int_{\mathcal{X}} |\eta_y(x) - r_y(x)| \, d\mu(x) \leq \nu/|\mathcal{Y}_\nu|. \quad (4.6)$$

Similarly to  $(\tilde{\eta}_y)_{y \in \mathcal{Y}_\nu}$ , define the piecewise constant functions  $(\tilde{r}_y)_{y \in \mathcal{Y}_\nu}$  by

$$\tilde{r}_y(x) = \mathbb{E}[r_y(X) \mid X \in I_{\mathcal{Y}}(x) \cap W] = \frac{\int_{I_{\mathcal{Y}}(x) \cap W} r_y(z) \, d\mu(z)}{\mu(I_{\mathcal{Y}}(x) \cap W)}.$$

We bound each integrand in (4.5) by

$$\begin{aligned}
 & |\eta_y(x) - \tilde{\eta}_y(x)| \\
 &\leq |\eta_y(x) - r_y(x)| + |r_y(x) - \tilde{r}_y(x)| + |\tilde{r}_y(x) - \tilde{\eta}_y(x)|.
 \end{aligned} \quad (4.7)$$

The integral of the first term in (4.7) is smaller than  $\nu/|\mathcal{Y}_\nu|$  by the definition of  $r_y$

in (4.6). For the integral of the third term in (4.7),

$$\begin{aligned}
& \int_W |\tilde{r}_y(x) - \tilde{\eta}_y(x)| \, d\mu(x) \\
&= \sum_{V \in \mathcal{V}} |\mathbb{E}[r_y(X)\mathbf{1}[X \in V \cap W]] - \mathbb{E}[\eta_y(X)\mathbf{1}[X \in V \cap W]]| \\
&= \sum_{V \in \mathcal{V}} \left| \int_{V \cap W} r_y(x) \, d\mu(x) - \int_{V \cap W} \eta_y(x) \, d\mu(x) \right| \\
&= \sum_{V \in \mathcal{V}} \left| \int_{V \cap W} (r_y(x) - \eta_y(x)) \, d\mu(x) \right| \\
&\leq \int_W |r_y(x) - \eta_y(x)| \, d\mu(x) \leq \nu/|\mathcal{Y}_\nu|.
\end{aligned}$$

Finally, for the integral of the second term in (4.7), we denote

$$\bar{\mathcal{V}} = \{V \cap W : \mu(V \cap W) \neq 0, V \in \mathcal{V}\}$$

and note that

$$\begin{aligned}
& \int_W |r_y(x) - \tilde{r}_y(x)| \, d\mu(x) \\
&= \sum_{\bar{V} \in \bar{\mathcal{V}}} \int_{\bar{V}} \left| r_y(x) - \frac{\mathbb{E}[r_y(X)\mathbf{1}[X \in \bar{V}]]}{\mu(\bar{V})} \right| \, d\mu(x) \\
&= \sum_{\bar{V} \in \bar{\mathcal{V}}} \frac{1}{\mu(\bar{V})} \int_{\bar{V}} |r_y(x)\mu(\bar{V}) - \mathbb{E}[r_y(X)\mathbf{1}[X \in \bar{V}]]| \, d\mu(x) \\
&= \sum_{\bar{V} \in \bar{\mathcal{V}}} \frac{1}{\mu(\bar{V})} \int_{\bar{V}} \left| r_y(x) \int_{\bar{V}} d\mu(z) - \int_{\bar{V}} r_y(z) \, d\mu(z) \right| \, d\mu(x) \\
&= \sum_{\bar{V} \in \bar{\mathcal{V}}} \frac{1}{\mu(\bar{V})} \int_{\bar{V}} \left| \int_{\bar{V}} (r_y(x) - r_y(z)) \, d\mu(z) \right| \, d\mu(x) \\
&\leq \sum_{\bar{V} \in \bar{\mathcal{V}}} \frac{1}{\mu(\bar{V})} \int_{\bar{V}} \int_{\bar{V}} |r_y(x) - r_y(z)| \, d\mu(x) \, d\mu(z).
\end{aligned}$$

Since  $|\mathcal{Y}_\nu| < \infty$  and any Lipschitz function is uniformly continuous on all of  $\mathcal{X}$ , the finite collection  $\{r_y : y \in \mathcal{Y}_\nu\}$  is equicontinuous. Namely, there exists a diameter  $\beta = \beta(\nu) > 0$  such that for any  $A \subseteq \mathcal{X}$  with  $\text{diam}(A) \leq \beta$ ,

$$\max_{y \in \mathcal{Y}_\nu} |r_y(x) - r_y(z)| \leq \nu/|\mathcal{Y}_\nu|$$

for every  $x, z \in A$  (note that  $\beta(\nu)$  does not depend on  $(\mathcal{V}, W)$ ). By condition (ii) in the lemma statement,  $\text{diam}(V \cap W) \leq \beta$  for all  $V \in \mathcal{V}$ . Hence,

$$\frac{1}{\mu(V \cap W)} \int_{V \cap W} \int_{V \cap W} |r_y(x) - r_y(z)| \, d\mu(x) \, d\mu(z) \leq \frac{\nu}{|\mathcal{Y}_\nu|} \mu(V \cap W).$$

Summing over all cells  $V \in \mathcal{V}$  with  $\mu(V \cap W) \neq 0$ , the integral of the second term in (4.7) satisfies

$$\int_W |r_y(x) - \tilde{r}_y(x)| d\mu(x) \leq \nu/|\mathcal{Y}_\nu|.$$

Combining the bounds for the three terms,

$$\sum_{y \in \mathcal{Y}_\nu} \int_W |\eta_y(x) - \tilde{\eta}_y(x)| d\mu(x) \leq \sum_{y \in \mathcal{Y}_\nu} \frac{3\nu}{|\mathcal{Y}_\nu|} = 3\nu.$$

Applying this bound to (4.5), we conclude  $\text{err}(h_{\mathcal{V},W}^*) - R^* \leq 5\nu$ .  $\square$

Next, we are going to invoke Lemma 4.3 to show that the generalization bound  $Q_n(\alpha_n^*, M_n^*)$  also approaches the Bayes error  $R^*$ , thus proving  $\limsup_{n \rightarrow \infty} T_{\Pi}(n) \leq 0$  almost surely. Given  $\varepsilon > 0$ , fix

$$\gamma = \gamma(\varepsilon) = \beta(\varepsilon/10)/4, \quad (4.8)$$

where  $\beta(\cdot)$  is furnished by Lemma 4.3.

Define  $\gamma_{\text{rand}} = \min\{\tilde{\gamma} \in \Gamma : \tilde{\gamma} \geq \gamma\}$  and observe that  $\mathbf{X}(\gamma_{\text{rand}})$ , the  $\gamma_{\text{rand}}$ -net calculated by the algorithm, is also a  $\gamma$ -net. Indeed, this clearly holds if  $\gamma_{\text{rand}} = \gamma$ . Otherwise,  $\gamma_{\text{rand}} > \gamma$ , and since  $\Gamma$  includes all interpoint distances in  $\mathbf{X}_n$ , there are no interpoint distances in  $[\gamma, \gamma_{\text{rand}})$ . In this case, any two points in  $\mathbf{X}(\gamma_{\text{rand}})$  are at least  $\gamma_{\text{rand}} > \gamma$  apart, as required by the definition of a  $\gamma$ -net. In addition, any point  $X_i \in \mathbf{X}_n$  has  $\rho(X_i, \mathbf{X}(\gamma_{\text{rand}})) < \gamma_{\text{rand}}$ . Since  $\rho(X_i, \mathbf{X}(\gamma_{\text{rand}})) \notin [\gamma, \gamma_{\text{rand}})$ , it follows that  $\rho(X_i, \mathbf{X}(\gamma_{\text{rand}})) < \gamma$ . Thus,  $\mathbf{X}(\gamma_{\text{rand}})$  is a maximal subset of  $\mathbf{X}_n$  in which any two points are at least  $\gamma$  apart, hence it is a  $\gamma$ -net. In accordance, from now on we use the notation  $\mathbf{X}(\gamma) \equiv \mathbf{X}(\gamma_{\text{rand}})$ , and similarly for  $\alpha_n(\gamma)$ ,  $M_n(\gamma)$ ,  $\mathbf{i}(\gamma)$ ,  $\mathbf{Y}'(\gamma)$ .

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

$$\mathbb{P}[Q_n(\alpha_n(\gamma), M_n(\gamma)) > R^* + \varepsilon] \leq Cne^{-cn\varepsilon^2} + 1/n^2. \quad (4.9)$$

Since Alg. 1 finds  $\gamma_n^*$  such that

$$\begin{aligned} Q_n(\alpha_n^*, M_n^*) &= \min_{\gamma' \in \Gamma} Q_n(\alpha_n(\gamma'), M_n(\gamma')) \\ &\leq Q_n(\alpha_n(\gamma_r), M_n(\gamma_r)) = Q_n(\alpha_n(\gamma), M_n(\gamma)), \end{aligned}$$

the bound in (4.9) implies that  $\forall n \geq N$ ,

$$\mathbb{P}[Q_n(\alpha_n^*, M_n^*) > R^* + \varepsilon] \leq Cne^{-cn\varepsilon^2} + 1/n^2. \quad (4.10)$$

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

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

Since  $\forall n, T_{\text{I}}(n) + T_{\text{II}}(n) \geq 0$ , this implies  $\lim_{n \rightarrow \infty} T_{\text{II}}(n) = 0$  almost surely, thus completing the proof of Theorem 4.2.

