

# Proper Learning, Helly Number, and an Optimal SVM Bound

**Olivier Bousquet**

*Google Research, Brain Team*

OBOUSQUET@GOOGLE.COM

**Steve Hanneke**

*Toyota Technological Institute at Chicago*

STEVE.HANNEKE@GMAIL.COM

**Shay Moran**

*Google Research, Brain Team*

SHAYMORAN@GOOGLE.COM

**Nikita Zhivotovskiy**

*Google Research, Brain Team*

ZHIVOTOVSKIY@GOOGLE.COM

**Editors:** Jacob Abernethy and Shivani Agarwal

## Abstract

The classical PAC sample complexity bounds are stated for any Empirical Risk Minimizer (ERM) and contain an extra multiplicative logarithmic factor  $\log \frac{1}{\varepsilon}$  which is known to be necessary for ERM in general. It has been recently shown by [Hanneke \(2016a\)](#) that the optimal sample complexity of PAC learning for any VC class  $\mathcal{C}$  does not include this log factor and is achieved by a particular *improper learning algorithm*, which outputs a specific majority-vote of hypotheses in  $\mathcal{C}$ . This leaves the question of when this bound can be achieved by *proper learning algorithms*, which are restricted to always output a hypothesis from  $\mathcal{C}$ .

In this paper we aim to characterize the classes for which the optimal sample complexity can be achieved by a proper learning algorithm. We identify that these classes can be characterized by the *dual Helly number*, which is a combinatorial parameter that arises in discrete geometry and abstract convexity. In particular, under general conditions on  $\mathcal{C}$ , we show that the dual Helly number is bounded if and only if there is a proper learner that obtains the optimal dependence on  $\varepsilon$ .

As further implications of our techniques we resolve a long-standing open problem posed by [Vapnik and Chervonenkis \(1974\)](#) on the performance of the *Support Vector Machine* in  $\mathbb{R}^n$  by proving that the sample complexity of SVM in the realizable case is

$$\Theta\left(\frac{n}{\varepsilon} + \frac{1}{\varepsilon} \log \frac{1}{\delta}\right).$$

This gives the first optimal PAC bound for Halfspaces in  $\mathbb{R}^n$  achieved by a proper learning algorithm, and moreover is computationally efficient.

**Keywords:** Statistical Learning Theory, PAC Learning, Sample Complexity, Proper Learning, SVM.

## 1. Introduction

In the literature on the theory of PAC learning, there has been much work discussing the important distinction between *proper vs improper* learning algorithms, where a proper learner is required to output a hypothesis from the concept class being learned, while an improper learner may output any classifier, not necessarily in the class. Most of this literature has focused on the *computational* separations between proper and improper learning (see e.g., [Kearns and Vazirani, 1994](#)). However,

it is also interesting to consider the effect on *sample complexity* of proper vs improper learning. While the optimal sample complexity of PAC learning was recently resolved by Hanneke (2016a), the proposed learning algorithm is *improper*: constructing its classifier based on a majority vote of well-chosen classifiers from the concept class. Furthermore, it follows from arguments analogous to the work of Daniely and Shalev-Shwartz (2014) that the optimal sample complexity of PAC learning is sometimes *not achievable* by proper learners (see also our Theorem 11). While the question of characterizing the best sample complexity achievable by proper learners has been resolved for several special-case concept classes (e.g., Auer and Ortner, 2007; Darnstädt, 2015; Hanneke, 2016b), the best known *general* results on the sample complexity of proper learning in the prior literature are still only the results that hold for *all* empirical risk minimization (ERM) algorithms (Vapnik and Chervonenkis, 1974; Blumer, Ehrenfeucht, Haussler, and Warmuth, 1989; Hanneke, 2016b; Zhivotovskiy and Hanneke, 2018). However, it is well known that there are many classes where specific proper learners can achieve better sample complexities (by a log factor) than the worst ERM learner (Auer and Ortner, 2007). Thus, it is important to go beyond the analysis of general ERM learners if we are to understand the best sample complexity achievable by proper learners in general.

In the present work, we aim to provide such a general analysis of the sample complexity of proper learning, applicable to *every* concept class, by identifying the relevant combinatorial complexity measures of the class. We specifically find that a quantity called the *dual Helly number* (previously proposed by Kane, Livni, Moran, and Yehudayoff, 2019 under the name *coVC dimension*) is of critical importance. In particular, when the dual Helly number is finite, the logarithmic factor in the well-known sample complexity bounds for ERM (Vapnik and Chervonenkis, 1974) may be replaced by a bounded quantity. The proper learner achieving this bound is a variant of the optimal PAC learner of Hanneke (2016a), but modified in several steps so that it remains proper.

As a further implication of the techniques we develop, we find that in the case of learning Halfspaces in  $\mathbb{R}^n$ , the well-known *support vector machine* (SVM) learning algorithm achieves the optimal sample complexity  $\Theta\left(\frac{n}{\epsilon} + \frac{1}{\epsilon} \log \frac{1}{\delta}\right)$ . This resolves a question that appeared in the seminal work of Vapnik and Chervonenkis (1974). Moreover, this also provides the first proof that Halfspaces are properly learnable with the optimal sample complexity: that is, sample complexity of the form  $\frac{n}{\epsilon} + \frac{1}{\epsilon} \log \frac{1}{\delta}$ . As a further implication, we find that Maximum classes of any given VC dimension  $d$  are also properly learnable with optimal sample complexity  $\Theta\left(\frac{d}{\epsilon} + \frac{1}{\epsilon} \log \frac{1}{\delta}\right)$ .

The known results on sample complexity are summarized in Figure 1 along with a (rough) statement of our new results.

**Notation.** To begin the formal discussion, we introduce some basic notation. Fix a space  $\mathcal{X}$  equipped with a  $\sigma$ -algebra specifying the measurable subsets. Let  $\mathcal{Y} = \{-1, 1\}$  denote the *label space*. A *classifier* is any measurable function  $h : \mathcal{X} \rightarrow \mathcal{Y}$ , and a *concept class* is any set  $\mathbb{C}$  of classifiers. To focus on nontrivial cases here, we will always suppose  $|\mathbb{C}| \geq 3$ . A *learning algorithm*  $\mathcal{A}$  maps any sequence (data set)  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  in  $\mathcal{X} \times \mathcal{Y}$ , of any length  $n$ , to a classifier  $\hat{h}_n$ ; the map  $\mathcal{A}$  may include randomization. A learning algorithm  $\mathcal{A}$  is called *proper* (for  $\mathbb{C}$ ) if  $\hat{h}_n$  is always an element of  $\mathbb{C}$ , for all possible data sets. Otherwise  $\mathcal{A}$  is called *improper*.

In the PAC learning problem, there is a *data distribution*  $\mathcal{P}$  (a probability measure on  $\mathcal{X}$ ), and a *target concept*  $f^* \in \mathbb{C}$ . For any classifier  $h$ , define  $\text{er}_{\mathcal{P}}(h; f^*) = \mathcal{P}(x : h(x) \neq f^*(x))$ . When  $\mathcal{P}$  and  $f^*$  are clear from the context, we simply write  $\text{er}(h)$ . In contexts where  $\mathcal{P}$  is specified, we let  $X_1, X_2, \dots$  denote an i.i.d. sequence of  $\mathcal{P}$ -distributed random variables. Define  $a \wedge b = \min\{a, b\}$  for  $a, b \in \mathbb{R}$ . Generally, for any sequence  $x_1, x_2, \dots$  and any  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , we use the notation

Bounds on the sample complexity of PAC learning		
Improper Learning	$\Theta\left(\frac{d}{\varepsilon} + \frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$	Hanneke, 2016a Ehrenfeucht et al., 1989
Any ERM	$O\left(\frac{d}{\varepsilon} \log\left(\frac{1}{\varepsilon} \wedge \frac{s}{d}\right) + \frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$ $\Omega\left(\frac{d}{\varepsilon} + \frac{1}{\varepsilon} \log\left(\frac{1}{\varepsilon} \wedge s\right) + \frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$	Hanneke, 2016b Vapnik and Chervonenkis, 1974
Proper Learning	$O\left(\frac{dk^2}{\varepsilon} \log(k) + \frac{k^2}{\varepsilon} \log \frac{1}{\delta}\right)$ $\Omega\left(\frac{d}{\varepsilon} + \frac{1}{\varepsilon} \log(k) + \frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$	New results in this work.
SVM / Halfspaces in $\mathbb{R}^n$	$\Theta\left(\frac{n}{\varepsilon} + \frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$	New result in this work.
Maximum Class (Proper)	$\Theta\left(\frac{d}{\varepsilon} + \frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$	New result in this work.

Figure 1: Summary of results on the sample complexity of  $(\varepsilon, \delta)$ -PAC learning, along with our new results.  $d$  denotes the *VC dimension* (Vapnik and Chervonenkis, 1971),  $s$  the *star number* (Hanneke and Yang, 2015), and  $k$  the *dual Helly number* (Kane, Livni, Moran, and Yehudayoff, 2019) discussed in this article. Specific definitions, conditions, and ranges of parameters for which the results hold are discussed below.

$1 : n = \{1, \dots, n\}$ ,  $x_{1:n} = \{x_1, \dots, x_n\}$ , and  $(x_{1:n}, f(x_{1:n})) = \{(x_1, f(x_1)), \dots, (x_n, f(x_n))\}$ . The *sample complexity*, the central quantity of study in this work, is defined as follows.

**Definition 1** For any  $\varepsilon, \delta \in (0, 1)$ , the sample complexity of  $(\varepsilon, \delta)$ -PAC learning, denoted  $\mathcal{M}(\varepsilon, \delta)$ , is defined as the smallest  $n \in \mathbb{N}$  for which there exists a learning algorithm  $\mathcal{A}$  such that, for every data distribution  $\mathcal{P}$  and every  $f^* \in \mathbb{C}$ , the (random) classifier  $\hat{h}_n = \mathcal{A}((X_{1:n}, f^*(X_{1:n})))$  satisfies

$$\Pr\left(\text{er}(\hat{h}_n) \leq \varepsilon\right) \geq 1 - \delta.$$

The sample complexity of  $(\varepsilon, \delta)$ -PAC proper learning, denoted by  $\mathcal{M}_{\text{prop}}(\varepsilon, \delta)$ , is defined identically, except that the learning algorithm  $\mathcal{A}$  is required to be proper: it always outputs an element of  $\mathbb{C}$ .

A fundamental quantity in characterizing the sample complexity is the *VC dimension* (Vapnik and Chervonenkis, 1971). We say  $\mathbb{C}$  *shatters* a sequence of points  $x_{1:n} \in \mathcal{X}^n$  if  $\forall y_{1:n} \in \mathcal{Y}^n, \exists h \in \mathbb{C}$  with  $h(x_{1:n}) = y_{1:n}$ . The VC dimension of  $\mathbb{C}$ , denoted by  $d$ , is the largest  $n \in \mathbb{N}$  for which there exists a sequence  $x_{1:n}$  shattered by  $\mathbb{C}$ ; otherwise if no such largest  $n$  exists, define  $d = \infty$ .

The sample complexity of (unrestricted) PAC learning was recently proven by Hanneke (2016a) to satisfy  $\mathcal{M}(\varepsilon, \delta) = \Theta\left(\frac{d}{\varepsilon} + \frac{1}{\varepsilon} \log\left(\frac{1}{\delta}\right)\right)$ , resolving a gap between the previously-known lower bound of this form (from Vapnik and Chervonenkis, 1974; Ehrenfeucht, Haussler, Kearns, and Valiant, 1989) and previous suboptimal upper bounds (Vapnik and Chervonenkis, 1974; Blumer, Ehrenfeucht, Haussler, and Warmuth, 1989; Haussler, Littlestone, and Warmuth, 1994; Simon, 2015). However, the optimal learning algorithm proposed by Hanneke (2016a) is *improper*.

Most of the work on the sample complexity of proper learning is based on the fact (due to Vapnik and Chervonenkis, 1974, Blumer, Ehrenfeucht, Haussler, and Warmuth, 1989) that any learning algorithm that outputs any  $\hat{h} \in \mathbb{C}$  making no mistakes on the training data (called *empirical risk*

*minimization*, or ERM) is guaranteed to achieve a sample complexity  $O(\frac{d}{\varepsilon} \log(\frac{1}{\varepsilon}) + \frac{1}{\varepsilon} \log(\frac{1}{\delta}))$ . This bound has been refined in some special cases (Hanneke, 2016b; Zhivotovskiy and Hanneke, 2018), but it is known that it cannot generally be improved while still holding for all ERM learners (Blumer and Littlestone, 1989; Haussler, Littlestone, and Warmuth, 1994; Auer and Ortner, 2007). On the other hand, for some special types of concept classes, it was observed that  $\mathcal{M}_{\text{prop}}(\varepsilon, \delta) = \Theta(\frac{d}{\varepsilon} + \frac{1}{\varepsilon} \log(\frac{1}{\delta}))$ : that is, proper learning is sometimes optimal. For instance, this was shown for any *intersection-closed* concept class, where the optimal sample complexity is achieved by a proper learner known as the *Closure* algorithm (Auer and Ortner, 2007; Darnstädt, 2015; Hanneke, 2016b). For the class of Halfspaces on  $\mathbb{R}^n$ , Vapnik and Chervonenkis (1974) found that the support vector machine (SVM) classifier (which is a proper learner) achieves the optimal dependence on  $d$  and  $\varepsilon$  for obtaining expected error  $\varepsilon$ , and they essentially asked the question of whether the exact optimal form  $\Theta(\frac{d}{\varepsilon} + \frac{1}{\varepsilon} \log(\frac{1}{\delta}))$  for PAC learning is achieved by SVM. This question has remained open since then, with a number of works investigating the question (e.g., Blumer and Littlestone, 1989; Balcan and Long, 2013; Zhivotovskiy, 2017; Hanneke and Kontorovich, 2019; Long and Long, 2020). We answer this question affirmatively, finding that indeed the SVM classifier achieves the optimal sample complexity for learning Halfspaces.