It remains to prove (4.9). For  $A \subseteq \mathcal{X}$ , we denote its  $\gamma$ -envelope by  $\text{UB}_\gamma(A) := \cup_{x \in A} B_\gamma(x)$  and consider the  $\gamma$ -missing mass of  $S_n$ , defined as the following random variable:

$$L_\gamma(S_n) := \mu(\mathcal{X} \setminus \text{UB}_\gamma(S_n)). \quad (4.11)$$

We bound the left-hand side of (4.9) by

$$\begin{aligned} & \mathbb{P}[Q_n(\alpha_n(\gamma), M_n(\gamma)) > R^* + \varepsilon] \quad (4.12) \\ & \leq \mathbb{P}\left[Q_n(\alpha_n(\gamma), M_n(\gamma)) > R^* + \varepsilon \wedge L_\gamma(S_n) \leq \frac{\varepsilon}{10} \wedge M_n(\gamma) \leq t_\gamma(n)\right] \\ & \quad + \mathbb{P}[L_\gamma(S_n) > \varepsilon/10] + \mathbb{P}[M_n(\gamma) > t_\gamma(n)] \\ & =: P_{\text{I}} + P_{\text{II}} + P_{\text{III}}, \end{aligned}$$

where  $t_\gamma(n) \in o(n)$  is a function specified by Lemma 4.1. This lemma implies that  $P_{\text{III}} \leq 1/n^2$ , while a bound on  $P_{\text{II}}$  is furnished by the following lemma, whose proof is given in Appendix A:

LEMMA 4.4. *Let  $(\mathcal{X}, \rho, \mu)$  be a separable metric probability space,  $\gamma > 0$  be fixed, and the  $\gamma$ -missing mass  $L_\gamma$  defined as in (4.11). Then, there exists a function  $u_\gamma(n) \in o(1)$ , such that for  $S_n \sim \mu^n$  and all  $t > 0$ ,*

$$\mathbb{P}[L_\gamma(S_n) \geq u_\gamma(n) + t] \leq \exp(-nt^2). \quad (4.13)$$

REMARK 4.5. A precursor to the present work, [36], assumed finite doubling dimension and diameter for  $\mathcal{X}$ , and bounded the  $\gamma$ -missing mass via covering numbers, as in [4, Theorem 8]. ◀

Taking  $n$  sufficiently large so that  $u_\gamma(n)$  — furnished by Lemma 4.4 — satisfies  $u_\gamma(n) \leq \varepsilon/20$ , and invoking Lemma 4.4 with  $t = \varepsilon/20$ , we have

$$P_{\text{II}} = \mathbb{P}[L_\gamma(S_n) > \varepsilon/10] \leq e^{-\frac{n\varepsilon^2}{400}}. \quad (4.14)$$

It remains to bound  $P_{\text{I}}$ . 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 4.6. *Fix  $\varepsilon > 0$  and let  $\gamma = \gamma(\varepsilon) = \beta(\varepsilon/10)/4$  as in (4.8). Under the conditions of Theorem 4.2, there exists an  $n_0$  such that for all  $n \geq n_0$ , and for all  $d \in [t_\gamma(n)]$ ,*

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

Applying Lemma 4.6 and summing over all  $1 \leq d \leq t_\gamma(n)$ , we have that, for  $n$  large enough so that  $t_\gamma(n) \leq n$ ,

$$P_1 \leq \sum_{d=1}^{t_\gamma(n)} p_d \leq t_\gamma(n) e^{-\frac{n\epsilon^2}{32}} \leq n e^{-\frac{n\epsilon^2}{32}}. \quad (4.15)$$

Plugging (4.15), (4.14), and (4.1) into (4.12), we get that (4.9) holds, which completes the proof of Theorem 4.2. The proof of Lemma 4.6 follows.

**PROOF OF LEMMA 4.6.** Let  $\mathbf{i} = \mathbf{i}(\gamma) \in [n]^d$  be the set of indices in the net  $\mathbf{X} = \mathbf{X}(\gamma)$  selected by the algorithm. Let  $\mathbf{Y}^* \in \mathcal{Y}^d$  be the true majority-vote labels with respect to the restricted partition  $\mathcal{V}(\mathbf{X}) \cap \text{UB}_{2\gamma}(\mathbf{X})$ ,

$$(\mathbf{Y}^*)_j = y^*(V_j \cap \text{UB}_{2\gamma}(\mathbf{X})), \quad j \in [d]. \quad (4.16)$$

We pair  $\mathbf{X}$  with the labels  $\mathbf{Y}^*$  to obtain the labeled set

$$S_n(\mathbf{i}, *) := S_n(\mathbf{i}, \mathbf{Y}^*) = (\mathbf{X}, \mathbf{Y}^*) \in (\mathcal{X} \times \mathcal{Y})^d. \quad (4.17)$$

Note that conditioned on  $\mathbf{X}$ ,  $S_n(\mathbf{i}, *)$  does not depend on the rest of  $S_n$ .

The induced 1-NN classifier  $h_{S_n(\mathbf{i}, *)}(x)$  can be expressed as  $h_{\mathcal{V}, W}^*(x) = y^*(I_{\mathcal{V}}(x) \cap W)$  with  $\mathcal{V} = \mathcal{V}(\mathbf{X})$  and  $W = \text{UB}_{2\gamma}(\mathbf{X})$  (see (4.4) for the definition of  $h_{\mathcal{V}, W}^*$ ). We now show that

$$L_\gamma(\mathbf{X}_n) \leq \frac{\epsilon}{10} \implies \text{err}(h_{S_n(\mathbf{i}, *)}) \leq R^* + \epsilon/2, \quad (4.18)$$

by showing that under the assumption  $L_\gamma(\mathbf{X}_n) \leq \frac{\epsilon}{10}$ , the conditions of Lemma 4.3 hold for  $\mathcal{V}, W$  as defined above. To this end, we bound the diameter of the partition  $\mathcal{V} \cap W = \mathcal{V} \cap \text{UB}_{2\gamma}(\mathbf{X})$ , and the measure of the missing mass  $\mu(\mathcal{X} \setminus W) = L_{2\gamma}(\mathbf{X})$  under the assumption.

To bound the diameter of the partition  $\mathcal{V} \cap \text{UB}_{2\gamma}(\mathbf{X})$ , let  $x \in V_j \cap \text{UB}_{2\gamma}(\mathbf{X})$ . Note that  $V_j$  is the Voronoi cell centered at  $x_{i_j} \in \mathbf{X}$ . Then  $\rho(x, x_{i_j}) = \min_{i \in \mathbf{i}} \rho(x, x_i)$  and, since  $x \in \text{UB}_{2\gamma}(\mathbf{X})$ ,  $\min_{i \in \mathbf{i}} \rho(x, x_i) \leq 2\gamma$ . Therefore

$$\text{diam}(\mathcal{V} \cap W) = \max_j \text{diam}(V_j \cap \text{UB}_{2\gamma}(\mathbf{X})) \leq 4\gamma.$$

To bound  $L_{2\gamma}(\mathbf{X})$  under the assumption  $L_\gamma(\mathbf{X}_n) \leq \frac{\epsilon}{10}$ , observe that for all  $z \in \text{UB}_\gamma(\mathbf{X}_n)$ , there is some  $i \in [n]$  such that  $z \in B_\gamma(x_i)$ . For this  $i$ , there is some  $j \in \mathbf{i}$  such that  $x_i \in B_\gamma(x_j)$ , since  $\mathbf{X}$  is a  $\gamma$ -net of  $\mathbf{X}_n$ . Therefore  $z \in B_{2\gamma}(x_j)$ . Thus,  $z \in \text{UB}_{2\gamma}(\mathbf{X})$ . It follows that  $\text{UB}_\gamma(\mathbf{X}_n) \subseteq \text{UB}_{2\gamma}(\mathbf{X})$ , thus  $L_{2\gamma}(\mathbf{X}) \leq L_\gamma(\mathbf{X}_n)$ . Under the assumption, we thus have  $L_{2\gamma}(\mathbf{X}) \leq \frac{\epsilon}{10}$ . Hence, by the choice of  $\gamma = \gamma(\epsilon)$  in the statement of the lemma, Lemma 4.3 implies (4.18).

To bound  $Q_n(\alpha_n(\gamma), M_n(\gamma))$ , we consider the relationship between the hypothetical true majority-vote classifier  $h_{S_n(\mathbf{i}, *)}$  and the actual classifier returned by the algorithm,  $h_{S_n(\mathbf{i}, \mathbf{Y}^\gamma)}$ . Note that

$$\alpha_n(\gamma) = \widehat{\text{err}}(h_{S_n(\mathbf{i}, \mathbf{Y}^\gamma)}, S_n) = \min_{\mathbf{Y} \in \mathcal{Y}^d} \widehat{\text{err}}(h_{S_n(\mathbf{i}, \mathbf{Y})}, S_n) \leq \widehat{\text{err}}(h_{S_n(\mathbf{i}, *)}, S_n),$$

and thus, from the monotonicity Property **Q2** of  $Q$ ,

$$Q_n(\alpha_n(\gamma), M_n(\gamma)) \leq Q_n(\widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n), M_n(\gamma)). \quad (4.19)$$

Combining (4.18) and (4.19) we have that

$$\begin{aligned} & \left\{ Q_n(\alpha_n(\gamma), M_n(\gamma)) > R^* + \varepsilon \wedge L_\gamma(\mathbf{X}_n) \leq \frac{\varepsilon}{10} \wedge M_n(\gamma) = d \right\} \\ & \implies \left\{ Q_n(\widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n), d) > \text{err}(h_{S_n(\mathbf{i},*)}) + \frac{\varepsilon}{2} \wedge |\mathbf{i}| = d \right\}. \end{aligned}$$

Hence, for all  $d \leq t_\gamma(n)$ ,

$$\begin{aligned} p_d & \leq \mathbb{P} \left[ Q_n(\widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n), d) > \text{err}(h_{S_n(\mathbf{i},*)}) + \frac{\varepsilon}{2} \wedge |\mathbf{i}| = d \right] \\ & \leq \mathbb{P} \left[ \exists \mathbf{i} \in [n]^d : Q_n(\widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n), d) > \text{err}(h_{S_n(\mathbf{i},*)}) + \frac{\varepsilon}{2} \right]. \quad (4.20) \end{aligned}$$

To bound the last expression, let  $\mathbf{i} \in [n]^d$  and denote

$$r_{d,n} = \sup_{\alpha \in (0,1)} (Q_n(\alpha, d) - \alpha).$$

We thus have,

$$Q_n(\widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n), d) \leq \widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n) + r_{d,n}.$$

Let  $\mathbf{i}' = \{1, \dots, n\} \setminus \mathbf{i}$  and note that

$$\widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n) \leq \frac{n-d}{n} \widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n(\mathbf{i}')) + \frac{d}{n}.$$

Combining the two inequalities above, we get

$$Q_n(\widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n), d) \leq \widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n(\mathbf{i}')) + \frac{d}{n} + r_{d,n}.$$

Recalling  $t_\gamma(n) \in o(n)$ , by Property **Q3**,

$$\lim_{n \rightarrow \infty} \frac{t_\gamma(n)}{n} + r_{t_\gamma(n),n} = 0.$$

In addition, by **Q2**, we have  $r_{d,n} \leq r_{t_\gamma(n),n}$  for all  $d \leq t_\gamma(n)$ . Hence, we take  $n$  sufficiently large so that for all  $d \leq t_\gamma(n)$ ,

$$\frac{d}{n} + r_{d,n} \leq \frac{\varepsilon}{4},$$

and thus

$$Q_n(\widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n), d) \leq \widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n(\mathbf{i}')) + \frac{\varepsilon}{4}.$$

Therefore, for such an  $n$ ,

$$\begin{aligned} Q_n(\widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n), d) &> \text{err}(h_{S_n(\mathbf{i},*)}) + \frac{\varepsilon}{2} \\ \implies \widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n(\mathbf{i}')) &> \text{err}(h_{S_n(\mathbf{i},*)}) + \frac{\varepsilon}{4}. \end{aligned}$$

Now,

$$\begin{aligned} &\mathbb{P} \left[ \widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n(\mathbf{i}')) > \text{err}(h_{S_n(\mathbf{i},*)}) + \frac{\varepsilon}{4} \right] \tag{4.21} \\ &= \mathbb{E}_{S_n(\mathbf{i})} \left[ \mathbb{P}_{S_n(\mathbf{i}') | S_n(\mathbf{i})} \left[ \widehat{\text{err}}(h_{S_n(\mathbf{i},*)}, S_n(\mathbf{i}')) > \text{err}(h_{S_n(\mathbf{i},*)}) + \frac{\varepsilon}{4} \right] \right]. \end{aligned}$$

Since  $\mathbb{P}_{S_n(\mathbf{i}') | S_n(\mathbf{i})}$  is a product distribution, by Hoeffding's inequality we have that (4.21) is bounded above by  $e^{-2(n-d)(\frac{\varepsilon}{4})^2}$ . Since  $h_{S_n(\mathbf{i},*)}$  is invariant to permutations of  $\mathbf{i}$ 's entries, bounding (4.20) by a union bound over  $\mathbf{i}$  yields

$$pd \leq \binom{n}{d} e^{-2(n-d)(\frac{\varepsilon}{4})^2} \leq e^{d \log(\frac{en}{d}) - 2(n-d)(\frac{\varepsilon}{4})^2},$$

where we used  $\binom{n}{d} \leq (\frac{en}{d})^d$ . Selecting  $n$  large enough so that for all  $d \leq t_\gamma(n)$  we have  $d \log(en/d) \leq (n-d)(\frac{\varepsilon}{4})^2$  and  $d \leq n/4$ , we get the statement of the lemma.  $\square$

### 5. Essential separability is necessary for universal Bayes consistency.

Recall that a metric space  $(\mathcal{X}, \rho)$  is *essentially separable* (ES) if for every probability measure  $\mu$  on the Borel  $\sigma$ -algebra  $\mathcal{B}$ , the metric probability space  $(\mathcal{X}, \rho, \mu)$  is separable. In Theorem 4.2 we established that OptiNet is indeed universally Bayes-consistent (UBC) on all such spaces. As such, essential separability of a metric space is sufficient for the existence of a UBC learning rule on 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" [17]. As it turns out, the question of whether non-ES metric spaces even exist is rather subtle, and in fact is *independent of the ZFC axioms of set theory* (see Section 5.1). In other words, assuming that ZFC is consistent, it can neither be proven nor disproven from the axioms of ZFC that a non-ES metric space exists. The main contribution of this section is to show that for any non-ES metric space (if one exists), no learning rule is UBC on it.

**THEOREM 5.1.** *Let  $(\mathcal{X}, \rho)$  be a non-ES metric space equipped with the Borel  $\sigma$ -algebra  $\mathcal{B}$ . Then no UBC algorithm exists on  $(\mathcal{X}, \rho)$ .*

As an immediate corollary of the preceding theorems, we have the following result, which reveals that OptiNet is *optimistically* UBC (adopting the terminology of [30]), in the sense that the only required assumption on  $(\mathcal{X}, \rho)$  is that universally Bayes-consistent learning is *possible*.

COROLLARY 5.2. *OptiNet is UBC in every metric space on which there exists a UBC learning rule.*

REMARK 5.3. Theorem 5.1 is somewhat unusual, in that it identifies a learning 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 [3] for another recent result where learnability is independent of ZFC. ◀

In the next section, we provide necessary preliminaries. Theorem 5.1 is proved in Section 5.2.

5.1. *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 5.5 below). A thorough treatment of the latter may be found in [32]. All throughout, we work under ZFC.

*Notation.* The cardinality of a set  $A$  is denoted by  $|A|$ . We denote the first infinite (countable) cardinal by  $\aleph_0 = \omega$  and the cardinality of the continuum by  $\mathfrak{c} = 2^{\aleph_0} > \aleph_0$ . We write  $[A]^n$  for 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  $A$ .

*Non-trivial 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  $\mu$  is *non-trivial* if it vanishes on singletons:  $\mu(\{x\}) = 0, \forall x \in \mathcal{X}$ . On countable  $\mathcal{X}$ , all measures are trivial.

For a cardinal  $\kappa$ , a measure  $\mu$  is  $\kappa$ -*additive* if for any  $\beta < \kappa$  and any pairwise disjoint family  $\{A_\alpha \subseteq \mathcal{X} : \alpha < \beta\}$ ,

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

By definition, any measure is  $\aleph_1$ -additive (commonly known as  $\sigma$ -additivity), where  $\aleph_1$  is the least cardinal larger than  $\aleph_0$ . The following lemma follows directly from the definitions above.

LEMMA 5.4. *Let  $\mu$  be a non-trivial  $\kappa$ -additive measure on  $\mathcal{B}$ . Then any set  $A \in \mathcal{B}$  with  $|A| < \kappa$  has  $\mu(A) = 0$ .*

*Real-valued measurable cardinals.* Let  $\mathcal{X}$  be of cardinality  $\kappa$  and consider the measurable space  $\mathfrak{X}_\kappa := (\mathcal{X}, 2^\mathcal{X})$ . As already noted, there are no non-trivial measures on  $\mathfrak{X}_\kappa$  for any  $\kappa \leq \aleph_0$ . The question of whether there exists a cardinal  $\kappa$  such that  $\mathfrak{X}_\kappa$  admits a non-trivial measure leads us into the study of large cardinals.

DEFINITION 5.5. A cardinal  $\kappa$  is *real-valued measurable* if there exists a non-trivial  $\kappa$ -additive measure on  $\mathfrak{X}_\kappa$ . Any such measure is called a *witnessing measure* on  $\mathfrak{X}_\kappa$ .

By now, it is well known that the existence of real-valued measurable cardinals is independent of ZFC: assuming the consistency of ZFC, it can neither be proven nor disproven from the axioms of ZFC that real-valued measurable cardinals exist [32, §12].

The following characterization of ES metric spaces in terms of real-valued measurable cardinals is taken from [8]. Recall that a set  $D \subseteq \mathcal{X}$  is *discrete* if for any  $x \in D$  there exists  $r_x > 0$  such that  $B_{r_x}(x) \cap D = \{x\}$ .

THEOREM 5.6 ([8, Appendix III, Theorem 2]<sup>6</sup>). *Let  $\kappa_{\min}$  be the least real-valued measurable cardinal (if one exists). Then a metric space  $(\mathcal{X}, \rho)$  is ES if and only if every discrete  $D \subseteq \mathcal{X}$  has  $|D| < \kappa_{\min}$ .*

REMARK 5.7. The independence of the existence of real-valued measurable cardinals from ZFC was not known at the time [8] was published. In particular, Theorem 5.6 is stated there only for the case  $\kappa_{\min} \leq \mathfrak{c}$ . However, one can readily verify that the proof extends essentially verbatim (by replacing “atomless” with “non-trivial”) to the case  $\kappa_{\min} > \mathfrak{c}$  as well. ◀

*Atomless measures, atomic measures, and Ulam’s dichotomy.* As will become clear in Section 5.2, the argument showing the impossibility of UBC in non-ES metric spaces depends on whether  $\kappa_{\min} \leq \mathfrak{c}$  or  $\kappa_{\min} > \mathfrak{c}$ . This is manifested by what is known as Ulam’s dichotomy, which we now present.