Related results are known for the multi-class setting (where we may have  $|\mathcal{Y}| > 2$ , or even infinite  $\mathcal{Y}$ ). In this case, Daniely, Sabato, Ben-David, and Shalev-Shwartz (2015) showed that different ERMs may have strikingly different sample complexities, and Daniely and Shalev-Shwartz (2014) showed that there exist classes (with  $|\mathcal{Y}| = \infty$ ) that are learnable but *not* properly learnable (i.e.,  $\mathcal{M}(\varepsilon, \delta) < \infty$  but  $\mathcal{M}_{\text{prop}}(\varepsilon, \delta) = \infty$ ). A similar separation between proper and improper learnability was also recently shown by Montasser, Hanneke, and Srebro (2019) for the problem of learning with adversarial robustness guarantees. Of course, these kinds of striking separations cannot happen in the binary classification setting ( $|\mathcal{Y}| = 2$ ) studied here, since the aforementioned result of Vapnik and Chervonenkis (1974) shows that ERM learners obtain sample complexities that are at most suboptimal by a factor  $O(\log \frac{1}{\varepsilon})$ . However, the argument used in the proofs of Daniely and Shalev-Shwartz (2014) and Montasser, Hanneke, and Srebro (2019) can be adapted to show that this logarithmic factor is sometimes *necessary*: that is, that there exist classes of any given VC dimension  $d$  for which  $\mathcal{M}_{\text{prop}}(\varepsilon, \delta) = \Omega(\frac{d}{\varepsilon} \log \frac{1}{\varepsilon} + \frac{1}{\varepsilon} \log \frac{1}{\delta})$ .<sup>1</sup>

One of the aims of this paper is to connect these scattered observations and explain the properties of the class  $\mathbb{C}$  responsible for the optimality or sub-optimality of proper PAC learners.

### Structure of the paper and main contributions

- Section 2 contains the definition of the dual Helly number and two variants of it: the *hollow star* number and the *projection number*. We show that these three parameters coincide whenever they are finite, and present some (infinite) classes where their values can differ.
- Section 3 contains an upper bound for proper learning classes with bounded dual Helly number. Specifically, we show that if the *projection number* is bounded, there is a particular ERM having a better sample complexity than arbitrary ERM. We also provide a proper learning algorithm achieving an even more-improved sample complexity, but this proper algorithm is not necessarily an ERM.
- Section 4 contains a lower bound on the sample complexity of proper algorithms when the hollow star number is large.

---

1. This result is unpublished, and essentially folklore, discovered independently by several different people familiar with the argument of Daniely and Shalev-Shwartz (2014).

- Section 5 presents a new upper bound for *stable compression schemes*: compression schemes whose choice of compression set is unaffected by removing points not in the compression set. In particular, as SVM can be viewed as a stable compression scheme, our result implies that SVM requires only  $O\left(\frac{n}{\varepsilon} + \frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$  examples in order to  $(\varepsilon, \delta)$ -PAC learn  $n$ -dimensional halfspaces. This resolves a long-standing open problem from [Vapnik and Chervonenkis \(1974\)](#). As a further implication of our general result for stable compression schemes, we also find that all Maximum classes are properly learnable with optimal sample complexity.

## 2. The dual Helly number

Helly’s Theorem is a fundamental result about convex sets ([Helly, 1923](#)). It asserts that for any finite sequence of convex sets  $C_1 \dots C_m \subseteq \mathbb{R}^n$  such that  $\bigcap_{i=1}^m C_i = \emptyset$  there is a subsequence  $C_{i_1} \dots C_{i_k}$ , with  $k \leq n + 1$  such that  $\bigcap_{j=1}^k C_{i_j} = \emptyset$ . This notion has been studied more abstractly in various settings (see, e.g., [Levi, 1951](#); [Danzer, Grünbaum, and Klee, 1963](#)); it is defined in an abstract manner as follows: let  $\mathcal{F}$  be a family of subsets over a domain  $\mathcal{X}$ . The *Helly number* of  $\mathcal{F}$  is the minimum integer  $k$  such that whenever  $\mathcal{C} \subseteq \mathcal{F}$  is a collection of sets whose intersection is empty then there is a subset  $\mathcal{C}' \subseteq \mathcal{C}$  of size at most  $k$  whose intersection is empty. That is, the empty intersection of the entire collection  $\mathcal{C}$  is witnessed by a subset of size at most  $k$ .

We adapt the Helly number (in a dual form) to our context, obtaining a parameter<sup>2</sup> critical to the proper sample complexity. For any  $S \subseteq \mathcal{X} \times \mathcal{Y}$ , define  $\mathbb{C}[S] = \{h \in \mathbb{C} : \forall (x, y) \in S, h(x) = y\}$ .

**Definition 2 (The dual Helly number)** *Define the dual Helly number of  $\mathbb{C}$ , denoted by  $k_w$ , as the smallest integer  $k$  such that, for any  $S \subseteq \mathcal{X} \times \mathcal{Y}$  such that  $\mathbb{C}[S] = \emptyset$ , there is a set  $W \subseteq S$  with  $|W| \leq k$  such that  $\mathbb{C}[W] = \emptyset$ . That is, for any unrealizable set of examples, there is an unrealizable subset of size at most  $k$ . If no such  $k$  exists, we define  $k_w = \infty$ .*

Observe that the dual Helly number is precisely the Helly number of the following family  $\mathcal{F}$ : for each  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  let  $\mathbb{C}_{(x,y)} := \{h \in \mathbb{C} : h(x) = y\}$ , and let  $\mathcal{F} := \{\mathbb{C}_{(x,y)} : (x, y) \in \mathcal{X} \times \mathcal{Y}\}$ . This definition is also related to the notion of a *teaching set* ([Goldman and Kearns, 1995](#)): recall that  $W \subseteq X$  is a teaching set for  $h$  with respect to  $\mathbb{C}$  if there exists no  $h' \in \mathbb{C} \setminus \{h\}$  which agrees with  $h$  on  $W$ . In particular, observe that any  $h \notin \mathbb{C}$  has a teaching set with respect to  $\mathbb{C}$  of size at most  $k_w$ .

We proceed with the second definition. We will need the following notation: two sequences (or samples)  $(x_1, y_1), \dots, (x_k, y_k) \in \mathcal{X} \times \mathcal{Y}$  and  $(x'_1, y'_1), \dots, (x'_k, y'_k) \in \mathcal{X} \times \mathcal{Y}$  are said to be *neighbors* if  $x_i = x'_i$  for all  $i \in 1 : k$  and there exists exactly one  $j \in 1 : k$  such that  $y_j \neq y'_j$ .

**Definition 3 (The hollow star number)** *Define the hollow star number of  $\mathbb{C}$ , denoted by  $k_o$ , as the largest integer  $k$  such that there is a sequence  $S = ((x_1, y_1), \dots, (x_k, y_k)) \in (\mathcal{X} \times \mathcal{Y})^k$  which is not realizable by  $\mathbb{C}$  (i.e.  $\mathbb{C}[S] = \emptyset$ ), however every sequence  $S'$  which is a neighbor of  $S$  is realizable by  $\mathbb{C}$  (i.e.,  $\mathbb{C}[S'] \neq \emptyset$ ). If no such largest  $k$  exists, define  $k_o = \infty$ .*

We refer to any unrealizable sequence  $S \in (\mathcal{X} \times \mathcal{Y})^*$ , such that every neighbor  $S'$  of  $S$  is realizable, as a *hollow star set*. Thus,  $k_o$  is the size of the largest finite hollow star set, or  $\infty$

2. We note that an equivalent parameter was introduced by the name “coVC dimension” by [Kane, Livni, Moran, and Yehudayoff \(2019\)](#), where it was used to characterize proper learning in a distributed setting.

if there exist hollow star sets of unbounded finite sizes. Equivalently, a hollow star set  $S = \{(x_1, y_1), \dots, (x_k, y_k)\}$  satisfies  $\mathbb{C}[S] = \emptyset$  and  $\forall i \leq k, \exists h_i \in \mathbb{C}$  s.t.  $\{j : h_i(x_j) \neq y_j\} = \{i\}$ .

The name *hollow star number* is chosen to stress the similarity with the *star number* defined by Hanneke and Yang (2015). The *star number*  $\mathfrak{s}$  is the maximum size of a realizable sequence  $S$  such that every sequence  $S'$  which is a neighbor of  $S$  is also realizable. Thus, the definitions of star number and hollow star number differ *only* in whether  $S$  is required to be realizable or unrealizable.

The star number  $\mathfrak{s}$  was defined by Hanneke and Yang (2015) to characterize the PAC sample complexity of active learning, and was shown by Hanneke (2016b) to also characterize the general rate of convergence of ERM (see Section 2.2 below). Interestingly, it can be shown (Hanneke and Yang, 2015) that the star number upper bounds the size of the *teaching set* of any classifier. It also has many other connections to various quantities arising in the learning theory literature (Hanneke and Yang, 2015; Hanneke, 2016b). From the definitions, we immediately have the simple inequality:

$$k_o - 1 \leq \mathfrak{s}.$$

However, as we discuss below, while the classes  $\mathbb{C}$  having  $\mathfrak{s} < \infty$  are very limited, classes with  $k_o < \infty$  are far more common (e.g., Halfspaces). Thus, this one small difference in the definition, requiring the star's *center*  $S$  to be unrealizable, significantly impacts the value of the quantity.

Our final definition is slightly more involved but it will play a key role in our upper bound. For a finite (multiset)  $\mathbb{C}' \subseteq \mathbb{C}$  let  $\text{Majority}(\mathbb{C}') : \mathcal{X} \rightarrow \{0, 1, ?\}$  denote the majority-vote classifier:

$$\text{Majority}(\mathbb{C}')(x) = \begin{cases} 0 & |\{c \in \mathbb{C}' : c(x) = 0\}| > \frac{|\mathbb{C}'|}{2}, \\ 1 & |\{c \in \mathbb{C}' : c(x) = 1\}| > \frac{|\mathbb{C}'|}{2}, \\ ? & \text{else.} \end{cases}$$

For  $\ell \geq 2$ , define the set  $\mathcal{X}_{\mathbb{C}', \ell} \subseteq \mathcal{X}$  of all the points  $x$  on which less than  $\frac{1}{\ell}$ -fraction of all classifiers in  $\mathbb{C}'$  disagree with the majority. That is, letting  $h_{\text{maj}} = \text{Majority}(\mathbb{C}')$ ,

$$\mathcal{X}_{\mathbb{C}', \ell} = \left\{ x \in \mathcal{X} : \sum_{h \in \mathbb{C}'} \mathbb{1}[h(x) \neq h_{\text{maj}}(x)] < \frac{|\mathbb{C}'|}{\ell} \right\}, \quad (1)$$

**Definition 4 (The projection number)** Define the projection number of  $\mathbb{C}$ , denoted by  $k_p$ , as the smallest integer  $k \geq 2$  such that, for any finite multiset  $\mathbb{C}' \subseteq \mathbb{C}$  there exists  $h \in \mathbb{C}$  that agrees with  $\text{Majority}(\mathbb{C}')$  on the entire set  $\mathcal{X}_{\mathbb{C}', k}$ . If no such integer  $k$  exists, define  $k_p = \infty$ .

This definition allows us to “project” the majority vote of any classifiers in  $\mathbb{C}$  to the class  $\mathbb{C}$ . Define

$$\text{Proj}_{\mathbb{C}}(\mathbb{C}') \text{ is any element in } \{h \in \mathbb{C} : h(x) = \text{Majority}(\mathbb{C}'), \text{ for all } x \in \mathcal{X}_{\mathbb{C}', k_p}\}.$$

The set used in this definition is always non-empty (by definition of  $k_p$ ) when  $k_p < \infty$ .

Following Kane, Livni, Moran, and Yehudayoff (2019), we say a class  $\mathbb{C}$  is “closed” if every  $S \subseteq \mathcal{X} \times \mathcal{Y}$  with  $\mathbb{C}[S] = \emptyset$  has a finite subset  $S' \subseteq S$  with  $\mathbb{C}[S'] = \emptyset$ . The following lemma connects these quantities. Its proof is included in Appendix A.

**Lemma 5**

- $k_o \leq k_p \leq k_w$ .
- If  $k_w < \infty$  or  $\mathbb{C}$  is closed, then  $k_o = k_p = k_w$ .

**Remark 6** Certainly if  $\mathcal{X}$  or  $\mathbb{C}$  is finite then  $k_w < \infty$ , so that  $k_o = k_p = k_w$ . However, in the general case, there are examples where each of these inequalities can be strict, due to one quantity being infinite and another finite. We discuss such examples in Section 2.1 below.

## 2.1. Some Examples

Let us now argue that many classes of interest have finite values for these complexity measures. Where appropriate, details of the examples are provided in Appendix A. We start with Halfspaces.

**Example 1** *Let  $\mathbb{C}$  be a class induced by halfspaces in  $\mathbb{R}^n$ . Then  $k_o = k_p = n + 2 = d + 1$ .*

Another simple example is the case of *intersection closed classes*. The class  $\mathbb{C}$  is intersection closed if the set  $\{\{x : h(x) = +1\} : h \in \mathbb{C}\}$  is closed under arbitrary intersections. Also, recall the notions of *Maximum* and *Extremal* classes (see e.g., Floyd and Warmuth, 1995; Lawrence, 1983; Bandelt, Chepoi, Dress, and Koolen, 2006; Moran and Warmuth, 2016). A class  $\mathbb{C}$  is a *Maximum class* if, for every integer  $m \geq d$  and every distinct  $x_1, \dots, x_m \in \mathcal{X}$ , we have  $|\{(h(x_1), \dots, h(x_m)) : h \in \mathbb{C}\}| = \sum_{i=0}^d \binom{m}{i}$ . A class  $\mathbb{C}$  is an *Extremal class* if, for every sequence  $x_1, \dots, x_m$  shattered by  $\mathbb{C}$ ,  $x_1, \dots, x_m$  is also shattered by a set of classifiers in  $\mathbb{C}$  that agree on all of  $\mathcal{X} \setminus \{x_1, \dots, x_m\}$ . It is known that every Maximum class is Extremal.

**Example 2** *If  $\mathbb{C}$  is intersection-closed or Extremal (e.g., any Maximum class), then  $k_o \leq d + 1$ .*

This also means any *closed* class that is intersection-closed or Extremal has  $k_o = k_p = k_w \leq d + 1$ .

As mentioned, all of the inequalities in Lemma 5 can be strict, due to the larger quantity being infinite while the smaller one is finite. We now discuss this issue, and how it relates to whether a class is closed (in terms of the definition of “closed” above, from Kane, Livni, Moran, and Yehudayoff, 2019). In particular, in all these examples, merely adding the limit cases into the class brings the complexity measures into agreement. We begin with a simple example where  $k_o = 2$  but  $k_p = \infty$ .