Let  $\mu$  be a measure on  $\mathcal{B}$ . A set  $A \subseteq \mathcal{X}$  is an *atom* of  $\mu$  if  $\mu(A) > 0$  and for every measurable  $B \subseteq A$ , either  $\mu(B) = 0$  or  $\mu(B) = \mu(A)$ . A measure  $\mu$  is *atomless* if it has no atoms. In this case, we have that for any  $A \in \mathcal{B}$  with  $\mu(A) > 0$  there exists a measurable  $B \subset A$  with  $0 < \mu(B) < \mu(A)$ . Conversely,  $\mu$  is *purely atomic* if every  $A \in \mathcal{B}$  with  $\mu(A) > 0$  contains an atom. A measure  $\mu$  is *two-valued* if  $\mu(A) \in \{0, 1\}$  for all  $A \in \mathcal{B}$ . Clearly, such a 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  $\mu(P_j) = 1$ .

Let  $\kappa$  be a real-valued measurable cardinal. We say that  $\kappa$  is *two-valued measurable* if there is a two-valued witnessing measure on  $\mathfrak{X}_\kappa$ ; we say that  $\kappa$  is *atomlessly measurable* if there is an atomless witnessing measure on  $\mathfrak{X}_\kappa$ . In 1930, Ulam established the following dichotomy (see [20, §543]).

<sup>6</sup>Appendix III in [8] was omitted from later editions.

THEOREM 5.8 (Ulam’s Dichotomy [58]). *Let  $\kappa$  be a real-valued measurable cardinal. Then*

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

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

5.2. *UBC is impossible in non-ES metric spaces.* In this section we prove the following theorem, from which Theorem 5.1 readily follows.

THEOREM 5.9. *Let  $(\mathcal{X}, \rho)$  be a non-ES metric space and let  $\text{Alg}$  be any (possibly random) learning algorithm mapping samples  $S \in (\mathcal{X} \times \{0, 1\})^{<\omega}$  to classifiers  $\text{Alg}(S) \in \{0, 1\}^\mathcal{X}$ . Let  $\kappa_{\min}$  be the least real-valued measurable cardinal.*

- (I) *If  $\kappa_{\min} \leq \mathfrak{c}$  then there exists a measure  $\bar{\mu}$  on  $\mathcal{X} \times \{0, 1\}$  (w.r.t. the Borel sets induced by  $\rho$ ), a measurable classifier  $h^* : \mathcal{X} \rightarrow \{0, 1\}$ , and an  $\varepsilon > 0$  such that, for all  $n \in \mathbb{N}$  and  $S \sim \bar{\mu}^n$ ,*

$$\text{err}_{\bar{\mu}}(h^*) = 0 < \varepsilon \leq \mathbb{E}[\text{err}_{\bar{\mu}}(\text{Alg}(S))]. \quad (5.1)$$

- (II) *If  $\kappa_{\min} > \mathfrak{c}$  then there are two distinct measures on  $\mathcal{X} \times \{0, 1\}$  such that  $\text{Alg}$  is not Bayes consistent under both.*

*Taken together, it follows that for any non-ES metric space, there does not exist a universally Bayes-consistent learning algorithm.*

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

REMARK 5.10. 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  $\mu$  is the marginal of  $\bar{\mu}$  on  $\mathcal{X}$ , and  $\mathcal{B}(L^1(\mu))$  is the Borel  $\sigma$ -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 criterion in Remark 5.10, 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  $\sigma$ -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  $\mathbb{E}\left[|\hat{h}_S(X) - f(X)| \middle| S\right]$  is well-defined and measurable. Therefore, for any  $\varepsilon > 0$ , the event that  $\mathbb{E}\left[|\hat{h}_S(X) - f(X)| \middle| S\right] \leq \varepsilon$  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  $\rho$  (which are measurable since, by definition,  $\rho$  induces the topology generating the Borel  $\sigma$ -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 5.9 are satisfied by OptiNet.  $\blacktriangleleft$

REMARK 5.11. In [10, Section 2.1], the authors define the metric space  $(\mathcal{X}, \rho)$ , 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  $\mu$ , 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 — and use the 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  $(\mathcal{X}, \rho, \mu)$  is a *metric probability space* — that is, that the  $\sigma$ -algebra supporting  $\mu$  is generated by the open sets of  $\rho$ . Indeed, the Borel  $\sigma$ -algebra generated by  $\rho$  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  $\sigma$ -algebra can be extended to all subsets of  $[0, 1]$  — a statement known to be equivalent to the existence of a real-valued measurable cardinal  $\leq \mathfrak{c}$  [58]. Another objection is that, under any reasonable notion of *dimension*, the metric space  $(\mathcal{X}, \rho)$  would be considered  $\mathfrak{c}$ -dimensional rather than finite-dimensional.  $\blacktriangleleft$

PROOF OF THEOREM 5.9. Since  $(\mathcal{X}, \rho)$  is non-ES, Theorem 5.6 implies that there exists a discrete set  $D$  with  $|D| = \kappa_{\min}$ . Let  $\mathfrak{X}_{|D|} := (D, 2^D)$  and note that by Lemma B.1,  $2^D \subseteq \mathcal{B}$ . Hence, it suffices to construct the required adversarial measures in case (I) and case (II) of Theorem 5.9 on  $D \times \{0, 1\}$ . That being the case, from now on we set without loss of generality  $\mathcal{X} := D$  and  $\mathcal{B} := 2^D$ . We proceed by considering the two cases (I) and (II) separately.

(I) *The case  $\kappa_{\min} \leq \mathfrak{c}$ .* By Ulam’s dichotomy in Theorem 5.8,  $|\mathcal{X}| = \kappa_{\min}$  is atomlessly measurable, so there exists an atomless witnessing measure  $\mu$  on  $\mathcal{B}$ . Fix such a  $\mu$  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  $(\mathcal{U}, \Delta)$ , where  $\mathcal{U} \subseteq \mathcal{B}$  is the quotient  $\sigma$ -algebra under the equivalence relation  $A \sim B \Leftrightarrow \Delta(A, B) = 0$ . The measure  $\mu$  induces the corresponding functional  $\tilde{\mu} : \mathcal{U} \rightarrow [0, 1]$  which agrees with  $\mu$  on the equivalence

classes. The following is proved in Appendix B by an application of Gitik-Shelah Theorem [22].

LEMMA 5.12. *Let  $\mathcal{X}$  be a set of atomlessly-measurable cardinality  $\kappa$  and let  $\mu$  be a witnessing measure on  $\mathfrak{X}_\kappa$ . Let  $(\mathcal{U}, \Delta)$  be as above. Then there exist  $\varepsilon > 0$  and a subset  $\mathcal{H}_\varepsilon \subseteq \mathcal{U}$  of cardinality  $|\mathcal{H}_\varepsilon| = \kappa$  that is  $\varepsilon$ -separated:*

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

Using Lemma 5.12, there exists  $\varepsilon > 0$  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  $\mathcal{H}_\varepsilon$  has cardinality  $\kappa_{\min}$ : that is, the same cardinality as  $\mathcal{X}$ . Since  $\kappa_{\min}$  is atomlessly-measurable, there exists an atomless witnessing measure  $\pi$  on  $(\mathcal{H}_\varepsilon, 2^{\mathcal{H}_\varepsilon})$ . We will consider constructing the distribution  $\bar{\mu}$  randomly, taking  $\mu$  fixed as the marginal on  $\mathcal{X}$  and taking  $h^*$  randomly  $\pi$ -distributed (independent of the  $x$  data, as 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

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

Now take  $\mathbf{S}_x = \{x_1, x_2, \dots\}$  to be an i.i.d.  $\mu$ -distributed sequence, and let  $h^* \sim \pi$  independently of  $\mathbf{S}_x$ . Let  $\bar{\mu}$  have marginal  $\mu$  over  $\mathcal{X}$  and  $\bar{\mu}(\{(x, h^*(x)) : x \in \mathcal{X}\}) = 1$ ; 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  $\hat{h}_{\mathbf{S}_x, \mathbf{y}}$  is a  $\mathcal{B}(L^1(\mu))$ -measurable random variable, and  $h^*$  is also  $\mathcal{B}(L^1(\mu))$ -measurable (its distribution is  $\pi$ , which is defined on this  $\sigma$ -algebra). Therefore, the value  $\mu(\{x' : \hat{h}_{\mathbf{S}_x, \mathbf{y}}(x') \neq h^*(x')\})$  is a measurable random variable, equal (a.s.) to  $\text{err}_{\bar{\mu}}(\hat{h}_{\mathbf{S}_x, \mathbf{y}})$ .

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

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

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

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

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

so that altogether we have

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

In particular, this also implies there exist fixed choices of  $h^*$  for which (5.1) holds. This completes the proof for the case (I).

(II) *The case  $\kappa_{\min} > \mathfrak{c}$ .* By Ulam's dichotomy in Theorem 5.8,  $|\mathcal{X}| = \kappa_{\min}$  is two-valued measurable, so there exists a two-valued witnessing measure  $\mu$  on  $\mathcal{B}$ . As the following lemma shows,  $\mu$  can be taken to further satisfy a key homogeneity property. As shown in Appendix C, the lemma follows by combining Theorems 10.20, 10.22 in [32] and Ulam's Theorem 5.8.

LEMMA 5.13. *Let  $\mathcal{X}$  be of two-valued measurable cardinality  $\kappa$  and let  $\mathfrak{X}_\kappa = (\mathcal{X}, 2^\mathcal{X})$ . Then, there is a witnessing measure  $\mu$  on  $\mathfrak{X}_\kappa$  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  $\mu$  be a two-valued witnessing measure on  $\mathcal{B} = 2^\mathcal{X}$  as furnished by Lemma 5.13. For a label  $y \in \{0, 1\}$  and any two-valued witnessing measure  $\phi$ , let  $\bar{\phi}_y$  be the measure over  $\mathcal{X} \times \{0, 1\}$  with  $\phi$  as its marginal over  $\mathcal{X}$  and

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

For  $\mu$  as above, and any other two-valued witnessing measure  $\phi$ , define the mixture measure

$$\lambda_\phi := \frac{2}{3}\bar{\phi}_1 + \frac{1}{3}\bar{\mu}_0. \tag{5.2}$$

We will show that there exists a two-valued witnessing measure  $\nu \neq \mu$  such that Alg cannot be Bayes consistent on both  $\lambda_\mu$  and  $\lambda_\nu$ . To this end we will use the following properties of the mixture  $\lambda_\phi$ , proved in Appendix C.

LEMMA 5.14. *Let  $\nu \neq \mu$  be any two distinct two-valued measures on  $(\mathcal{X}, \mathcal{B})$  and let  $\lambda_\phi$  with  $\phi \in \{\mu, \nu\}$  be as in (5.2). Let  $h^* : \mathcal{X} \rightarrow \{0, 1\}$  be any Bayes-optimal classifier on  $\lambda_\phi$ .*

(i) If  $\phi = \mu$  then  $h^*$  achieves the Bayes-optimal  $\text{err}_{\lambda_\mu}(h^*) = \frac{1}{3}$  iff

$$\mathbb{E}_{X \sim \mu} [h^*(X)] = 1.$$

(ii) If  $\phi = \nu$  then  $h^*$  achieves the Bayes-optimal  $\text{err}_{\lambda_\nu}(h^*) = 0$  iff

$$\mathbb{E}_{X \sim \mu} [h^*(X)] = 0 \quad \text{and} \quad \mathbb{E}_{X \sim \nu} [h^*(X)] = 1.$$

Let  $\text{Alg} : (\mathcal{X} \times \{0, 1\})^{<\omega} \rightarrow 2^{\mathcal{X}}$  be any (possibly randomized) learning algorithm, and denote  $\hat{h}_S$  the classifier output for data set  $S$ ; for  $S$  and  $X$  independent samples from Borel measures on  $\mathcal{X}$ , we suppose  $\hat{h}_S(X)$  is a measurable random variable (by definition of learning algorithm; see Remark 5.10). Let  $\nu \neq \mu$  be a two-valued witnessing measure to be chosen below. Consider the quantity

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

In the case of randomized Alg, also add an innermost expectation over the independent randomness of Alg in the above expression. By Lemma 5.14, for Alg to be Bayes consistent on both  $\lambda_\nu$  and  $\lambda_\mu$  we must have

$$L_n^\phi \xrightarrow{n \rightarrow \infty} \delta_{\mu, \phi}, \quad \phi \in \{\mu, \nu\}. \quad (5.3)$$

So to prove the claim it suffices to show that we can choose  $\nu := \nu(\mu, \text{Alg})$  such that (5.3) 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

$$n_1 := n_1(\mathbf{Y}_n) = \sum_{i=1}^n Y_i \quad \text{and} \quad n_0 := n - n_1$$

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  $\mu$  and  $\phi$ . 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

$$\begin{aligned} L_n^\phi &= \mathbb{E}_{S_n \sim (\lambda_\phi)^n} \left[ \mathbb{E}_{X \sim \mu} [\hat{h}_{S_n}(X)] \right] \\ &= \mathbb{E}_{\mathbf{Y}_n} \mathbb{E}_{\mathbf{X}_n | \mathbf{Y}_n} \left[ \mathbb{E}_{X \sim \mu} [\hat{h}_{(\mathbf{X}_n, \mathbf{Y}_n)}(X)] \right] \\ &= \mathbb{E}_{\mathbf{Y}_n} \mathbb{E}_{\mathbf{X}_n^1 \sim \phi^{n_1}} \mathbb{E}_{\mathbf{X}_n^0 \sim \mu^{n_0}} \left[ \mathbb{E}_{X \sim \mu} [\hat{h}_{((\mathbf{X}_n^0, \mathbf{X}_n^1), \mathbf{Y}_n)}(X)] \right]. \end{aligned} \quad (5.4)$$

Towards applying Lemma 5.13, 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  $\phi$  vanishes on singletons, all instances in  $\mathbf{X}$  are distinct with probability one,

$$\mathbb{P}_{\mathbf{X} \sim \phi^k} [|\mathbf{W}_{\mathbf{X}}| = k] = 1. \quad (5.5)$$

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  $\pi$  of the instances in  $W$ . Then, by (5.5) and the fact that  $\phi^k$  is a product measure, we have that for any measurable function  $f : \mathcal{X}^k \rightarrow [0, 1]$  the following symmetrization holds,

$$\mathbb{E}_{\mathbf{X} \sim \phi^k} [f(\mathbf{X})] = \mathbb{E}_{\mathbf{X} \sim \phi^k} \left[ \mathbb{E}_{\mathbf{X}' \sim \Pi(W_{\mathbf{X}})} [f(\mathbf{X}')] \mid |\mathbf{W}_{\mathbf{X}}| = k \right]. \quad (5.6)$$

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

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

In the case of randomized Alg, we also include an innermost conditional expectation over the value of  $\hat{h}_{((\mathbf{X}_n^0, \mathbf{X}^1), \mathbf{Y}_n)}(X)$  given  $\mathbf{X}_n^0, \mathbf{X}^1, \mathbf{Y}_n, X$ . Putting this in (5.4) while using (5.5) and (5.6),

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

By the choice of  $\mu$ , Lemma 5.13 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

$$U = \bigcap_{n \in \mathbb{N}} \bigcap_{\mathbf{Y}_n \in \{0, 1\}^n} U_{\mathbf{Y}_n}. \quad (5.7)$$

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

$$F_{\mathbf{Y}_n}(W) = C_{\mathbf{Y}_n}, \quad \forall n \in \mathbb{N}, \quad \forall \mathbf{Y}_n \in \{0, 1\}^n, \quad \forall W \in [U]^{n_1}. \quad (5.8)$$

In addition, by Lemma C.1,  $\mu(U) = 1$ .

We are now in position to choose  $\nu := \nu(\mu, \text{Alg})$ . By Lemma 5.4, we may split  $U$  in (5.7) into two disjoint sets  $B$  and  $U \setminus B$  such that  $|B| = |U \setminus B| = |\mathcal{X}|$ . Since  $\mu$  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  $\nu'$  on  $(B, 2^B)$  with  $\nu'(B) = 1$ . Extend  $\nu'$  to a measure  $\nu$  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  $\forall k \in \mathbb{N}$ ,

$$\Pr_{\mathbf{X} \sim \phi^k} [W_{\mathbf{X}} \in [U]^k \mid |\mathbf{W}_{\mathbf{X}}| = k] = 1.$$

So, for  $\phi \in \{\mu, \nu\}$ ,

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

where we used (5.8) 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  $\phi$ , we conclude that  $L_n^\mu = L_n^\nu$  for all  $n \in \mathbb{N}$ . However by (5.3), for Alg to be Bayes consistent on  $\lambda_\mu$  and  $\lambda_\nu$  we must have

$$L_n^\mu \xrightarrow{n \rightarrow \infty} 1 \quad \text{and} \quad L_n^\nu \xrightarrow{n \rightarrow \infty} 0.$$

Thus Alg cannot be Bayes consistent on both  $\lambda_\mu$  and  $\lambda_\nu$ .  $\square$

**6. Discussion and future work.** We showed that OptiNet, which is a computationally efficient multiclass learning algorithm, is universally strongly Bayes consistent in all metric spaces that admit such a learner. As such, it is optimistically universal (in the terminology of [30]). A natural open problem is whether similar guarantees can be provided for OptiNet for the case of real-valued regression. Here, we note that if the regression loss minimizer can be achieved on the sample using labels occurring in the sample, as in the case of the  $\ell_1$  loss, and if in addition the label set is bounded, then our techniques carry over seamlessly, by setting the label in each Voronoi cell to the label that minimizes the sample loss in the cell. However, this does not necessarily hold in the general case. In particular, for agnostic regression with the  $\ell_2$  loss, the compression size cannot be controlled [15], nor can this be done for any  $\ell_p$  losses for  $p \notin \{1, \infty\}$ . In addition, our proof of almost-sure convergence, which combines an exponential inequality with the Borel-Cantelli lemma, does not go through for unbounded losses. Thus, solving this open problem might require new techniques.

*Acknowledgments.* We thank Robert Furber, Iosif Pinelis, and Menachem Kojman for providing some set-theoretical background. 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).

#### APPENDIX A: AUXILIARY LEMMAS FOR SECTION 4

LEMMA A.1. *For every metric probability space  $(\mathcal{X}, \rho, \mu)$ , the set of Lipschitz functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  is dense in  $L_1(\mu) = \{f : \int |f| d\mu < \infty\}$ . In other words, for any  $\varepsilon > 0$  and  $f \in L_1(\mu)$ , there is an  $L < \infty$  and an  $L$ -Lipschitz function  $g \in L_1(\mu)$  such that  $\int |f - g| d\mu < \varepsilon$ .*

PROOF. We believe this fact to be classical, but were unable to locate an appropriate citation, so for completeness we include a brief proof. The proof follows

closely that of a weaker result from [34, Section 37, Theorem 2]. It relies on the fact that, for any probability measure  $\mu$  on a Borel  $\sigma$ -algebra  $\mathcal{B}$ ,  $\mu$  is *regular* [33, Theorem 17.10]. In particular, for every  $A \in \mathcal{B}$ ,  $\mu(A) = \sup_{F \in \mathcal{F}: F \subseteq A} \mu(F)$ , where  $\mathcal{F}$  is

the *closed* sets (under the topology that generates  $\mathcal{B}$ ).