**Example 3** *Consider  $\mathcal{X} = \mathbb{N}$  and  $\mathbb{C} = \{2\mathbb{1}_{\{t\}} - 1 : t \in \mathcal{X}\}$  the class of singletons. The only finite hollow star sets are of size 2: namely, sets  $\{(x, 1), (x', 1)\}$  and  $\{(x, 1), (x, -1)\}$ . Thus,  $k_o = 2$ . However, for any finite set  $\mathbb{C}' \subset \mathbb{C}$  of size at least 3, Majority( $\mathbb{C}'$ ) is  $-1$  everywhere on  $\mathcal{X}$ , but any  $\ell < |\mathbb{C}'|$  has  $\mathcal{X}_{\mathbb{C}', \ell} = \mathcal{X}$ , so that we must have  $k_p > \ell$ ; thus,  $k_p = \infty$ .*

*Note that the constant function  $x \mapsto -1$  is a pointwise limit of functions in  $\mathbb{C}$ . By adding just this one extra function, the modified class  $\mathbb{C} \cup \{x \mapsto -1\}$  is closed and has  $k_w = k_p = k_o = 2$ .*

Next we describe a simple example where  $k_p = 2$  but  $k_w = \infty$ .

**Example 4** *Consider  $\mathcal{X} = [0, \infty)$  and  $\mathbb{C} = \{2\mathbb{1}_{[t, \infty)} - 1 : t \in [0, \infty)\}$  the class of threshold functions. This class is Maximum with  $d = 1$ . For any finite  $\mathbb{C}' \subset \mathbb{C}$ , the majority vote classifier is a median threshold from  $\mathbb{C}'$ , so there is always an  $h \in \mathbb{C}$  that agrees with Majority( $\mathbb{C}'$ ) on all of  $\mathcal{X}$ . Therefore,  $k_p = 2$ , the smallest possible value of  $k_p$ . On the other hand, the set  $S = \{(x, -1) : x \in \mathcal{X}\}$  is unrealizable by  $\mathbb{C}$ , but there is no finite set witnessing this fact, and therefore  $k_w = \infty$ .*

*In this case, there are an infinite number of functions that are pointwise limits of functions in  $\mathbb{C}$ : namely,  $x \mapsto -1$ , and every open threshold  $2\mathbb{1}_{(t, \infty)} - 1$ ,  $t \in \mathcal{X}$ . By adding these functions, the class  $\mathbb{C} \cup \{x \mapsto -1\} \cup \{2\mathbb{1}_{(t, \infty)} - 1 : t \in \mathcal{X}\}$  has  $k_w = k_p = k_o = 2$ .*

These distinctions can occur in more extreme forms as well, as the following example illustrates.

**Example 5** *Consider  $\mathcal{X} = \mathbb{R}$  and  $\mathbb{C} = \{2\mathbb{1}_{[a, b]} - 1 : a, b \in \mathbb{R}\}$  the class of bounded closed intervals. This class is intersection-closed (and Maximum), with  $d = 2$ , and has  $k_o = k_p = 3$ , but the set  $S = \{(x, 1) : x \in \mathcal{X}\}$  has no finite subset witnessing its non-realizability, so  $k_w = \infty$ .*

*In this case, to bring the complexity measures into agreement, we must increase  $\mathbb{C}$  to be the set of all intervals (including unbounded intervals, closed intervals, open intervals, and half-open half-closed intervals); this expanded class is closed and has  $k_w = k_p = k_o = 3$ .*

## 2.2. Star number and sample complexity of ERM

We finish this section by recalling the known relations between ERM and the star number. We use the following notation (also used in the results below). For  $x \geq 0$ , define  $\text{Log}(x) = \max\{\ln(x), 1\}$ .

**Definition 7 (The worst-case sample complexity of ERM)** *For any  $\varepsilon, \delta \in (0, 1)$ , the worst-case sample complexity of ERM, denoted by  $\mathcal{M}_{\text{ERM}}(\varepsilon, \delta)$ , is the smallest integer  $n$  such that for every possible data distribution  $\mathcal{P}$  and every  $f^* \in \mathbb{C}$ , for  $S = (X_{1:n}, f^*(X_{1:n}))$  with  $X_{1:n} \sim \mathcal{P}^n$ ,*

$$\Pr \left( \sup_{h \in \mathbb{C}[S]} \text{er}(h, f^*) \leq \varepsilon \right) \geq 1 - \delta.$$

The following bounds were shown by (Hanneke, 2016b):

$$\frac{1}{\varepsilon} \left( d + \text{Log} \left( \mathfrak{s} \wedge \frac{1}{\varepsilon} \right) + \text{Log} \left( \frac{1}{\delta} \right) \right) \lesssim \mathcal{M}_{\text{ERM}}(\varepsilon, \delta) \lesssim \frac{1}{\varepsilon} \left( d \text{Log} \left( \frac{\mathfrak{s}}{d} \wedge \frac{1}{\varepsilon} \right) + \text{Log} \left( \frac{1}{\delta} \right) \right).$$

In particular, having  $\mathfrak{s} < \infty$  is necessary and sufficient for the existence of a distribution-free bound on the error rates of all ERMs converging at a rate inversely proportional to the sample size. Our results below show that a much weaker assumption  $k_p < \infty$  implies an analog of this property but only for *some* proper learners instead of *all* ERMs. It is important to notice that the projection number is finite for many expressive VC classes. In contrast, the star number, while implying an upper bound for  $k_o$ , is infinite except for some relatively simple classes: e.g.,  $\mathfrak{s} = \infty$  for Halfspaces in  $\mathbb{R}^n$  if  $n \geq 2$ , and for many intersection-closed and maximum classes (Hanneke and Yang, 2015).

## 3. Upper bounds

Our upper bound will be established for the following algorithm, a modification of the optimal PAC learner of Hanneke (2016a). The main modifications compared to the original algorithm involve using  $k_p + 1$  recursive calls, rather than 3 calls, and using the projection operator (introduced above) to replace a majority vote classifier with an element of  $\mathbb{C}$  so that the algorithm is a proper learner.

$\mathbb{A}(S; T)$

1. If  $|S| < 4$ , Return  $\text{ERM}(S \cup T)$
2. Let  $S_0$  be the first  $\lceil |S|/2 \rceil$  points in  $S$
3. Let  $S_1, \dots, S_{k_p+1}$  be independent uniform subsamples of  $S \setminus S_0$  of size  $\lfloor |S|/4 \rfloor$
4. Let  $h_i = \mathbb{A}(S_0; T \cup S_i)$  for each  $i = 1, \dots, k_p + 1$
5. Return  $\hat{h} = \text{Proj}_{\mathbb{C}}(h_1, \dots, h_{k_p+1})$ .

To be precise, in Step 3,  $S_i$  is sampled without replacement. For this algorithm, we have the following theorem, representing one of the main results of this article.<sup>3</sup> We present its proof in Appendix B.

**Theorem 8** *For any class  $\mathbb{C}$ , the proper sample complexity (achieved by  $\mathbb{A}(S; \emptyset)$ ) satisfies*

$$\mathcal{M}_{\text{prop}}(\varepsilon, \delta) = O \left( \frac{k_p^2}{\varepsilon} \left( d \text{Log}(k_p) + \text{Log} \left( \frac{1}{\delta} \right) \right) \right).$$

3. We implicitly assume that  $\mathbb{C}$  satisfies conditions so that all of the relevant random variables in the analysis are measurable. This is always the case if  $\mathcal{X}$  is countable. For the uncountable case, we refer the reader to discussions by Blumer, Ehrenfeucht, Haussler, and Warmuth (1989); van der Vaart and Wellner (1996); van Handel (2013).

Interestingly, algorithm  $\mathbb{A}$  above is not necessarily an ERM: it is not always consistent with the training sample. However, it is also possible to define a sample-consistent variant of the algorithm, only losing a factor of  $k_p$  in the sample complexity bound. Note that the standard VC bounds for ERM apply to *any* ERM classifier, whereas the bound we establish for this algorithm only holds for this *specific* choice of ERM classifier, which is therefore sometimes significantly better than the worst ERM. We proceed with the definition and theorem, the proof of which is given in Appendix B.

$\mathbb{A}_{\text{ERM}}(S)$

1. If  $|S| < k_p + 1$ , Return  $\text{ERM}(S)$
2. Split  $S$  into  $k_p + 1$  disjoint subsets  $S_i$  of size at least  $\lfloor |S|/(k_p + 1) \rfloor$
3. Set  $h_i = \mathbb{A}_{\text{ERM}}(\bigcup_{i' \neq i} S_{i'})$  for  $i = 1, \dots, k_p + 1$ .
4. Return  $\hat{h} = \text{Proj}_{\mathbb{C}}(h_1, \dots, h_{k_p+1})$ .

To be precise, in Step 2,  $S$  is split into the subsets  $S_i$  based purely on the indices: for instance, take  $S_1$  as the first  $\lfloor |S|/(k_p + 1) \rfloor$  data points in the sequence,  $S_2$  as the next  $\lfloor |S|/(k_p + 1) \rfloor$  points in the sequence, and so on, with  $S_{k_p+1}$  as the last  $|S| - k_p \lfloor |S|/(k_p + 1) \rfloor$  data points in the sequence.

**Theorem 9** *For any class  $\mathbb{C}$ , the sample complexity  $\mathcal{M}_{\mathbb{A}_{\text{ERM}}}$  of  $\mathbb{A}_{\text{ERM}}(S)$  satisfies*

$$\mathcal{M}_{\mathbb{A}_{\text{ERM}}}(\varepsilon, \delta) = O\left(\frac{k_p^3}{\varepsilon} \left( d\text{Log}(k_p) + \text{Log}\left(\frac{1}{\delta}\right) \right)\right).$$

#### 4. Lower Bounds

For the purpose of a lower bound, we will use the hollow star number  $k_o$ . From Lemma 5, for many classes we have  $k_o = k_p = k_w$ , in which case this result indicates that the appearance of  $k_p$  in Theorem 8 is unavoidable (though there may be room to improve the specific dependence on  $k_p$ ).

**Theorem 10** *Every class with  $k_o < \infty$  has proper sample complexity*

$$\mathcal{M}_{\text{prop}}(\varepsilon, \delta) = \Omega\left(\frac{d}{\varepsilon} + \frac{1}{\varepsilon} \text{Log}\left(\frac{1}{\delta}\right) + \frac{1}{\varepsilon} \text{Log}(k_o) \mathbb{1}[\varepsilon \leq 1/k_o]\right).$$

Also, if  $k_o = \infty$  then

$$\mathcal{M}_{\text{prop}}(\varepsilon, \delta) \neq o\left(\frac{1}{\varepsilon} \text{Log}\left(\frac{1}{\varepsilon}\right)\right).$$

In other words, there is a sequence  $\varepsilon_i \rightarrow 0$  with  $\mathcal{M}_{\text{prop}}(\varepsilon_i, \delta) \geq \frac{c}{\varepsilon_i} \log\left(\frac{1}{\varepsilon_i}\right)$  for a constant  $c > 0$ .

The proof, which is inspired by the arguments in (Daniely and Shalev-Shwartz, 2014), is deferred to Appendix C. Let us only provide some intuition behind the proof. Assume, for simplicity, that we want to lower bound the sample complexity in a particular regime where  $\varepsilon = \frac{1}{2(k_o-1)}$ ,  $k_o \geq 2$ . Let  $S = \{(x_1, y_1), \dots, (x_{k_o}, y_{k_o})\}$  be a hollow star set, and  $\forall i \in 1 : k_o$  let  $h_i \in \mathbb{C}$  be such that  $\{j : h_i(x_j) \neq y_j\} = \{i\}$ . We set  $f^* = h_{i^*}$  for some  $i^* \in 1 : k_o$  and set  $\mathcal{P}(\{x_{i^*}\}) = 0$  and  $\mathcal{P}(\{x_i\}) = 2\varepsilon$  for  $i \neq i^*$ . Observe that the only way for the learner to output  $\hat{h} \in \mathbb{C}$  having  $\text{er}_{\mathcal{P}}(\hat{h}, f^*) \leq \varepsilon$  is for  $\hat{h}$  to agree with  $h_{i^*}$  on *all* of  $S$ . However, the learner will not be able to identify

the corresponding point  $x_{i^*}$  having zero mass before it observes  $x_i$  for every  $i \neq i^*$ . By the standard coupon collector argument (see Lemma 19 in Appendix C) due to the fact that there will be some copies in the training sample we will need the sample size of order  $\Omega(k_o \text{Log}(k_o)) = \Omega(\frac{1}{\varepsilon} \text{Log}(k_o))$ . The formal application of this idea (including extension to any  $\varepsilon \leq 1/k_o$ ) is given in Appendix C.

Our second lower bound provides stronger guarantees. However, it is less general. In what follows, given  $d$  and  $k_w$  we present a *particular* class  $\mathbb{C}$  and a space  $\mathcal{X}$  such that the desired lower bound holds. The proof of this result uses the same logic and the technical details are deferred to Appendix C.

**Theorem 11** *There is a numerical constant  $c > 0$  such that, for any value of  $d \geq 1$  and  $2 \leq k_w < \infty$ , there exists a space  $\mathcal{X}$  and a class  $\mathbb{C}$  with VC dimension  $d$  and dual Helly number  $k_w$ , for which the proper sample complexity satisfies*

$$\mathcal{M}_{\text{prop}}(\varepsilon, \delta) \geq \frac{c}{\varepsilon} \left( d \text{Log} \left( \frac{k_w}{d} \wedge \frac{1}{\varepsilon} \right) + \text{Log} \left( \frac{1}{\delta} \right) \right),$$

for every  $\varepsilon \in (0, 1/8)$  and  $\delta \in (0, 1/100)$ . Furthermore, for any  $d \geq 1$  there exists  $\mathcal{X}$  and a space  $\mathbb{C}$  with VC dimension  $d$  and hollow star number  $k_o = \infty$  and

$$\mathcal{M}_{\text{prop}}(\varepsilon, \delta) \geq \frac{c}{\varepsilon} \left( d \text{Log} \left( \frac{1}{\varepsilon} \right) + \text{Log} \left( \frac{1}{\delta} \right) \right).$$

## 5. Stable compression schemes and optimality of SVM for PAC learning Halfspaces

In this section we establish a new generalization bound for a special type of compression scheme, referred to as a *stable* compression scheme, which removes a log factor (which is known to not be removable for general compression schemes). We apply the result to obtain new tighter bounds for several quantities of interest in the learning theory literature.

As our main application, we apply this new bound to resolve a long-standing open question: namely, showing that the well-known *support vector machine* (SVM) learning algorithm for Halfspaces achieves the optimal sample complexity. This resolves a question posed by Vapnik and Chervonenkis (1974), which has received considerable attention in the literature (e.g., Blumer and Littlestone, 1989; Balcan and Long, 2013; Zhivotovskiy, 2017; Hanneke and Kontorovich, 2019; Long and Long, 2020). According to a note by A. Chervonenkis (see Chapter I in Vovk, Papadopoulos, and Gammernan, 2015) the in-expectation version of the risk bound for SVM had been proven in 1966, even before the renowned uniform law of large numbers was announced (Vapnik and Chervonenkis, 1968). However, obtaining a high-probability version of the bound, without introducing additional log factors, remained an open problem since then. The following theorem resolves this problem. Furthermore, this is also the first proof that the optimal sample complexity for Halfspaces is achievable by *some* proper learner.