For any  $A \in \mathcal{B}$  and  $\varepsilon > 0$ , regularity implies that there is an  $F \in \mathcal{F}$  with  $F \subseteq A$  and  $\mu(A \setminus F) < \varepsilon/2$ . Now denote  $G_r = \bigcup_{x \in F} B_r(x)$ . Since  $\mathcal{X} \setminus F$  is open, for any  $x' \notin F$ , there is an  $r > 0$  with  $B_r(x') \subseteq \mathcal{X} \setminus F$ , and hence  $x' \notin G_r$ . Together with monotonicity of  $G_r$  in  $r$ , this implies  $G_r \setminus F \rightarrow \emptyset$  as  $r \rightarrow 0$ . Thus, by continuity of probability measures, there is an  $r > 0$  such that  $\mu(G_r \setminus F) < \varepsilon/2$ . Furthermore, for this  $r$ ,  $G_r \supseteq F$  and  $G_r$  is a union of open sets, hence open. Thus, denoting  $F_r = \mathcal{X} \setminus G_r$ ,  $F_r$  is a closed set, disjoint from  $F$ , and (by definition of  $G_r$ ) satisfies  $\inf_{x \in F, x' \in F_r} \rho(x, x') \geq r > 0$ .

Now define

$$g_{A,\varepsilon}(x) = \frac{\inf_{x' \in F_r} \rho(x', x)}{\inf_{x' \in F_r} \rho(x', x) + \inf_{x' \in F} \rho(x', x)}.$$

In particular, note that  $g_{A,\varepsilon}(x) = 1$  for  $x \in F$ ,  $g_{A,\varepsilon}(x) = 0$  for  $x \in F_r$ , and every other  $x$  has  $g_{A,\varepsilon}(x) \in [0, 1]$ . This implies  $\{x : g_{A,\varepsilon}(x) < 1, \mathbf{1}_A(x) = 1\} \subseteq A \setminus F$  and  $\{x : g_{A,\varepsilon}(x) > 0, \mathbf{1}_A(x) = 0\} \subseteq (\mathcal{X} \setminus F_r) \setminus A = G_r \setminus A \subseteq G_r \setminus F$ , so that

$$\int |\mathbf{1}_A - g_{A,\varepsilon}| d\mu \leq \mu(A \setminus F) + \mu(G_r \setminus F) < \varepsilon.$$

Furthermore, since  $F$  and  $F_r$  are  $r$ -separated,  $g_{A,\varepsilon}$  is  $\frac{1}{r}$ -Lipschitz, and since  $g_{A,\varepsilon}$  is bounded we also have  $g_{A,\varepsilon} \in L_1(\mu)$ . Thus, we have established the desired result for indicator functions.

To extend this to all of  $L_1(\mu)$ , we use the “standard machinery” technique. By definition of Lebesgue integration, for any  $f \in L_1(\mu)$  and  $\varepsilon > 0$ , there exists a finite simple function  $f_\varepsilon$  with  $\int |f - f_\varepsilon| d\mu < \varepsilon/2$ : that is, there is an  $n \in \mathbb{N}$ ,  $a_1, \dots, a_n \in \mathbb{R}$ , and  $A_1, \dots, A_n \in \mathcal{B}$  with  $f_\varepsilon(x) = \sum_{i=1}^n a_i \mathbf{1}_{A_i}(x)$ . Now let  $a^* = \max\{|a_1|, \dots, |a_n|, 1\}$  and denote  $\varepsilon' = \varepsilon/(2na^*)$ . By the above, for each  $i \in \{1, \dots, n\}$ , there exists a Lipschitz function  $g_{A_i, \varepsilon'} \in L_1(\mu)$  with  $\int |\mathbf{1}_{A_i} - g_{A_i, \varepsilon'}| d\mu < \varepsilon'$ . Therefore, denoting  $g = \sum_{i=1}^n a_i g_{A_i, \varepsilon'}$ , we have

$$\int |f_\varepsilon - g| d\mu \leq a^* \sum_{i=1}^n \int |\mathbf{1}_{A_i} - g_{A_i, \varepsilon'}| d\mu < a^* n \varepsilon' = \varepsilon/2.$$

Together we have that  $\int |f - g| d\mu \leq \int |f - f_\varepsilon| d\mu + \int |f_\varepsilon - g| d\mu < \varepsilon$ . Since a finite linear combination of Lipschitz functions is still Lipschitz, this establishes the claim for all  $f \in L_1(\mu)$ .  $\square$

**PROOF OF LEMMA 4.1.** Let  $\mathbf{X}(\gamma)$  be any  $\gamma$ -net of  $S_n = (X_1, \dots, X_n)$  and let  $\mathcal{A} = \{A_1, A_2, \dots\}$  be a fixed countable partition of  $\mathcal{X}'$  (for separable  $\mathcal{X}' \subseteq \mathcal{X}$  of  $\mu(\mathcal{X}') = 1$ ) with

$$\text{diam}(\mathcal{A}) = \sup_{i \in \mathbb{N}} (\text{diam}(A_i)) < \gamma,$$

which exists by the separability assumption. Denote the number of occupied cells in  $\mathcal{A}$  by

$$U_n(X_1, \dots, X_n) = \sum_{A_i \in \mathcal{A}} \mathbf{1}[S_n \cap A_i \neq \emptyset].$$

Since  $\mathbf{X}(\gamma)$  is a  $\gamma$ -net of  $S_n$ , any cell  $A_i \in \mathcal{A}$  contains at most one  $X \in \mathbf{X}(\gamma)$ . Hence,

$$M_n(\gamma) = |\mathbf{X}(\gamma)| \leq U_n(X_1, \dots, X_n).$$

So it suffices to bound  $U_n$ . To this end, denote by  $i(X)$  the cell in  $\mathcal{A}$  such that  $X \in A_{i(X)}$ . Then,

$$\begin{aligned} \mathbb{E}[U_n(X_1, \dots, X_n)] &= \sum_{j=1}^n \mathbb{P}[A_{i(X_j)} \cap \{X_1, \dots, X_{j-1}\} = \emptyset] \\ &= \sum_{j=1}^n \mathbb{E}[1 - \mu(\cup_{k=1}^{j-1} A_{i(X_k)})]. \end{aligned}$$

Since

$$\lim_{j \rightarrow \infty} \mathbb{E} \left[ \mu \left( \bigcup_{k=1}^j A_{i(X_k)} \right) \right] = \lim_{n \rightarrow \infty} \mathbb{E} \left[ \mu \left( \bigcup_{A_i \in \mathcal{A}: A_i \cap S_n \neq \emptyset} A_i \right) \right] = 1, \quad (\text{A.1})$$

we have

$$\mathbb{E}[U_n(X_1, \dots, X_n)] \in o(n).$$

Let

$$t_\gamma(n) = \mathbb{E}[U_n(X_1, \dots, X_n)] + \sqrt{n \log n} \in o(n).$$

Since  $U_n$  is 1-Lipschitz with respect to the Hamming distance, McDiarmid's inequality implies

$$\mathbb{P}[U_n(X_1, \dots, X_n) \geq t_\gamma(n)] \leq 1/n^2,$$

concluding the proof.  $\square$

**PROOF OF LEMMA 4.4.** Let  $\mathcal{A} = \{A_1, A_2, \dots\}$  be a fixed countable partition of  $\mathcal{X}$  with  $\text{diam}(\mathcal{A}) < \gamma$  as in the proof of Lemma 4.1. Consider the random variable

$$F_{\mathcal{A}}(S_n) = 1 - \mu(\cup_{A_i \in \mathcal{A}: A_i \cap S_n \neq \emptyset} A_i),$$

corresponding to the total mass of all cells not hit by the sample  $S_n$ . Since  $\text{diam}(\mathcal{A}) < \gamma$ , we have that  $A_i \subseteq B_\gamma(x)$  for all  $x \in A_i$ . Hence, with probability 1

$$L_\gamma(S_n) = 1 - \mu(\text{UB}_\gamma(S_n)) \leq F_{\mathcal{A}}(S_n),$$

whence

$$\mathbb{P}[L_\gamma(S_n) \geq \mathbb{E}[F_{\mathcal{A}}(S_n)] + t] \leq \mathbb{P}[F_{\mathcal{A}}(S_n) \geq \mathbb{E}[F_{\mathcal{A}}(S_n)] + t].$$

Invoking the concentration bound for the missing mass in [5, Theorem 1],

$$\mathbb{P}[F_{\mathcal{A}}(S_n) \geq \mathbb{E}[F_{\mathcal{A}}(S_n)] + t] \leq \exp(-nt^2)$$

and observing that, by (A.1),  $\lim_{n \rightarrow \infty} \mathbb{E}[F_{\mathcal{A}}(S_n)] = 0$ , shows that the choice  $u_\gamma(n) := \mathbb{E}[F_{\mathcal{A}}(S_n)]$  verifies the properties claimed.  $\square$

#### APPENDIX B: AUXILIARY LEMMAS FOR SECTION 5.2 – CASE (I)

LEMMA B.1. *Suppose that  $U$  is a discrete subset of a metric space  $(\mathcal{X}, \rho)$ . Then every  $E \subseteq U$  is Borel.*

PROOF. By discreteness, for each  $x \in U$  there is an  $r_x > 0$  such that  $B_{r_x}(x) \cap U = \{x\}$ . The latter property is satisfied by any other  $0 < r < r_x$ . Further, for any  $n \in \mathbb{N}$ , the set  $V_n := \cup_{x \in E} B_{r_x/n}$  is open and

$$E = \bigcap_{n \in \mathbb{N}} V_n,$$

whence  $E$  is Borel.  $\square$

PROOF OF LEMMA 5.12. For  $\varepsilon > 0$ , consider the family of all  $\varepsilon$ -separated subsets of  $\mathcal{U}$ ,

$$\mathcal{F}_\varepsilon = \{F \subseteq \mathcal{U} : \Delta(A, B) \geq \varepsilon, A \neq B \in F\}.$$

Let  $\mathcal{F}_\varepsilon^{\max}$  be the set of maximal elements in  $\mathcal{F}_\varepsilon$ , which by Zorn's lemma is non-empty. Any  $F \in \mathcal{F}_\varepsilon^{\max}$  is an  $\varepsilon$ -net of  $\mathcal{U}$  by maximality. Let  $\{\varepsilon_i\}_{i \in \mathbb{N}}$  be a sequence such that  $\varepsilon_i > 0$  and  $\lim_{i \rightarrow \infty} \varepsilon_i = 0$  and let  $\{D_i\}_{i \in \mathbb{N}}$  be such that  $D_i \in \mathcal{F}_{\varepsilon_i}^{\max}$  for all  $i \in \mathbb{N}$ . Clearly, the set  $D = \cup_{i \in \mathbb{N}} D_i$  is dense in  $\mathcal{U}$ . Moreover, by [2, Lemma 2],

$$|D| = \sup_{i \in \mathbb{N}} |D_i| = d(\mathcal{U}),$$

where  $d(\mathcal{U})$  is the density of  $(\mathcal{U}, \Delta)$  (namely, the smallest cardinality of any subset of  $\mathcal{U}$  which is dense for the metric space). Hence, for any cardinal  $\alpha < d(\mathcal{U})$ , there is a finite  $i \in \mathbb{N}$  such that  $D_i$  is an  $\varepsilon_i$ -separated set with cardinality  $|D_i| \geq \alpha$ . Thus, to prove the Lemma it suffices to show that  $d(\mathcal{U}) > \kappa$ .