**Theorem 12** *The PAC sample complexity of SVM in  $\mathbb{R}^n$  is*

$$\mathcal{M}_{\text{SVM}}(\varepsilon, \delta) = \Theta \left( \frac{n}{\varepsilon} + \frac{1}{\varepsilon} \log \frac{1}{\delta} \right).$$

The proof is presented below, after the abstract result that implies it. This bound improves the sample complexity bound  $\mathcal{M}_{\text{SVM}}(\varepsilon, \delta) = O \left( \frac{n}{\varepsilon} \log \frac{1}{\delta} \right)$  shown in (Zhivotovskiy, 2017) and the general bound  $\mathcal{M}_{\text{ERM}}(\varepsilon, \delta) = O \left( \frac{n}{\varepsilon} \log \frac{1}{\varepsilon} + \frac{1}{\varepsilon} \log \frac{1}{\delta} \right)$  holding for any ERM (Vapnik and Chervonenkis, 1974; Blumer, Ehrenfeucht, Haussler, and Warmuth, 1989).

### 5.1. Stable compression schemes

We present here our new generalization bound for stable compression schemes. Let us recall some standard definitions. A *sample compression scheme* consists of two functions: a *compression function*  $\kappa : (\mathcal{X} \times \mathcal{Y})^* \rightarrow (\mathcal{X} \times \mathcal{Y})^*$  and a *reconstruction function*  $\rho : (\mathcal{X} \times \mathcal{Y})^* \rightarrow \mathcal{Y}^{\mathcal{X}}$ . The following definition is due to [Littlestone and Warmuth \(1986\)](#).

**Definition 13 (Sample compression scheme)** *The pair of functions  $(\kappa, \rho)$  define a sample compression scheme of size  $\ell$  if for any  $m \in \mathbb{N}$  and  $h \in \mathbb{C}$  and any sample  $x_{1:m} \in \mathcal{X}^m$  it holds that  $\kappa((x_{1:m}, h(x_{1:m}))) \subseteq (x_{1:m}, h(x_{1:m}))$  and  $|\kappa((x_{1:m}, h(x_{1:m})))| \leq \ell$ , and the classifier  $\hat{h} = \rho(\kappa((x_{1:m}, h(x_{1:m}))))$  satisfies  $\hat{h}(x_{1:m}) = h(x_{1:m})$ : that is, it recovers  $h$ 's classifications on  $x_{1:m}$ .*

We work with the following natural definition taking its roots in [\(Vapnik and Chervonenkis, 1974\)](#).

**Definition 14 (Stable compression scheme)** *A sample compression scheme  $(\kappa, \rho)$  is called stable if for any  $m \in \mathbb{N}$ ,  $h \in \mathbb{C}$ ,  $x_{1:m} \in \mathcal{X}^m$  and any  $(x, h(x)) \in (x_{1:m}, h(x_{1:m})) \setminus \kappa((x_{1:m}, h(x_{1:m})))$ ,*

$$\kappa((x_{1:m}, h(x_{1:m})) \setminus (x, h(x))) = \kappa((x_{1:m}, h(x_{1:m}))).$$

This definition means that removing any  $(x, h(x))$  *not* belonging to the compression set, the compression set of the sub-sample remains the same. Finally, we say that the sample compression scheme  $(\kappa, \rho)$  is *proper* if the image of the reconstruction function  $\rho$  is contained in  $\mathbb{C}$ .

It is known that for any stable compression scheme  $(\kappa, \rho)$  of any size  $\ell$ , for any  $\mathcal{P}$  and  $f^* \in \mathbb{C}$ , for any  $m \in \mathbb{N}$ ,  $\mathbb{E}[\text{er}(\rho(\kappa((X_{1:m}, f^*(X_{1:m})))))] \leq \frac{\ell}{m+1}$ . This follows from the leave-one-out analysis of [Vapnik and Chervonenkis \(1974\)](#) (see also [Haussler, Littlestone, and Warmuth, 1994](#); [Zhivotovskiy, 2017](#)). The best known PAC generalization bound on  $\text{er}(\rho(\kappa(X_{1:m}, f^*(X_{1:m}))))$  (holding with probability  $1 - \delta$ ) valid for *any* sample compression scheme of a size  $\ell$  is due to [Littlestone and Warmuth \(1986\)](#) (see also [Floyd and Warmuth, 1995](#)):  $O(\frac{1}{m} (\ell \log(m) + \log(\frac{1}{\delta})))$ , where the  $\log(m)$  factor improves to  $\text{Log}(\frac{m}{\ell})$  if  $\rho$  is permutation-invariant. [Floyd and Warmuth \(1995\)](#) showed that there exist spaces  $\mathbb{C}$  and compression schemes for  $\mathbb{C}$  for which this log factor cannot be improved. However, in the special case of *stable* compression schemes, [Zhivotovskiy \(2017\)](#) established a bound  $O(\frac{\ell}{m} \log(\frac{1}{\delta}))$ , which is sometimes better.

As one of the main contributions of this work, the following result improves this PAC generalization bound by completely removing the log factor from the bound of [Littlestone and Warmuth \(1986\)](#). A simple proof of this result is provided in [Appendix D](#).

**Theorem 15** *Assume that  $\mathbb{C}$  has a stable sample compression scheme  $(\kappa, \rho)$  of size  $\ell$ . Then, for any  $\mathcal{P}$  and any  $f^* \in \mathbb{C}$ , for any integer  $m > 2\ell$ , given an i.i.d. sample  $S = (X_{1:m}, f^*(X_{1:m}))$  of size  $m$ , for any  $\delta \in (0, 1)$ , we have with probability at least  $1 - \delta$ ,*

$$\text{er}(\rho(\kappa(S))) < \frac{2}{m - 2\ell} \left( \ell \ln(4) + \ln\left(\frac{1}{\delta}\right) \right).$$

This result immediately yields the following proof of [Theorem 12](#) establishing optimality of SVM for PAC learning Halfspaces.

**Proof of [Theorem 12](#)** Since SVM may be expressed as a stable compression scheme of size  $n + 1$  (see [\(Long and Long, 2020\)](#) for a transparent proof of this fact, originally proven by [Vapnik and Chervonenkis 1974](#)), the upper bound is immediate from [Theorem 15](#). The lower bound follows from [\(Ehrenfeucht, Haussler, Kearns, and Valiant, 1989\)](#). ■

**Maximum classes.** As a second application of Theorem 15 to proper learning, consider any *maximum class*  $\mathbb{C}$  (Floyd and Warmuth, 1995; recall the definition from Section 2.1). Every maximum class  $\mathbb{C}$  is known to have a proper stable compression scheme of size  $d$ ; this follows from Theorem 5.1 and condition (R2) of Theorem 6.1 from the paper of Chalopin, Chepoi, Moran, and Warmuth (2019) (specifically, for their  $(\kappa, \rho)$ , every realizable data set contains exactly one compression set for a concept in  $\mathbb{C}$  consistent with the data, and hence removing any point not in that compression set cannot change the identity of this unique compression set). Hence, we have the following corollary of Theorem 15, establishing (for the first time) that maximum classes are properly learnable with sample complexity of the same order as the optimal PAC sample complexity.

**Corollary 16** For  $\mathbb{C}$  which is a Maximum class,

$$\mathcal{M}_{\text{prop}}(\varepsilon, \delta) = \Theta\left(\frac{d}{\varepsilon} + \frac{1}{\varepsilon} \log\left(\frac{1}{\delta}\right)\right).$$

**Remark 17** We note that Theorem 15 is also able to recover the optimal PAC sample complexity for the Closure algorithm for intersection-closed classes (Helmbold, Sloan, and Warmuth, 1990); this sample complexity result was already known via different arguments (Darnstädt, 2015; Hanneke, 2016b; Auer and Ortner, 2007), though Theorem 15 offers improved numerical constants.

**Remark 18** A further implication of Theorem 15 and Theorem 10 together is that any class  $\mathbb{C}$  with  $k_o = \infty$  does not have a proper stable compression scheme of bounded size.

## 6. Remaining Gaps and Open Questions

Although our upper and lower bounds give a description of proper PAC sample complexity in many cases, there are still situations not explained by our general bounds. Consider Example 3 (singletons  $x \mapsto 2\mathbb{1}_{\{x\}}(x) - 1$ ). Since  $k_o = 2$  but  $k_p = \infty$ , our upper bound (Theorem 8) does not match our lower bound (Theorem 10), and also does not match the optimal (improper) sample complexity  $\mathcal{M}(\varepsilon, \delta) = \Theta\left(\frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$ . However, it turns out there *does* exist a proper stable compression scheme of size 1 for this  $\mathbb{C}$ , and therefore Theorem 15 implies  $\mathcal{M}_{\text{prop}}(\varepsilon, \delta) = \Theta\left(\frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$ , matching the optimal (improper) sample complexity. Specifically, for any data set  $S$ , if  $\exists(x, 1) \in S$ , define  $\kappa(S) = \{(x, 1)\}$  and  $\rho(\{(x, 1)\}) = 2\mathbb{1}_{\{x\}} - 1$ ; if  $S$  is all negative examples, define  $\kappa(S) = \{(x, -1)\}$  for the *largest*  $x$  in  $S$ , and  $\rho(\{(x, -1)\}) = 2\mathbb{1}_{\{x+1\}} - 1$ . Alternatively, we may simply note that  $\mathbb{C}$  is a *Maximum class*, hence Corollary 16 applies. It is therefore important to ask whether  $\mathcal{M}_{\text{prop}}(\varepsilon, \delta)$  is perhaps *always* characterized by  $k_o$  (and  $d, \varepsilon, \delta$ ). It is also interesting to consider whether the optimal size of a proper stable compression scheme is also always characterized by  $k_o$  (and  $d$ ). Another interesting direction is sharpening the dependence on  $k_p$  in our upper bounds.

Finally, we remark that the result of Section 5 is related to the question of obtaining generalization bounds for learning algorithms that are stable with respect to small perturbations in the training sample. In the case of uniformly stable algorithms, recent results (Feldman and Vondrak, 2019; Bousquet, Klochkov, and Zhivotovskiy, 2019) provide sharp high-probability bounds, and the proofs are based on a sub-sampling argument: the learner is tested on carefully chosen parts of the training sample. Similarly to uniformly stable algorithms, stable compression schemes easily provide sharp in-expectation error bounds (Haussler, Littlestone, and Warmuth, 1994). The challenging part, already pointed out by Vapnik and Chervonenkis (1974), is to prove that these algorithms admit sharp high-probability bounds. Our result, based on related arguments, is the first to prove that an optimal high-probability error bound holds for stable compression schemes, including SVM.

## References

- P. Auer and R. Ortner. A new PAC bound for intersection-closed concept classes. *Machine Learning*, 66(2-3):151–163, 2007.
- M.-F. Balcan and P. M. Long. Active and passive learning of linear separators under log-concave distributions. In *Proceedings of the 26<sup>th</sup> Conference on Learning Theory*, 2013.
- H.-J. Bandelt, V. Chepoi, A. W. M. Dress, and J. H. Koolen. Combinatorics of lopsided sets. *Eur. J. Comb.*, 27(5):669–689, 2006.
- A. Blumer and N. Littlestone. Learning faster than promised by the Vapnik-Chervonenkis dimension. *Discrete Applied Mathematics*, 24(1-3):47–53, 1989.
- A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *Journal of the Association for Computing Machinery*, 36(4):929–965, 1989.
- O. Bousquet, Y. Klochkov, and N. Zhivotovskiy. Sharper bounds for uniformly stable algorithms. *arXiv preprint arXiv:1910.07833*, 2019.
- M. Braverman, G. Kol, S. Moran, and R. R. Saxena. Convex set disjointness, distributed learning of halfspaces, and LP feasibility. *arXiv:1909.03547*, 2019.
- J. Chalopin, V. Chepoi, S. Moran, and M. K. Warmuth. Unlabeled sample compression schemes and corner peelings for ample and maximum classes. In *Proceedings of the 46<sup>th</sup> International Colloquium on Automata, Languages and Programming*, 2019.
- A. Daniely and S. Shalev-Shwartz. Optimal learners for multiclass problems. In *Proceedings of the 27<sup>th</sup> Conference on Learning Theory*, 2014.
- A. Daniely, S. Sabato, S. Ben-David, and S. Shalev-Shwartz. Multiclass learnability and the ERM principle. *The Journal of Machine Learning Research*, 16(1):2377–2404, 2015.
- L. Danzer, B. Grünbaum, and V. Klee. *Helly’s theorem and its relatives*. Proceedings of symposia in pure mathematics: Convexity. American Mathematical Society, 1963.
- M. Darnstädt. The optimal PAC bound for intersection-closed concept classes. *Information Processing Letters*, 115(4):458–461, 2015.
- A. Ehrenfeucht, D. Haussler, M. Kearns, and L. Valiant. A general lower bound on the number of examples needed for learning. *Information and Computation*, 82(3):247–261, 1989.
- V. Feldman and J. Vondrak. High probability generalization bounds for uniformly stable algorithms with nearly optimal rate. In *Proceedings of the 32<sup>nd</sup> Conference on Learning Theory*, pages 1270–1279, 2019.
- S. Floyd and M. Warmuth. Sample compression, learnability, and the Vapnik-Chervonenkis dimension. *Machine Learning*, 21(3):269–304, 1995.