To show that  $d(\mathcal{U}) > \kappa$ , consider the measure algebra  $(\mathcal{U}, \tilde{\mu})$ . Since  $\tilde{\mu}$  is totally finite, the topology of the measure algebra is the same topology generated by the metric space  $(\mathcal{U}, \Delta)$  [20, 323A(d)]. So the density of the measure algebra topology is  $d(\mathcal{U})$  as well. The Maharam type  $\tau(\mathcal{U})$  of  $(\mathcal{U}, \tilde{\mu})$  is defined as the smallest cardinality of any subset of  $\mathcal{U}$  which generates the topology of  $(\mathcal{U}, \tilde{\mu})$  [20, 331E-F]. Since  $\mathcal{U}$  is infinite, [20, 521O(ii)] implies  $d(\mathcal{U}) = \tau(\mathcal{U})$ . Since  $\mu$  is a finite, atomless, and  $\kappa$ -additive measure on  $(\mathcal{X}, 2^{\mathcal{X}})$ , Gitik-Shelah Theorem [20, 543F] implies that  $(\mathcal{U}, \tilde{\mu})$  has Maharam type  $\tau(\mathcal{U}) > \kappa$ . Hence,  $d(\mathcal{U}) = \tau(\mathcal{U}) > \kappa$ .  $\square$

## APPENDIX C: AUXILIARY LEMMAS FOR SECTION 5.2 – CASE (II)

LEMMA C.1. *Let  $\mu$  be a measure. For any countable family  $\{U_i\}_{i=1}^\infty$  with  $\mu(U_i) = 1, \forall i \in \mathbb{N}$ , we have  $\mu(\bigcap_{i=1}^\infty U_i) = 1$ .*

PROOF OF LEMMA C.1. Let  $U_i^c = \mathcal{X} \setminus U_i$  and  $V_i = U_i^c \setminus \{\bigcup_{j<i} U_j^c\}$ . Then  $\mu(V_i) = 0$ ,  $\{V_i\}_{i=1}^\infty$  are pairwise disjoint, and  $(\bigcap_{i=1}^\infty U_i)^c = \bigcup_{i=1}^\infty U_i^c = \bigcup_{i=1}^\infty V_i$ . Thus,  $\mu(\bigcap_{i=1}^\infty U_i) = 1 - \mu(\bigcup_{i=1}^\infty V_i) = 1 - \sum_{i=1}^\infty \mu(V_i) = 1$ .  $\square$

LEMMA C.2. *Let  $\nu \neq \mu$  be any distinct two-valued measures on  $(\mathcal{X}, \mathcal{B})$ . Then there exists  $B \subseteq \mathcal{X}$  such that  $\nu(B) = \mu(\mathcal{X} \setminus B) = 1$ .*

PROOF OF LEMMA C.2. By definition,  $\mu$  and  $\nu$  are distinct if  $\exists A \subseteq \mathcal{X}$  such that  $\mu(A) \neq \nu(A)$ . Since  $\mu$  and  $\nu$  are two-valued, we must have that  $\mu(A) = 1 - \nu(A) \in \{0, 1\}$ . In addition, either  $\mu(A) = 1$  or  $\mu(\mathcal{X} \setminus A) = 1$ . Assuming without loss of generality that  $\mu(A) = 1$ , the set  $B = \mathcal{X} \setminus A$  satisfies the required properties.  $\square$

PROOF OF LEMMA 5.13. The required measure  $\mu$  is taken as a witnessing measure with the additional property of being *normal*.<sup>7</sup> For our needs, it suffices that a normal measure  $\mu$  exists on  $\mathfrak{X}_\kappa$ , a fact proved in [32, Theorem 10.20]. To establish the homogeneity property of  $\mu$  we apply [32, Theorem 10.22], which in the terminology of current paper takes the following form.

THEOREM C.3 ([32], Theorem 10.22). *Let  $\mathcal{X}$  be of two-valued measurable cardinality  $\kappa$ , let  $\mu$  be a normal measure on  $\mathfrak{X}_\kappa$ , and let  $F : [\mathcal{X}]^{<\omega} \rightarrow \mathcal{R}$  with  $|\mathcal{R}| < \kappa$ . Then, there exists a set  $U \subseteq \mathcal{X}$  with  $\mu(U) = 1$  that is homogeneous for  $F$ .*

Since  $\kappa$  is two-valued measurable, Ulam's dichotomy implies  $\kappa > \mathfrak{c} = |\mathbb{R}|$ . An application of Theorem C.3 with  $\mathcal{R} = \mathbb{R}$  completes the proof of the Lemma.  $\square$

PROOF OF LEMMA 5.14. Assume first that  $\phi = \mu$  and let  $(X, Y) \sim \lambda_\mu$ . Note that  $Y \sim \text{Bernoulli}(2/3)$  and  $X \sim \mu$  is independent of  $Y$ . Thus, for any classifier  $h : \mathcal{X} \rightarrow \{0, 1\}$ ,

$$\begin{aligned} \text{err}_{\lambda_\mu}(h) &= \mathbb{P}[h(X) \neq Y] \\ &= \frac{2}{3} \cdot \mathbb{P}[h(X) = 0 \mid Y = 1] + \frac{1}{3} \cdot \mathbb{P}[h(X) = 1 \mid Y = 0] \\ &= \frac{2}{3} \cdot \mu(h(X) = 0) + \frac{1}{3} \cdot \mu(h(X) = 1) \\ &= \frac{2}{3} - \frac{1}{3} \cdot \mu(h(X) = 1) \geq \frac{1}{3}. \end{aligned}$$

Hence, the Bayes-optimal error is 1/3 (as demonstrated by the classifier  $h^*(x) = 1$ ) and is achieved if and only if  $\mu(h(X) = 1) = 1$ .

<sup>7</sup>In the terminology of [32, Chapter 10], a two-valued witnessing measure  $\mu$  on  $\mathcal{X}$  is equivalent to a  $\kappa$ -complete non-principal ultrafilter on  $2^\mathcal{X}$  consisting of all sets with measure 1 under  $\mu$ . The latter is normal if the ultrafilter is also closed under diagonal intersection.

Assume now that  $\phi = \nu \neq \mu$ . Then,

$$X|Y \sim \begin{cases} \nu, & \text{if } Y = 1; \\ \mu, & \text{if } Y = 0. \end{cases}$$

Thus, the error of a classifier  $h$  is

$$\begin{aligned} \text{err}_{\lambda_\nu}(h) &= \frac{2}{3} \cdot \mathbb{P}[h(X) = 0 | Y = 1] + \frac{1}{3} \cdot \mathbb{P}[h(X) = 1 | Y = 0] \\ &= \frac{2}{3} \cdot \nu(h(X) = 0) + \frac{1}{3} \cdot \mu(h(X) = 1). \end{aligned}$$

Since both  $\mu$  and  $\nu$  are two-valued, Lemma C.2 implies  $\exists B \subseteq \mathcal{X}$  such that

$$\nu(B) = \mu(\mathcal{X} \setminus B) = 1.$$

Thus, the Bayes-optimal error is 0 (as demonstrated by  $h^*(x) = \mathbf{1}[x \in B]$ ) and is achieved if and only if  $\nu(h(X) = 0) = 0$  and  $\mu(h(X) = 1) = 0$ .  $\square$

## 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](#)
- [2] BARBATI, A., COSTANTINI, C. et al. (1997). On the density of the hyperspace of a metric space. *Commentationes Mathematicae Universitatis Carolinae* **38** 349–360.
- [3] BEN-DAVID, S., HRUBES, P., MORAN, S., SHPILKA, A. and YEHUDAYOFF, A. (2019). Learnability can be undecidable. *Nature Machine Intelligence* **1** 44–48.
- [4] BEREND, D. and KONTOROVICH, A. (2012). The missing mass problem. *Statistics & Probability Letters* **82** 1102–1110.
- [5] BEREND, D. and KONTOROVICH, A. (2013). On the concentration of the missing mass. *Electronic Communications in Probability* **18** 1–7.
- [6] BIAU, G., BUNEA, F. and WEGKAMP, M. H. (2005). Functional classification in Hilbert spaces. *IEEE Trans. Inform. Theory* **51** 2163–2172. [MR2235289](#)
- [7] BIAU, G., CÉROU, F. and GUYADER, A. (2010). Rates of convergence of the functional  $k$ -nearest neighbor estimate. *IEEE Trans. Inform. Theory* **56** 2034–2040. [MR2654492](#)
- [8] BILLINGSLEY, P. (1968). *Convergence of probability measures*, first ed. *Wiley Series in Probability and Statistics: Probability and Statistics*. John Wiley & Sons Inc., New York. A Wiley-Interscience Publication. [MR1700749 \(2000e:60008\)](#)
- [9] BOIMAN, O., SHECHTMAN, E. and IRANI, M. (2008). In defense of Nearest-Neighbor based image classification. In *CVPR*.
- [10] CÉROU, F. and GUYADER, A. (2006). Nearest neighbor classification in infinite dimension. *ESAIM: Probability and Statistics* **10** 340–355.
- [11] CHAUDHURI, K. and DASGUPTA, S. (2014). Rates of Convergence for Nearest Neighbor Classification. In *NIPS*.
- [12] CHEN, G. H. and SHAH, D. (2018). Explaining the Success of Nearest Neighbor Methods in Prediction. *Foundations and Trends® in Machine Learning* **10** 337–588.
- [13] 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. Proceedings of a meeting held 6-9 December 2010, Vancouver, British Columbia, Canada*. 406–414.
- [14] COVER, T. M. and HART, P. E. (1967). Nearest neighbor pattern classification. *IEEE Transactions on Information Theory* **13** 21–27.

- [15] DAVID, O., MORAN, S. and YEHUDAYOFF, A. (2016). Supervised learning through the lens of compression. In *Advances in Neural Information Processing Systems 29* 2784–2792.
- [16] DEVROYE, L., GYÖRFI, L. and LUGOSI, G. (1996). *A probabilistic theory of pattern recognition*. Springer-Verlag New York, Inc.
- [17] DUDLEY, R. M. (1999). *Uniform Central Limit Theorems*. *Cambridge Studies in Advanced Mathematics*. Cambridge University Press.
- [18] FIX, E. and HODGES, J. J. L. (1989). Discriminatory Analysis. Nonparametric Discrimination: Consistency Properties. *International Statistical Review / Revue Internationale de Statistique* **57** pp. 238-247.
- [19] FORZANI, L., FRAIMAN, R. and LLOP, P. (2012). Consistent nonparametric regression for functional data under the Stone–Besicovitch conditions. *IEEE Transactions on Information Theory* **58** 6697–6708.
- [20] FREMLIN, D. H. (2000). *Measure theory* **1-5**. Torres Fremlin.
- [21] GADAT, S., KLEIN, T. and MARTEAU, C. (2016). Classification in general finite dimensional spaces with the  $k$ -nearest neighbor rule. *Ann. Statist.* **44** 982–1009.
- [22] GITIK, M. and SHELAH, S. (1989). Forcings with ideals and simple forcing notions. *Israel Journal of Mathematics* **68** 129–160.
- [23] GOTTLIEB, L., KONTOROVICH, A. and KRAUTHGAMER, R. (2014). Efficient Classification for Metric Data (extended abstract COLT 2010). *IEEE Transactions on Information Theory* **60** 5750–5759.
- [24] GOTTLIEB, L.-A., KONTOROVICH, A. and NISNEVITCH, P. (2017). Nearly optimal classification for semimetrics (extended abstract AISTATS 2016). *Journal of Machine Learning Research* **18** 1-22.
- [25] GOTTLIEB, L., KONTOROVICH, A. and NISNEVITCH, P. (2018). Near-Optimal Sample Compression for Nearest Neighbors (extended abstract: NIPS 2014). *IEEE Trans. Information Theory* **64** 4120–4128.
- [26] GRAEPEL, T., HERBRICH, R. and SHAWE-TAYLOR, J. (2005). PAC-Bayesian compression bounds on the prediction error of learning algorithms for classification. *Machine Learning* **59** 55–76.
- [27] GYÖRFI, L., KOHLER, M., ZAK, A. K. and WALK, H. (2002). *A Distribution-Free Theory of Nonparametric Regression*. Springer-Verlag New York, Inc.
- [28] HALL, P. and KANG, K.-H. (2005). Bandwidth choice for nonparametric classification. *Ann. Statist.* **33** 284–306.
- [29] HALL, P., PARK, B. U. and SAMWORTH, R. J. (2008). Choice of Neighbor Order in Nearest-Neighbor Classification. *The Annals of Statistics* **36** 2135–2152.
- [30] HANNEKE, S. (2017). Learning Whenever Learning is Possible: Universal Learning under General Stochastic Processes. *CoRR* [abs/1706.01418](#).
- [31] HANNEKE, S. and KONTOROVICH, A. (2019). A Sharp Lower Bound for Agnostic Learning with Sample Compression Schemes. In *Proceedings of the 30th International Conference on Algorithmic Learning Theory*.
- [32] JECH, T. (2003). *Set theory*. *Springer Monographs in Mathematics*. Springer-Verlag, Berlin The third millennium edition, revised and expanded. [MR1940513](#)
- [33] KECHRIS, A. S. (1995). *Classical Descriptive Set Theory*. Springer-Verlag New York.
- [34] KOLMOGOROV, A. N. and FOMIN, S. V. (1970). *Introductory Real Analysis*. Prentice Hall.
- [35] KONTOROVICH, A., SABATO, S. and URNER, R. (2017). Active Nearest-Neighbor Learning in Metric Spaces (extended abstract: NIPS 2016). *Journal of Machine Learning Research* **18** 195:1–195:38.
- [36] KONTOROVICH, A., SABATO, S. and WEISS, R. (2017). Nearest-Neighbor Sample Compression: Efficiency, Consistency, Infinite Dimensions. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA* 1572–1582.
- [37] KONTOROVICH, A., SABATO, S. and WEISS, R. (2017). Nearest-Neighbor Sample Compression: Efficiency, Consistency, Infinite Dimensions. *CoRR* [abs/1705.08184](#).

- [38] KONTOROVICH, A. and WEISS, R. (2014). A Bayes consistent 1-NN classifier. In *Artificial Intelligence and Statistics (AISTATS 2015)*.
- [39] KONTOROVICH, A. and WEISS, R. (2014). Maximum margin multiclass nearest neighbors. In *International Conference on Machine Learning (ICML 2014)*.
- [40] KPOTUFE, S. (2009). Fast, smooth and adaptive regression in metric spaces. In *Advances in Neural Information Processing Systems 22*.
- [41] KRAUTHGAMER, R. and LEE, J. R. (2004). Navigating nets: Simple algorithms for proximity search. In *15th Annual ACM-SIAM Symposium on Discrete Algorithms* 791–801.
- [42] KULKARNI, S. R. and POSNER, S. E. (1995). Rates of convergence of nearest neighbor estimation under arbitrary sampling. *IEEE Trans. Inform. Theory* **41** 1028–1039. [MR1366756](#)
- [43] LITTLESTONE, N. and WARMUTH, M. K. (1986). Relating Data Compression and Learnability. unpublished.
- [44] PÁLMASSON, H., JÓNSSON, B. P., AMSALEG, L., SCHEDL, M. and KNEES, P. (2017). On Competitiveness of Nearest-Neighbor-Based Music Classification: A Methodological Critique. In *Similarity Search and Applications* (C. BEECKS, F. BORUTTA, P. KRÖGER and T. SEIDL, eds.) 275–283. Springer International Publishing, Cham.
- [45] PAPADIMITRIOU, C. H. and STEIGLITZ, K. (1998). *Combinatorial optimization : algorithms and complexity*. Prentice Hall.
- [46] PELILLO, M. (2014). Alhazen and the nearest neighbor rule. *Pattern Recognition Letters* **38** 34–37.
- [47] PESTOV, V. (2000). On the geometry of similarity search: dimensionality curse and concentration of measure. *Inform. Process. Lett.* **73** 47–51. [MR1741506](#)
- [48] PESTOV, V. (2013). Is the  $k$ -NN classifier in high dimensions affected by the curse of dimensionality? *Comput. Math. Appl.* **65** 1427–1437. [MR3061714](#)
- [49] PREISS, D. (1979). Invalid Vitali theorems. *Abstracts. 7th Winter School on Abstract Analysis* 58–60.
- [50] PREISS, D. (1981). Gaussian measures and the density theorem. *Comment. Math. Univ. Carolin.* **22** 181–193. [MR609946](#)
- [51] PSALTIS, D., SNAPP, R. R. and VENKATESH, S. S. (1994). On the finite sample performance of the nearest neighbor classifier. *IEEE Transactions on Information Theory* **40** 820–837.
- [52] SAMWORTH, R. J. (2012). Optimal weighted nearest neighbour classifiers. *Ann. Statist.* **40** 2733–2763.
- [53] SCHERVISH, M. J. (1995). *Theory of statistics. Springer Series in Statistics*. Springer-Verlag, New York. [MR1354146](#)
- [54] SHALEV-SHWARTZ, S. and BEN-DAVID, S. (2014). *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press.
- [55] SNAPP, R. R. and VENKATESH, S. S. (1998). Asymptotic expansions of the  $k$  nearest neighbor risk. *Ann. Statist.* **26** 850–878. [MR1635410](#) ([99h:62062](#))
- [56] STONE, C. J. (1977). Consistent Nonparametric Regression. *The Annals of Statistics* **5** 595–620.
- [57] TIŠER, J. (2003). Vitali covering theorem in Hilbert space. *Trans. Amer. Math. Soc.* **355** 3277–3289. [MR1974687](#)
- [58] ULAM, S. M. (1930). *Zur Masstheorie in der allgemeinen Mengenlehre*. Uniwersytet, seminarjum matematyczne.
- [59] WEINBERGER, K. Q. and SAUL, L. K. (2009). Distance Metric Learning for Large Margin Nearest Neighbor Classification. *Journal of Machine Learning Research* **10** 207–244.
- [60] ZHAO, L. C. (1987). Exponential bounds of mean error for the nearest neighbor estimates of regression functions. *J. Multivariate Anal.* **21** 168–178. [MR877849](#) ([88e:62110](#))