- S. A. Goldman and M. J. Kearns. On the complexity of teaching. *Journal of Computer and System Sciences*, 50:20–31, 1995.
- S. Hanneke. *Theoretical Foundations of Active Learning*. PhD thesis, Machine Learning Department, School of Computer Science, Carnegie Mellon University, 2009.
- S. Hanneke. The optimal sample complexity of PAC learning. *Journal of Machine Learning Research*, 17(38):1–15, 2016a.
- S. Hanneke. Refined error bounds for several learning algorithms. *Journal of Machine Learning Research*, 17(135):1–55, 2016b.
- S. Hanneke and A. Kontorovich. Optimality of SVM: Novel proofs and tighter bounds. *Theoretical Computer Science*, 796:99–113, 2019.
- S. Hanneke and L. Yang. Minimax analysis of active learning. *Journal of Machine Learning Research*, 16(12):3487–3602, 2015.
- D. Haussler, N. Littlestone, and M. Warmuth. Predicting  $\{0, 1\}$ -functions on randomly drawn points. *Information and Computation*, 115(2):248–292, 1994.
- E. Helly. Über mengen konvexer körper mit gemeinschaftlichen punkte. *Jahresbericht der Deutschen Mathematiker-Vereinigung*, 32:175–176, 1923.
- D. Helmbold, R. Sloan, and M. Warmuth. Learning nested differences of intersection-closed concept classes. *Machine Learning*, 5(2):165–196, 1990.
- W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- D. Kane, R. Livni, S. Moran, and A. Yehudayoff. On communication complexity of classification problems. In *Proceedings of the 32<sup>nd</sup> Conference on Learning Theory*, pages 1903–1943, 2019.
- M. Kearns and U. Vazirani. *An Introduction to Computational Learning Theory*. The MIT Press, 1994.
- J. Lawrence. Lopsided sets and orthant-intersection of convex sets. *Pacific J. Math.*, 104:155–173, 1983.
- F. W. Levi. On Helly’s theorem and the axioms of convexity. *J. Indian Math. Soc.*, 15:65–76, 1951.
- N. Littlestone and M. Warmuth. Relating data compression and learnability. *Unpublished manuscript*, 1986.
- P. Long and R. Long. On the complexity of proper distribution-free learning of linear classifiers. In *Proceedings of the 31<sup>st</sup> International Conference on Algorithmic Learning Theory*, 2020.
- O. Montasser, S. Hanneke, and N. Srebro. VC classes are adversarially robustly learnable, but only improperly. In *Proceedings of the 32<sup>nd</sup> Conference on Learning Theory*, 2019.

- S. Moran and M. Warmuth. Labeled compression schemes for extremal classes. In *Proceedings of the 27<sup>th</sup> International Conference on Algorithmic Learning Theory*, pages 34–49. Springer, 2016.
- R. Motwani and P. Raghavan. *Randomized algorithms*. Chapman & Hall/CRC, 2010.
- H. Simon. An almost optimal PAC algorithm. In *Proceedings of the 28<sup>th</sup> Conference on Learning Theory*, 2015.
- A. W. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes*. Springer, 1996.
- R. van Handel. The universal Glivenko–Cantelli property. *Probability Theory and Related Fields*, 155(3-4):911–934, 2013.
- V. Vapnik and A. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. In *Proc. USSR Acad. Sci.*, 1968.
- V. Vapnik and A. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264–280, 1971.
- V. Vapnik and A. Chervonenkis. *Theory of Pattern Recognition*. Nauka, Moscow, 1974.
- V. Vovk, H. Papadopoulos, and A. Gammernan. *Measures of Complexity*. Springer, 2015.
- N. Zhivotovskiy. Optimal learning via local entropies and sample compression. In *Proceedings of the 30<sup>th</sup> Conference on Learning Theory*, pages 2023–2065, 2017.
- N. Zhivotovskiy and S. Hanneke. Localization of VC classes: Beyond local Rademacher complexities. *Theoretical Computer Science*, 742:27–49, 2018.

## Appendix A. Omitted proofs of Section 2

**Proof of Lemma 5** We note that many of these arguments were given in some form in the course of the proofs of Kane, Livni, Moran, and Yehudayoff (2019). However, we provide the proofs explicitly here, particularly since our context is slightly different.

First, on a technical note, we remark that because of our assumption  $|\mathbb{C}| \geq 3$  stated initially, all of  $k_o, k_p, k_w$  are at least 2 (so that the restriction to  $k_p \geq 2$  in its definition does not affect the claims).

For the first claimed inequalities, given the  $k_o$  classifiers  $\mathbb{C}'$  in  $\mathbb{C}$  witnessing a hollow star, for any  $\ell < k_o$  the region  $\mathcal{X}_{\mathbb{C}', \ell}$  contains the hollow star, and hence the majority vote of the  $\mathbb{C}'$  classifiers is unrealizable on  $\mathcal{X}_{\mathbb{C}', \ell}$ . Therefore,  $k_p \geq k_o$ .

If  $k_w = \infty$ , then trivially  $k_p \leq k_w$ , so suppose  $k_w < \infty$ . Suppose some finite multiset  $\mathbb{C}' \subseteq \mathbb{C}$  has no  $h \in \mathbb{C}$  that coincides with  $\text{Majority}(\mathbb{C}')$  on  $\mathcal{X}_{\mathbb{C}', k_w}$ . Then  $S = \{(x, \text{Majority}(\mathbb{C}')(x)) : x \in \mathcal{X}_{\mathbb{C}', k_w}\}$  is an unrealizable set. Therefore, it contains a subset  $W$  of size at most  $k_w$  that is also unrealizable. But (by definition of  $\mathcal{X}_{\mathbb{C}', k_w}$ ) each point  $(x, y)$  in  $W$  contradicts strictly fewer than  $|\mathbb{C}'|/k_w$  elements in  $\mathbb{C}'$ , so that there must be at least one  $h \in \mathbb{C}'$  that survives: a contradiction. Therefore,  $k_p \leq k_w$ .

It remains to show that these quantities are all equal when  $k_w < \infty$ . Let  $S$  be an unrealizable set such that the smallest unrealizable subset  $W$  has size  $k_w < \infty$ . If  $W$  is not a hollow star set, then there exists a point  $(x, y) \in W$  such that  $(W \setminus \{(x, y)\}) \cup \{(x, -y)\}$  is also not realizable, which implies  $W \setminus \{(x, y)\}$  is also not realizable: a contradiction. Therefore,  $W$  is a hollow star (indeed, any unreducible unrealizable set is a hollow star), and hence  $k_w \leq k_o$ . The equalities then follow from the first claim, established above.

For the remaining claim, if  $\mathbb{C}$  is closed and yet  $k_w = \infty$ , it implies there is a sequence  $S_i$  of sets with  $\mathbb{C}[S_i] = \emptyset$  for which the smallest  $W_i \subseteq S_i$  with  $\mathbb{C}[W_i] = \emptyset$  has  $\lim_{i \rightarrow \infty} |W_i| = \infty$ , yet each  $W_i$  is finite (due to the ‘‘closedness’’ assumption). As above, these minimum-size  $W_i$  sets must be hollow star sets, which implies there is no finite bound on the size of all finite hollow star sets. Therefore  $k_o = \infty$ , and the inequality established above then implies  $k_p = \infty$  as well. ■

**Proof for Example 1** This follows from Proposition 2.8 in (Braverman, Kol, Moran, and Saxena, 2019). Indeed, this implies that  $k_o \leq n + 2$  since if a finite unrealizable sample  $S$  is a hollow star set, then in particular every proper subsample of  $S$  must be realizable; however, by Proposition 2.8 in (Braverman, Kol, Moran, and Saxena, 2019) the set  $S$  must contain an unrealizable subsample of size at most  $n + 2$ , and hence it must be that  $|S| \leq n + 2$ . To see that  $k_o \geq n + 2$ , pick  $x_1, \dots, x_{n+1} \in \mathbb{R}^n$  to be the vertices of a simplex, and choose

$$S = \left\{ (x_1, +1), \dots, (x_{n+1}, +1), \left( \frac{x_1 + \dots + x_{n+1}}{n+1}, -1 \right) \right\}.$$

We leave it to the reader to verify that the above  $S$  witnesses that  $k_o \geq |S| = n + 2$ .

It remains to show that  $k_p = n + 2$ . By Lemma 5 it suffices to show that  $k_p \leq n + 2$ . Let  $\mathbb{C}'$  be a finite collection of halfspaces in  $\mathbb{R}^n$ . We need to show that there exists a halfspace  $h$  which agrees with  $\text{Majority}(\mathbb{C}')$  on the set  $\mathcal{X}_{\mathbb{C}', n+2}$ . Let  $\mathcal{X}_+ \subseteq \mathcal{X}_{\mathbb{C}', n+2}$  denote the set of all points  $x \in \mathcal{X}_{\mathbb{C}', n+2}$  such that  $\text{Majority}(\mathbb{C}')(x) = +1$  and similarly let  $\mathcal{X}_- \subseteq \mathcal{X}_{\mathbb{C}', n+2}$  denote the set of all points  $x \in \mathcal{X}_{\mathbb{C}', n+2}$  such that  $\text{Majority}(\mathbb{C}')(x) = -1$ . We first claim that the convex hulls  $\text{conv}(\mathcal{X}_+)$  and  $\text{conv}(\mathcal{X}_-)$  are disjoint; indeed, otherwise by Proposition 2.8 in (Braverman, Kol, Moran, and Saxena, 2019) there exist  $S_+ \subseteq \mathcal{X}_+, S_- \subseteq \mathcal{X}_-$  such that  $\text{conv}(S_-) \cap \text{conv}(S_+) \neq \emptyset$  and  $|S_-| + |S_+| \leq n + 2$ . However, every  $x \in S_- \cup S_+$  is classified correctly by more than  $1 - \frac{1}{n+2}$  fraction of the halfspaces in  $\mathbb{C}'$  and hence there must be a halfspace in  $\mathbb{C}'$  that classifies correctly all points in  $S_+ \cup S_-$  and so  $\text{conv}(S_-) \cap \text{conv}(S_+) = \emptyset$ , which is a contradiction.

Having established that  $\text{conv}(S_-) \cap \text{conv}(S_+) = \emptyset$ , we are ready to finish the proof. Indeed, by the *Hyperplane Separation Theorem* there exists a linear function  $L : \mathbb{R}^n \rightarrow \mathbb{R}$  and a value  $v$  such that  $L(x) \leq v$  for every  $x \in S_-$  and  $L(x) \geq v$  for every  $x \in S_+$ . However, note that in fact  $L(x) < v$  for every  $x \in S_-$ : this follows because  $S_-$  is an open set (indeed, it can be written as a union of (finite) intersections of sets of the form  $h^{-1}(-1)$ , where  $h \in \mathbb{C}'$ , each of which is an open set). Thus, the halfspace  $\{x : L(x) \geq v\}$  agrees with  $\text{Majority}(\mathbb{C}')$  on the set  $\mathcal{X}_{\mathbb{C}', n+2}$ , as required.<sup>4</sup> ■

**Proof for Example 2** Let us sketch the proof. The bound  $k_o \leq d + 1$  on the dual Helly number for extremal classes follows from the fact that: (i) the complement of an extremal class is extremal, and

4. Note that we use the definition of Halfspaces where the positive side of each halfspace is closed: i.e., every halfspace is of the form  $\text{sign}(L(x) - v)$ , where  $\text{sign}(0) = +1$ .

(ii) the *one-inclusion graph* on an extremal graph projected to any data set is connected (we refer the reader to [Bandelt, Chepoi, Dress, and Koolen, 2006](#) and [Moran and Warmuth, 2016](#) for these definitions and results). In particular, it follows that the one-inclusion graph of the complement of  $\mathbb{C}$  projected to any data set is connected.

Since a hollow star corresponds to a one-inclusion graph where the star's center is an isolated vertex in the one-inclusion graph of the complement of  $\mathbb{C}$  projected to the points, it must be the only vertex in the complement, which means any strict subset of the points is shattered by  $\mathbb{C}$ . This immediately implies  $k_o - 1 \leq d$ .

The case of intersection-closed classes is also straightforward. Let  $S = \{(x_1, y_1), \dots, (x_k, y_k)\}$  be a finite hollow star set, and let  $h_1, \dots, h_k$  be elements of  $\mathbb{C}$  such that  $\{j : h_i(x_j) \neq y_j\} = \{i\}$  for each  $i$ . Denote by  $m$  the number of  $y_i$  values equal  $-1$ . We will first show that  $m \leq 1$ . Indeed, if there are at least two values  $y_i$  and  $y_j$  ( $i \neq j$ ) both equal to  $-1$ , then letting  $h_0$  be the classifier in  $\mathbb{C}$  with  $\{x : h_0(x) = 1\} = \{x : h_i(x) = 1\} \cap \{x : h_j(x) = 1\}$  (which exists since  $\mathbb{C}$  is intersection-closed), we would have  $h_0$  correct on all of  $S$ : a contradiction to  $S$  being unrealizable. Next we argue that there are at most  $d$  values  $y_i$  equal  $1$ : that is,  $k - m \leq d$ . Suppose  $y_{i_1}, \dots, y_{i_{k-m}}$  are equal  $1$ . Then for any  $y'_{i_1}, \dots, y'_{i_{k-m}} \in \mathcal{Y}$ , there exists a classifier  $h \in \mathbb{C}$  with  $\{x : h(x) = 1\} = \bigcap \{\{x : h_{i_j}(x) = 1\} : y'_{i_j} = -1\}$ , and this  $h$  has  $(h(x_{i_1}), \dots, h(x_{i_{k-m}})) = (y'_{i_1}, \dots, y'_{i_{k-m}})$ ; thus, the set  $\{x_{i_1}, \dots, x_{i_{k-m}}\}$  is shattered by  $\mathbb{C}$ , and hence has size at most  $d$ . Altogether, we have that  $k = (k - m) + m \leq d + 1$ . Since this applies to *any* finite hollow star set  $S$ , we conclude that  $k_o \leq d + 1$ .  $\blacksquare$

**Proof for Example 5** We provide the details for the final claim that the augmented class has  $k_w = k_p = k_o = 3$ . Since the class is intersection-closed with  $d = 2$ , Example 2 implies  $k_o \leq 3$ . Thus, since  $k_o \geq 3$  (witnessed by the hollow star set  $\{(1, 1), (2, -1), (3, 1)\}$ ), Lemma 5 implies it suffices to show the class is closed. Let  $S$  be an infinite unrealizable set; we aim to show it must contain a finite unrealizable subset. If  $S$  contains  $(x, 1)$  and  $(x, -1)$  for some  $x$ , then  $\{(x, 1), (x, -1)\}$  is a finite unrealizable subset. Otherwise, suppose no such  $x$  exists. Notice that  $S$  must contain at least one point having label  $1$ , since otherwise  $S$  would be realizable by the constant classifier  $h_{-1} = -1$ , which is in the class. Let  $\underline{x} = \inf\{x : (x, 1) \in S\}$  and  $\bar{x} = \sup\{x : (x, 1) \in S\}$ . Note that every  $(x, y) \in S$  with  $x > \bar{x}$  has  $y = -1$  and similarly every  $(x, y) \in S$  with  $x < \underline{x}$  has  $y = -1$ . In particular, this implies that if every  $(x, y) \in S$  with  $\underline{x} < x < \bar{x}$  has  $y = 1$ , then  $S$  would be realizable by a classifier corresponding to one of the intervals  $(\underline{x}, \bar{x})$  (if  $(\underline{x}, 1) \notin S$  and  $(\bar{x}, 1) \in S$ ),  $[\underline{x}, \bar{x})$  (if  $(\underline{x}, 1) \in S$  and  $(\bar{x}, 1) \notin S$ ),  $[\underline{x}, \bar{x}]$  (if  $(\underline{x}, 1) \in S$  and  $(\bar{x}, 1) \in S$ ), or  $(\underline{x}, \bar{x}]$  (if  $(\underline{x}, 1) \notin S$  and  $(\bar{x}, 1) \notin S$ ): a contradiction. Therefore, there exists  $(x_2, -1) \in S$  with  $\underline{x} < x_2 < \bar{x}$ . By the definitions of  $\underline{x}$  and  $\bar{x}$ , this implies there exist finite  $x_1, x_3$  with  $x_1 < x_2 < x_3$  such that  $(x_1, 1)$  and  $(x_3, 1)$  are both in  $S$ . Then we have that the set  $\{(x_1, 1), (x_2, -1), (x_3, 1)\} \subset S$  is a finite unrealizable subset. Since this applies to *any* unrealizable infinite set  $S$ , we conclude that the class is closed.  $\blacksquare$

## Appendix B. Proofs of the upper bounds

**Proof of Theorem 8** Fix any target concept  $f^* \in \mathbb{C}$  and distribution  $\mathcal{P}$ . We first argue that, for any finite data sets  $S$  and  $T$  with labels consistent with  $f^*$ , the proper learner  $\mathbb{A}(S; T)$  outputs

$\hat{h} \in \mathbb{C}$  such that  $\hat{h}$  is correct on  $T$ . Note that this is trivially true if  $|S| < 4$ , since Step 1 returns  $\text{ERM}(S \cup T)$ , which is correct on  $T$  by definition. Now, for induction, suppose  $S$  is a correctly labeled finite data set such that, for any strict subset  $S' \subset S$ , and any correctly labeled finite data set  $T'$ ,  $\mathbb{A}(S'; T')$  returns a classifier in  $\mathbb{C}$  that is correct on  $T'$ . Now we extend this property to the full set  $S$ . Fix any correctly labeled finite data set  $T$ . Recalling the notation from the algorithm, define  $h_{\text{maj}}(x) = \text{Majority}(h_1(x), \dots, h_{k_p+1}(x))$  (breaking ties to favor label  $-1$ , say), and recalling the notation (1) let

$$\mathcal{X}_0 = \mathcal{X}_{\{h_1, \dots, h_{k_p+1}\}, k_p}.$$

By definition (from Step 5 in  $\mathbb{A}(S; T)$ ), the classifier  $\hat{h} = \mathbb{A}(S; T) \in \mathbb{C}$  has  $\hat{h}(x) = h_{\text{maj}}(x)$  on every  $x \in \mathcal{X}_0$ . Furthermore, since  $h_i = \mathbb{A}(S_0; T \cup S_i)$  for each  $i$ , where  $S_0 \subset S$ , the inductive hypothesis implies  $h_i$  is correct on  $T$ . Therefore, all  $h_i$  agree on the labels in  $T$ , and hence the set of points in  $T$  is contained in  $\mathcal{X}_0$ , which implies  $\hat{h}$  is correct on  $T$  as well. By induction, this implies that for any correctly labeled finite data sets  $S$  and  $T$ ,  $\hat{h} = \mathbb{A}(S; T)$  is correct on  $T$ .

Next we argue that, for any  $m_0 \in \mathbb{N}$  and any  $\delta_0 \in (0, 1)$ , if  $T$  is any correctly labeled finite data set and  $S$  is an i.i.d. labeled data set of size  $m_0$ , with  $\mathcal{P}$  marginal distribution and  $f^*$  labels, then with probability at least  $1 - \delta_0$ , we have that  $\hat{h}$  satisfies

$$\text{er}(\hat{h}) \leq \frac{c \cdot k_p^2}{m_0} \left( \text{dLog}(k_p) + \text{Log}\left(\frac{1}{\delta_0}\right) \right), \quad (2)$$

where  $c \geq 1$  is an appropriate finite numerical constant. Note that this would imply Theorem 8, since setting  $\delta_0 = \delta$  and  $m_0$  of size proportional to the claimed bound on  $\mathcal{M}_{\text{prop}}(\varepsilon, \delta)$  from Theorem 8, the bound (2) is less than  $\varepsilon$ .

If  $m_0 < 4$ , the claim trivially holds, as the bound is greater than 1. In particular, this will be our base case in an inductive argument. Now, for induction, suppose  $m \geq 4$ , and that for any  $\delta_0 \in (0, 1)$  and any  $m_0 < m$ , if  $S$  is an i.i.d. data set of size  $m_0$  (with  $\mathcal{P}$  marginal distribution and  $f^*$  labels) and  $T$  is any finite data set with  $f^*$  labels, then with probability at least  $1 - \delta_0$  the inequality (2) holds for  $\hat{h} = \mathbb{A}(S; T)$ .

Next we extend this claim to hold for  $m_0 = m$ . Fix any  $\delta_0 \in (0, 1)$ . If  $m < 160 \text{Log}\left(\frac{6k_p^2}{\delta_0}\right)$ , the inequality trivially holds since the bound is greater than 1 (for sufficiently large choice of  $c$ ), so suppose  $m \geq 160 \text{Log}\left(\frac{6k_p^2}{\delta_0}\right)$ . Consider the sets  $S_i$  and classifiers  $h_i$  as defined in the specification of the algorithm  $\mathbb{A}(S; T)$  above Theorem 8, and with a slight abuse of notation we also use  $S_i$  to denote the *unlabeled* portion of the set  $S_i$  (i.e., the points  $x$  such that  $(x, f^*(x)) \in S_i$ ). As argued above, each  $h_i$  is correct on  $T \cup S_i$ .

For any  $h$ , define  $\text{ER}(h) = \{x : h(x) \neq f^*(x)\}$ . We claim that

$$\text{ER}(\hat{h}) \subseteq \bigcup_{i,j:i \neq j} \text{ER}(h_i) \cap \text{ER}(h_j). \quad (3)$$

To see this, note that since  $\hat{h}$  agrees with  $h_{\text{maj}}$  on  $\mathcal{X}_0$ , we have

$$\text{ER}(\hat{h}) \subseteq (\mathcal{X} \setminus \mathcal{X}_0) \cup (\mathcal{X}_0 \cap \text{ER}(h_{\text{maj}})). \quad (4)$$

Furthermore, for any  $x \in \mathcal{X} \setminus \mathcal{X}_0$ , at least two values of  $i$  have  $h_i(x)$  different from the majority of the values  $h_1(x), \dots, h_{k_p+1}(x)$ , which means there are at least two classifiers predicting each label

in  $\mathcal{Y}$ , and hence there are at least two classifiers  $h_i$  with  $h_i(x) \neq f^*(x)$ . Furthermore, any  $x$  with  $h_{\text{maj}}(x) \neq f^*(x)$  certainly also has at least two  $h_i$  classifiers with  $h_i(x) \neq f^*(x)$ . Therefore, the set on the right hand side of (4) is contained within  $\bigcup_{i,j:i \neq j} \text{ER}(h_i) \cap \text{ER}(h_j)$ , and (3) follows.

In particular, (3) implies

$$\text{er}(\hat{h}) = \mathcal{P}(\text{ER}(\hat{h})) \leq \mathcal{P}\left(\bigcup_{i,j:i < j} \text{ER}(h_i) \cap \text{ER}(h_j)\right) \leq \sum_{i,j:i < j} \mathcal{P}(\text{ER}(h_i) \cap \text{ER}(h_j)). \quad (5)$$

The remainder of the proof will establish that each term  $\mathcal{P}(\text{ER}(h_i) \cap \text{ER}(h_j))$  is small with high probability. This is achieved using a ‘‘Win-Win’’ argument showing that for every distinct  $i, j$ , either  $\mathcal{P}(\text{ER}(h_i)) = \text{er}(h_i)$  is small, or else  $\mathcal{P}(\text{ER}(h_j) | \text{ER}(h_i))$  is small. In either case, it will follow that  $\mathcal{P}(\text{ER}(h_i) \cap \text{ER}(h_j))$  is small.

Specifically, we follow a ‘‘conditioning’’ argument.<sup>5</sup> We claim that, with probability at least  $1 - \delta_0/3$ , for every pair  $i, j$  with  $i < j$ , either

$$\text{er}(h_i) < \frac{320}{m} \ln\left(\frac{6k_p^2}{\delta_0}\right) \quad (6)$$

or else

$$|\text{ER}(h_i) \cap (S_j \setminus S_i)| \geq \text{er}(h_i)m/80, \quad (7)$$

where here the notation  $S_j \setminus S_i$  denotes the set of samples from  $S$  that are in  $S_j$  but not  $S_i$  (distinguished by their *indices*, so that if  $S$  contains two copies of some  $x \in \mathcal{X}$  and one is in  $S_i$  and the other in  $S_j$ , then the latter will still appear in  $S_j \setminus S_i$ ).

Toward establishing the above claim, note that for each distinct  $i, j$  we have

$$\begin{aligned} & \Pr\left(|\text{ER}(h_i) \cap (S_j \setminus S_i)| < \text{er}(h_i)m/80 \text{ and } \text{er}(h_i) \geq \frac{320}{m} \ln\left(\frac{6k_p^2}{\delta_0}\right)\right) \\ & \leq \Pr\left(|\text{ER}(h_i) \cap (S_j \setminus S_i)| < (1/2)\text{er}(h_i) | S_j \setminus S_i \text{ and } \text{er}(h_i) \geq \frac{320}{m} \ln\left(\frac{6k_p^2}{\delta_0}\right) \text{ and } |S_j \setminus S_i| \geq \frac{m}{40}\right) \\ & \quad + \Pr\left(|S_j \setminus S_i| < \frac{m}{40}\right). \end{aligned} \quad (8)$$

We begin with bounding the second term. Note that  $\mathbb{E}[|S_j \setminus S_i| | S_i] = \left(1 - \frac{\lfloor m/4 \rfloor}{\lfloor m/2 \rfloor}\right) \lfloor m/4 \rfloor \geq m/20$  (noting that, since  $m \geq 20$ , we have  $\lfloor m/4 \rfloor \geq m/5$  and  $\lfloor m/2 \rfloor \geq m/3$ ). Therefore, by a multiplicative Chernoff bound<sup>6</sup> conditioned on  $S_i$ , and the law of total probability, we have  $\Pr(|S_j \setminus S_i| < m/40) \leq e^{-m/160} \leq \frac{\delta_0}{6k_p^2}$ .

5. Conditioning arguments of this type originate in the work of [Hanneke \(2009\)](#) on ERM bounds and active learning, and were later used to analyze several different learning algorithms (e.g., [Darnstädt, 2015](#); [Hanneke, 2016b](#); [Zhivotovskiy and Hanneke, 2018](#)). Notably, the argument was used by [Simon \(2015\)](#) to analyze majority votes of independent ERMs, and was used by [Hanneke \(2016a\)](#) in the proof of the optimal PAC sample complexity. Its use in the present proof most closely follows this latter work.

6. As proven by [Hoeffding, 1963](#), the moment generating function for sampling without replacement is upper bounded by the moment generating function for sampling with replacement, and hence the usual Chernoff bounds for sampling with replacement also hold for sampling without replacement.

Next we bound the first term in (8). We note that, conditioned on  $|S_j \setminus S_i|$ , the samples in  $S_j \setminus S_i$  are i.i.d. (with distribution  $\mathcal{P}$ ) and independent of  $h_i$ . Thus, a multiplicative Chernoff bound implies (almost surely)

$$\Pr\left(|\text{ER}(h_i) \cap (S_j \setminus S_i)| < (1/2)\text{er}(h_i)|S_j \setminus S_i| \mid h_i, |S_j \setminus S_i|\right) \leq e^{-(1/8)\text{er}(h_i)|S_j \setminus S_i|}.$$

Therefore, the first term in (8) is bounded by

$$\mathbb{E}\left[e^{-(1/8)\text{er}(h_i)|S_j \setminus S_i|} \mathbb{1}\left[\text{er}(h_i) \geq \frac{320}{m} \ln\left(\frac{6k_p^2}{\delta_0}\right) \text{ and } |S_j \setminus S_i| \geq m/40\right]\right] \leq \frac{\delta_0}{6k_p^2}.$$

Altogether, we have that (8) is at most  $\frac{\delta_0}{3k_p^2}$ . Finally, by a union bound over all pairs  $i, j$  with  $i < j$ , we conclude that with probability at least  $1 - \frac{\delta_0}{3}$ , for every  $i, j$  with  $i < j$ , at least one of (6) or (7) holds, as claimed.

Since  $h_j$  is correct on  $S_j$ , it is certainly correct on  $\text{ER}(h_i) \cap (S_j \setminus S_i)$ . Also note that the samples in  $\text{ER}(h_i) \cap (S_j \setminus S_i)$  are conditionally i.i.d. given  $h_i$  and  $|\text{ER}(h_i) \cap (S_j \setminus S_i)|$ , with conditional distribution  $\mathcal{P}(\cdot | \text{ER}(h_i))$ . We can therefore apply the classic PAC bound for ERM (Vapnik and Chervonenkis, 1974; Blumer, Ehrenfeucht, Haussler, and Warmuth, 1989), under the conditional distribution given  $h_i$  and  $|\text{ER}(h_i) \cap (S_j \setminus S_i)|$ , together with the law of total probability, to obtain that, for any distinct  $i, j$ , with probability at least  $1 - \frac{\delta_0}{3k_p^2}$ ,

$$\mathcal{P}(\text{ER}(h_j) | \text{ER}(h_i)) \leq \frac{2/\ln(2)}{|\text{ER}(h_i) \cap (S_j \setminus S_i)|} \left( d \text{Log}\left(\frac{2e|\text{ER}(h_i) \cap (S_j \setminus S_i)|}{d}\right) + \text{Log}\left(\frac{6k_p^2}{\delta_0}\right) \right).$$

(interpreting the bound to be infinite in the case  $|\text{ER}(h_i) \cap (S_j \setminus S_i)| = 0$ ). By the union bound, this holds simultaneously for all  $i, j$  with  $i < j$  with probability at least  $1 - \frac{\delta_0}{3}$ . Combining this with the event established above, by the union bound and monotonicity of  $x \mapsto (1/x)\text{Log}(ax)$ , with probability at least  $1 - \frac{2}{3}\delta_0$ , every  $i, j$  with  $i < j$  either satisfy (6) or

$$\mathcal{P}(\text{ER}(h_j) | \text{ER}(h_i)) \leq \frac{160/\ln(2)}{\text{er}(h_i)m} \left( d \text{Log}\left(\frac{\text{er}(h_i)m e}{40d}\right) + \text{Log}\left(\frac{6k_p^2}{\delta_0}\right) \right). \quad (9)$$

Since  $S_0$  is i.i.d. (with marginal distribution  $\mathcal{P}$  and  $f^*$  labels) with  $m/2 \leq |S_0| < m$ , the inductive hypothesis and the union bound imply that, with probability at least  $1 - \frac{\delta_0}{3}$ , every  $h_i$  has

$$\text{er}(h_i) \leq \frac{ck_p^2}{m/2} \left( d \text{Log}(k_p) + \text{Log}\left(\frac{3(k_p + 1)}{\delta_0}\right) \right).$$

Plugging this into the log in (9), by the union bound we have that, with probability at least  $1 - \delta_0$ , every  $i, j$  with  $i < j$  either satisfy (6) or  $\mathcal{P}(\text{ER}(h_j) | \text{ER}(h_i))$  is upper bounded by

$$\frac{160/\ln(2)}{\text{er}(h_i)m} \left( d \text{Log}\left(ck_p^2 \left( \text{Log}(k_p) + \frac{1}{d} \text{Log}\left(\frac{3(k_p + 1)}{\delta_0}\right) \right) \right) + \text{Log}\left(\frac{6k_p^2}{\delta_0}\right) \right).$$

In either case (i.e., whether (6) holds or not), on this event

$$\begin{aligned} \mathcal{P}(\text{ER}(h_i) \cap \text{ER}(h_j)) &= \text{er}(h_i) \mathcal{P}(\text{ER}(h_j) | \text{ER}(h_i)) \\ &\leq \frac{320}{m} \left( \text{dLog} \left( ck_p^2 \left( \text{Log}(k_p) + \frac{1}{d} \text{Log} \left( \frac{3(k_p+1)}{\delta_0} \right) \right) \right) + \text{Log} \left( \frac{6k_p^2}{\delta_0} \right) \right). \end{aligned}$$

Combining this with (5) we have that, with probability at least  $1 - \delta_0$ ,

$$\text{er}(\hat{h}) \leq \binom{k_p+1}{2} \frac{320}{m} \left( \text{dLog} \left( ck_p^2 \left( \text{Log}(k_p) + \frac{1}{d} \text{Log} \left( \frac{3(k_p+1)}{\delta_0} \right) \right) \right) + \text{Log} \left( \frac{6k_p^2}{\delta_0} \right) \right).$$

Finally, simplifying the expression, and noting that the constant  $c$  only appears in a logarithmic term, one can verify that for a sufficiently large choice of numerical constant  $c$  (e.g., any  $c \geq 2^{17}$  would suffice), the right hand side is at most

$$\frac{ck_p^2}{m} \left( \text{dLog}(k_p) + \text{Log} \left( \frac{1}{\delta_0} \right) \right),$$

which extends the inductive hypothesis to  $m_0 = m$ . The result now follows by the principle of induction.  $\blacksquare$

**Proof of Theorem 9** Fix any target concept  $f^* \in \mathbb{C}$  and any distribution  $\mathcal{P}$  on  $\mathcal{X}$ . We will argue that, for any finite labeled data set  $S$  with  $f^*$  labels, the proper learner  $\mathbb{A}_{\text{ERM}}(S)$  outputs  $\hat{h} \in \mathbb{C}$  correct on  $S$ , and in the case that  $S$  is  $m_0$  i.i.d. training examples (with  $\mathcal{P}$  marginal distribution and  $f^*$  labels), then for any  $\delta_0 \in (0, 1)$ , with probability at least  $1 - \delta_0$ , the classifier  $\hat{h} = \mathbb{A}(S)$  satisfies

$$\text{er}(\hat{h}) \leq \frac{ck_p^3}{m_0} \left( \text{dLog}(k_p) + \text{Log} \left( \frac{1}{\delta_0} \right) \right), \quad (10)$$

where  $c \geq 1$  is an appropriate finite numerical constant. Note that Theorem 9 would immediately follow from this (since it holds for any  $\mathcal{P}$  and any  $f^* \in \mathbb{C}$ ), taking  $\delta_0 = \delta$ , and noting that  $m_0$  of size proportional to the stated bound on  $\mathcal{M}_{\mathbb{A}_{\text{ERM}}}(\varepsilon, \delta)$  makes the right hand side of (10) less than  $\varepsilon$ .

If  $m_0 < k_p + 1$ , the algorithm returns  $\hat{h} = \text{ERM}(S)$  in Step 1, which is an element of  $\mathbb{C}$  that is correct on  $S$  by definition; furthermore, the inequality trivially holds in this case, as the right hand side is greater than 1. These values of  $m_0$  will serve as our base case in an inductive argument. Now, for induction, suppose  $m \geq k_p + 1$ , and that for any  $\delta_0 \in (0, 1)$  and  $m_0 < m$ , for any correctly labeled data set  $S$  of size  $m_0$ , the classifier  $\hat{h}$  returned by  $\mathbb{A}_{\text{ERM}}(S)$  is in  $\mathbb{C}$  and is correct on  $S$ , and in the case that  $S$  is i.i.d. (with  $\mathcal{P}$  marginal distribution and  $f^*$  labels), then with probability at least  $1 - \delta_0$  (10) holds.

Next we extend this claim to  $m_0 = m$ . Consider a run of  $\mathbb{A}_{\text{ERM}}(S)$  with a correctly labeled finite data set  $S$  of size  $m$ . Consider the sets  $S_i$  and classifiers  $h_i$  as defined in the algorithm, and with a slight abuse of notation we also use  $S_i$  to denote the unlabeled portion of  $S_i$  (i.e., the points  $x$  such that  $(x, f^*(x)) \in S_i$ ). Define  $h_{\text{maj}}(x) = \text{Majority}(h_1(x), \dots, h_{k_p+1}(x))$  (breaking ties to favor label  $-1$ , say), and as before using the notation (1),

$$\mathcal{X}_0 = \mathcal{X}_{\{h_1, \dots, h_{k_p+1}\}, k_p}.$$

Since each  $h_i$  is in  $\mathbb{C}$  (by the inductive hypothesis), the classifier  $\hat{h} = \text{Proj}_{\mathbb{C}}(h_1, \dots, h_{k_p+1})$  in Step 4 is well defined, and by definition,  $\hat{h} \in \mathbb{C}$  and has  $\hat{h}(x) = h_{\text{maj}}(x)$  on every  $x \in \mathcal{X}_0$ . Note that since every  $(x, y) \in S$  is included in just one set  $S_i$ , and hence is in every set  $\bigcup_{j' \neq j} S_{j'}$  except  $j = i$ , and by the inductive hypothesis every  $j \neq i$  has  $h_j$  correct on  $\bigcup_{j' \neq j} S_{j'}$ , we see that every  $(x, y) \in S$  has  $x \in \mathcal{X}_0$ , and  $h_{\text{maj}}(x) = y$ , so that this extends the claim that  $\hat{h}$  is in  $\mathbb{C}$  and is correct on  $S$  for the inductive proof, and all that remains is to extend the bound on the error rate to hold for  $m_0 = m$ .

Toward this end, consider the case that  $S$  is an i.i.d. data set of size  $m$  (with marginal distribution  $\mathcal{P}$  and  $f^*$  labels). Fix any  $\delta_0 \in (0, 1)$ . By the inductive hypothesis each  $h_i$  is correct on  $\bigcup_{j \neq i} S_j$  and has  $h_i \in \mathbb{C}$ . We will follow a similar ‘‘conditioning’’ argument to that used in the proof of Theorem 8. As in that proof, define  $\text{ER}(h) = \{x : h(x) \neq f^*(x)\}$  for any classifier  $h$ . Since  $\hat{h}$  agrees with  $h_{\text{maj}}$  on  $\mathcal{X}_0$ , we have

$$\text{ER}(\hat{h}) \subseteq (\mathcal{X} \setminus \mathcal{X}_0) \cup (\mathcal{X}_0 \cap \text{ER}(h_{\text{maj}})). \quad (11)$$

Furthermore, for any  $x \in \mathcal{X} \setminus \mathcal{X}_0$ , at least two values of  $i$  have  $h_i(x)$  different from the majority of the values  $h_1(x), \dots, h_{k_p+1}(x)$ , which means there are at least two classifiers predicting each label in  $\mathcal{Y}$ , and hence there are at least two classifiers  $h_i$  with  $h_i(x) \neq f^*(x)$ . Furthermore, any  $x$  with  $h_{\text{maj}}(x) \neq f^*(x)$  certainly also has at least two  $h_i$  classifiers with  $h_i(x) \neq f^*(x)$ . Therefore, the set on the right hand side of (11) is contained within  $\bigcup_{i,j:i \neq j} \text{ER}(h_i) \cap \text{ER}(h_j)$ . In particular, this implies

$$\text{er}(\hat{h}) = \mathcal{P}(\text{ER}(\hat{h})) \leq \mathcal{P}\left(\bigcup_{i,j:i < j} \text{ER}(h_i) \cap \text{ER}(h_j)\right) \leq \sum_{i,j:i < j} \mathcal{P}(\text{ER}(h_i) \cap \text{ER}(h_j)). \quad (12)$$

The remainder of the proof will establish that each term  $\mathcal{P}(\text{ER}(h_i) \cap \text{ER}(h_j))$  is small with high probability.

For any  $i$ , we have

$$\begin{aligned} & \Pr\left(|\text{ER}(h_i) \cap S_i| < (1/2)\text{er}(h_i)|S_i| \text{ and } \text{er}(h_i) \geq \frac{8}{|S_i|} \ln\left(\frac{3(k_p+1)}{\delta_0}\right)\right) \\ &= \mathbb{E}\left[\Pr\left(|\text{ER}(h_i) \cap S_i| < (1/2)\text{er}(h_i)|S_i| \mid \text{er}(h_i)\right) \mathbf{1}\left[\text{er}(h_i) \geq \frac{8}{|S_i|} \ln\left(\frac{3(k_p+1)}{\delta_0}\right)\right]\right]. \end{aligned} \quad (13)$$

Since  $S_i$  is excluded from the training set  $\bigcup_{j \neq i} S_j$  producing  $h_i$ , we have that  $S_i$  and  $h_i$  are independent random variables. Therefore, a multiplicative Chernoff bound implies that (almost surely)

$$\Pr\left(|\text{ER}(h_i) \cap S_i| < (1/2)\text{er}(h_i)|S_i| \mid \text{er}(h_i)\right) \leq \exp(-\text{er}(h_i)|S_i|/8),$$

so that (13) is at most  $\frac{\delta_0}{3(k_p+1)}$ . In other words, with probability at least  $1 - \frac{\delta_0}{3(k_p+1)}$ , either

$$\text{er}(h_i) < \frac{8}{|S_i|} \log\left(\frac{3(k_p+1)}{\delta_0}\right) \quad (14)$$

or else

$$|\text{ER}(h_i) \cap S_i| \geq (1/2)\text{er}(h_i)|S_i|. \quad (15)$$

By the union bound, with probability at least  $1 - \frac{\delta_0}{3}$ , every  $i$  satisfies at least one of (14) or (15).

For distinct  $j, i$ , since  $h_j$  is correct on  $\bigcup_{i' \neq j} S_{i'} \supseteq S_i$  it is certainly correct on  $\text{ER}(h_i) \cap S_i$ . Also note that the samples in  $\text{ER}(h_i) \cap S_i$  are conditionally i.i.d. given  $h_i$  and  $|\text{ER}(h_i) \cap S_i|$ , with conditional distribution  $\mathcal{P}(\cdot | \text{ER}(h_i))$ . Thus, by the classic PAC bound for ERM (Vapnik and Chervonensis, 1974; Blumer, Ehrenfeucht, Haussler, and Warmuth, 1989) (applied under the conditional distribution given  $h_i$  and  $|\text{ER}(h_i) \cap S_i|$ ) and the law of total probability, with probability at least  $1 - \frac{\delta_0}{3k_p^2}$ ,

$$\mathcal{P}(\text{ER}(h_j) | \text{ER}(h_i)) \leq \frac{2/\ln(2)}{|\text{ER}(h_i) \cap S_i|} \left( \text{dLog} \left( \frac{2e|\text{ER}(h_i) \cap S_i|}{d} \right) + \text{Log} \left( \frac{6k_p^2}{\delta_0} \right) \right).$$

By the union bound (over  $i, j$  pairs with  $i < j$ , and combining with the event above) and monotonicity of  $x \mapsto (1/x)\text{Log}(ax)$ , with probability at least  $1 - \frac{2}{3}\delta_0$ , every pair  $i, j$  with  $i < j$  have either (14) or

$$\mathcal{P}(\text{ER}(h_j) | \text{ER}(h_i)) \leq \frac{4/\ln(2)}{\text{er}(h_i)|S_i|} \left( \text{dLog} \left( \frac{e \text{er}(h_i)|S_i|}{d} \right) + \text{Log} \left( \frac{6k_p^2}{\delta_0} \right) \right). \quad (16)$$

Since each  $h_i$  is correct on  $\bigcup_{j \neq i} S_j$ , which is itself an i.i.d. data set of size strictly smaller than  $m$  and no smaller than  $k_p \lfloor m/(k_p + 1) \rfloor \geq \frac{m}{3}$  (using the fact that  $k_p \geq 2$ , from the definition), the inductive hypothesis and the union bound imply that with probability at least  $1 - \frac{\delta_0}{3}$ , every  $h_i$  has

$$\text{er}(h_i) \leq \frac{3ck_p^3}{m} \left( \text{dLog}(k_p) + \text{Log} \left( \frac{3(k_p + 1)}{\delta_0} \right) \right).$$

Plugging this into the log in (16) above, together with the fact that  $|S_i| \geq \lfloor m/(k_p + 1) \rfloor \geq (1/2)m/(k_p + 1)$ , we have by the union bound that with probability at least  $1 - \delta_0$ , every pair  $i, j$  with  $i < j$  either have (14) or have that  $\mathcal{P}(\text{ER}(h_j) | \text{ER}(h_i))$  is upper bounded by

$$\frac{(8/\ln(2))(k_p + 1)}{\text{er}(h_i)m} \left( \text{dLog} \left( \frac{3}{2}eck_p^2 \left( \text{Log}(k_p) + \frac{1}{d} \text{Log} \left( \frac{6k_p^2}{\delta_0} \right) \right) \right) + \text{Log} \left( \frac{6k_p^2}{\delta_0} \right) \right).$$

In either case (i.e., whether (14) holds or not), on this event we have

$$\begin{aligned} \mathcal{P}(\text{ER}(h_i) \cap \text{ER}(h_j)) &= \text{er}(h_i) \mathcal{P}(\text{ER}(h_j) | \text{ER}(h_i)) \\ &\leq \frac{16(k_p + 1)}{m} \left( \text{dLog} \left( \frac{3}{2}eck_p^2 \left( \text{Log}(k_p) + \frac{1}{d} \text{Log} \left( \frac{6k_p^2}{\delta_0} \right) \right) \right) + \text{Log} \left( \frac{6k_p^2}{\delta_0} \right) \right). \end{aligned}$$

Combining this with (12), we conclude that on the above event of probability at least  $1 - \delta_0$ ,

$$\text{er}(\hat{h}) \leq \binom{k_p + 1}{2} \frac{16(k_p + 1)}{m} \left( \text{dLog} \left( \frac{3}{2}eck_p^2 \left( \text{Log}(k_p) + \frac{1}{d} \text{Log} \left( \frac{6k_p^2}{\delta_0} \right) \right) \right) + \text{Log} \left( \frac{6k_p^2}{\delta_0} \right) \right).$$

By simplifying the expression on the right hand side and noting that the constant  $c$  only appears in a logarithm, one can verify that for a sufficiently large choice of numerical constant  $c$  (e.g., any  $c \geq e^8$  would suffice), the right hand side is at most

$$\frac{ck_p^3}{m} \left( \text{dLog}(k_p) + \text{Log} \left( \frac{1}{\delta_0} \right) \right),$$

which extends the inductive hypothesis to  $m_0 = m$ . The result now follows by the principle of induction.  $\blacksquare$

### Appendix C. Proofs of lower bounds

We start with the analysis of the *coupon collector's problem* and corresponding lower bounds (see e.g., (Motwani and Raghavan, 2010)). Since we need a slightly more general result (not appearing in the standard textbooks to the best of our knowledge), we present a short proof for the sake of completeness. We remark that a similar argument was used in (Simon, 2015).

**Lemma 19 (Generalized coupon collector's problem)** *Let  $m \leq k \in \mathbb{N}$ . Consider a sequence  $x_1, x_2, \dots$  of independent uniform draws from a set of size  $k$ . Assume  $z \in \mathbb{N}$  satisfies that with probability at least  $1/2$ , the number of distinct elements among  $x_1, \dots, x_z$  is at least  $k - m$ . Then,*

$$z \geq k \left( \ln \frac{k}{m} - 1 - \sqrt{\frac{2}{m}} \right).$$

**Proof** Let  $Z$  denote the random variable that counts the number of independent draws until at least  $k - m$  distinct elements are present among  $x_1, x_2, \dots$ . We may write  $Z = \sum_{i=1}^{k-m} Z_i$ , where  $Z_i$  represents the (random) number of draws after  $i - 1$  distinct elements were observed and up to and including the first draw when  $i$  distinct elements have been observed. Observe that  $Z_i$ -s are independent, each having the geometric distribution with parameter  $p_i = \frac{k-i+1}{k}$ . Thus,  $\mathbb{E}Z_i = \frac{1}{p_i}$  and  $\text{Var}(Z_i) = \frac{1-p_i}{p_i^2}$ , which implies

$$\mathbb{E}Z = \sum_{i=1}^{k-m} \mathbb{E}Z_i = \sum_{i=1}^{k-m} \frac{k}{k-i+1} = k(H_k - H_m),$$

where  $H_p = \sum_{i=1}^p \frac{1}{i}$  stands for the  $p$ -th Harmonic number. Further, we have

$$\text{Var}(Z) = \sum_{i=1}^{k-m} \frac{k(i-1)}{(k-i+1)^2} = \sum_{j=m+1}^k \frac{k(k-j)}{j^2} \leq k^2 \sum_{j=m+1}^k \frac{1}{j^2} \leq \frac{k^2}{m}.$$

Finally, by Chebyshev's inequality and the relation  $\ln p \leq H_p \leq \ln p + 1$  we have, with probability at least  $\frac{1}{2}$ ,

$$Z \geq \mathbb{E}Z - k\sqrt{\frac{2}{m}} \geq k \left( \ln \frac{k}{m} - 1 - \sqrt{\frac{2}{m}} \right).$$

The claim follows.  $\blacksquare$

**Remark 20** *We will often use the following handy corollary of Lemma 19: under the conditions of this result,  $z \geq \frac{k}{2} \ln \frac{k}{m}$ , provided that  $1 + \sqrt{\frac{2}{m}} \leq \frac{1}{2} \ln \frac{k}{m}$ .*

**Proof of Theorem 10** A lower bound  $\Omega\left(\frac{d}{\varepsilon} + \frac{1}{\varepsilon} \log\left(\frac{1}{\delta}\right)\right)$  holds for all learning algorithms (Vapnik and Chervonenkis, 1974; Ehrenfeucht, Haussler, Kearns, and Valiant, 1989), so we focus only on establishing a lower bound  $\Omega\left(\frac{1}{\varepsilon} \log(k_o) \mathbb{1}[\varepsilon \leq 1/k_o]\right)$  for proper learners. Without loss of generality we may assume that  $k_o \geq 128$  since for smaller values of  $k_o$  the lower bound is automatically established by choosing a small enough numerical constant factor. We in fact establish a stronger result, which also implies the second claim: namely, for any  $k \geq 128$  such that there exists a hollow star of size  $k$ , any  $\varepsilon \leq 1/k$  has  $\mathcal{M}_{\text{prop}}(\varepsilon, \delta) \geq \frac{c}{\varepsilon} \log(k)$  (for a numerical constant  $c > 0$ ). Note that both of the claimed lower bounds will follow from this, since when  $k_o < \infty$  there exists a hollow star of size  $k_o$ , and when  $k_o = \infty$  there exists a sequence  $k_i \rightarrow \infty$  for which there exist hollow stars of each size  $k_i$ , so that choosing  $\varepsilon_i = 1/k_i$  the lower bound  $\frac{c}{\varepsilon_i} \log\left(\frac{1}{\varepsilon_i}\right)$  holds for each  $\varepsilon_i$ .

Fix any  $k \geq 128$  such that there exists a hollow star set  $S = \{(x_1, y_1), \dots, (x_k, y_k)\}$ , and for each  $i \in \{1, \dots, k\}$  let  $h_i \in \mathbb{C}$  be such that  $\{j : h_i(x_j) \neq y_j\} = \{i\}$ . Fix any proper learning algorithm  $\mathbb{A}'$ . We construct a target function  $f^*$  and distribution  $\mathcal{P}$  to witness the lower bound via the probabilistic method. Let  $\varepsilon \leq 1/k$  and choose  $i^* \sim \text{Uniform}(\{2, \dots, k\})$ , and set  $f^* = h_{i^*}$  and  $\mathcal{P}(\{x_i\}) = \varepsilon/(1-\varepsilon)$  for  $i \in \{2, \dots, k\} \setminus \{i^*\}$ ,  $\mathcal{P}(\{x_1\}) = 1 - (k-2)\varepsilon/(1-\varepsilon)$  (which is greater than  $\varepsilon$ ), and  $\mathcal{P}(\{x_{i^*}\}) = 0$ . Consider running  $\mathbb{A}'$  with a data set  $D_n$  (conditionally i.i.d. given  $i^*$ , with each point  $(X, Y)$  having  $X \sim \mathcal{P}$  and  $Y = f^*(X)$ ), of some size  $n < \frac{1}{8} \frac{1-\varepsilon}{\varepsilon} \ln(k-2)$ , and let  $\hat{h}$  be the classifier it outputs. Since  $\mathbb{A}'$  is proper ( $\hat{h} \in \mathbb{C}$ ) and  $\mathbb{C}[S] = \emptyset$  ( $S$  being a hollow star), we know that  $\hat{h}$  cannot realize the  $y_i$  classification of every  $x_i$ , so there must be a non-empty set  $\hat{I} = \{i : \hat{h}(x_i) \neq y_i\} \neq \emptyset$ . If any of these  $\hat{i} \in \hat{I}$  are not equal  $i^*$ , then  $\text{er}(\hat{h}) > \varepsilon$ . Denote by  $\hat{n}$  the number of the  $n$  data points in  $D_n$  falling in  $\{(x_2, y_2), \dots, (x_k, y_k)\} \setminus \{(x_{i^*}, y_{i^*})\}$ , and denote by  $\hat{n}_1$  the number of *distinct* elements of  $\{(x_2, y_2), \dots, (x_k, y_k)\} \setminus \{(x_{i^*}, y_{i^*})\}$  observed in the data set  $D_n$ . By Markov's inequality, with probability at least  $\frac{3}{4}$ , we have  $\hat{n} \leq 4 \frac{n(k-2)\varepsilon}{1-\varepsilon} < \frac{1}{2}(k-2) \ln(k-2)$ . Furthermore, note that conditioned on  $\hat{n}$  and  $i^*$ , the  $\hat{n}$  samples in  $D_n$  falling in  $\{(x_2, y_2), \dots, (x_k, y_k)\} \setminus \{(x_{i^*}, y_{i^*})\}$  are conditionally independent, with conditional distribution uniform on this set. Therefore, on the event that  $\hat{n} < \frac{1}{2}(k-2) \ln(k-2)$ , Lemma 19 implies that  $\Pr(\hat{n}_1 < k-3 | \hat{n}, i^*) > \frac{1}{2}$  (noting that  $k \geq 128$  implies  $1 + \sqrt{2} \leq \frac{1}{2} \ln(k-2)$ ). Finally, note that conditioned on  $D_n$ , the variable  $i^*$  has conditional distribution uniform on the  $k-1-\hat{n}_1$  values  $i \in \{2, \dots, k\}$  with  $(x_i, y_i) \notin D_n$ . Thus, on the event that  $\hat{n}_1 < k-3$ , we have that  $\Pr(\hat{I} \neq \{i^*\} | \hat{I}, D_n) \geq \frac{(k-1-\hat{n}_1)-1}{k-1-\hat{n}_1} \geq \frac{1}{2}$ . Altogether, we have

$$\begin{aligned} \Pr(\hat{I} \neq \{i^*\}) &\geq \mathbb{E}\left[\Pr(\hat{I} \neq \{i^*\} | \hat{I}, D_n) \mathbb{1}[\hat{n}_1 < k-3]\right] \geq \frac{1}{2} \Pr(\hat{n}_1 < k-3) \\ &\geq \frac{1}{2} \mathbb{E}\left[\Pr(\hat{n}_1 < k-3 | \hat{n}, i^*) \mathbb{1}[\hat{n} < \frac{1}{2}(k-2) \ln(k-2)]\right] \geq \frac{1}{4} \Pr(\hat{n} < \frac{1}{2}(k-2) \ln(k-2)) \geq \frac{3}{16}. \end{aligned}$$

In particular, this implies that for any proper learning algorithm  $\mathbb{A}'$ , if  $n < \frac{1}{8} \frac{1-\varepsilon}{\varepsilon} \ln(k-2)$ , there exist fixed choices of  $f^* \in \mathbb{C}$  and  $\mathcal{P}$  such that, with probability at least  $3/16$ , the classifier  $\hat{h}$  returned by  $\mathbb{A}'$  has  $\text{er}(\hat{h}) > \varepsilon$ . The claim follows.  $\blacksquare$

**Proof of Theorem 11** Fix any  $d, k_w$ . Since there is already a known lower bound  $\frac{c'}{\varepsilon} (d + \text{Log}(\frac{1}{\delta}))$  from (Ehrenfeucht, Haussler, Kearns, and Valiant, 1989; Blumer, Ehrenfeucht, Haussler, and Warmuth, 1989; Vapnik and Chervonenkis, 1974), for PAC learning in general (for some numerical constant  $c' > 0$ ), we focus on showing a lower bound  $\frac{cd}{\varepsilon} \text{Log}(\frac{k_w}{d} \wedge \frac{1}{\varepsilon})$  for some numerical constant  $c > 0$ . Furthermore, since  $\text{Log}(x) \geq 1$  (from its definition above), this lower bound again already

follows from the lower bound of [Ehrenfeucht, Haussler, Kearns, and Valiant \(1989\)](#) in the case that  $k_w < 126d$  (for instance, taking  $\mathcal{X} = \{1, \dots, d-1+k_w\}$  and  $\mathbb{C} = \{x \mapsto 2\mathbb{1}[x \in I] - 1 : I \subseteq \mathcal{X}, |I \cap \{1, \dots, k_w\}| = 1\}$ , which one can easily verify has VC dimension  $d$  and dual Helly number  $k_w$ ). To address the remaining case, suppose  $k_w \geq 126d$ .

In the special case  $d = 1$ , simply take  $\mathcal{X} = \{1, \dots, k_w\}$  and  $\mathbb{C} = \{x \mapsto 2\mathbb{1}[x = t] - 1 : t \in \mathcal{X}\}$  the singleton classifiers. It is an easy exercise to verify that the VC dimension of  $\mathbb{C}$  is indeed 1, and that  $\{(x, -1) : x \in \mathcal{X}\}$  is a hollow star set of size  $k_w$ , so that [Lemma 5](#) (together with the fact that this is clearly the largest possible hollow star set, and that  $\mathbb{C}$  is finite and therefore closed) implies the dual Helly number is indeed  $k_w$ . The claimed lower bound for this case then follows from [Theorem 10](#). To address the remaining case, for the rest of the proof suppose  $d \geq 2$ .

Let  $\mathcal{X} = \{(i, j) : i \in \{d-1, \dots, k_w+d-2\}, j \in \{1, \dots, i\}\}$ . For each  $i \in \{d-1, \dots, k_w+d-2\}$  and each  $J \subseteq \{1, \dots, i\}$  with  $|J| = d-1$ , define a classifier  $h_{i,J}(i', j) = 1 - 2\mathbb{1}[i' = i]\mathbb{1}[j \notin J]$ : that is,  $h_{i,J}$  classifies as 1 everything that does not have first coordinate equal  $i$ , and exactly  $d-1$  of the points that do have first coordinate equal  $i$ . Set

$$\mathbb{C} = \{h_{i,J} : i \in \{d-1, \dots, k_w+d-2\}, J \subseteq \{1, \dots, i\}, |J| = d-1\}.$$

We first show that the VC dimension of  $\mathbb{C}$  is indeed  $d$ . To see that  $d$  points can be shattered, simply take  $i = 2d-1$  (which has  $i \leq k_w+d-2$  since  $k_w \geq d+1$ ) and we claim that the points  $(i, 1), \dots, (i, d)$  are shattered: for any strict subset  $J \subset \{1, \dots, d\}$ , we can realize a labeling with  $\{(i, j) : j \in J\}$  positive and the other  $d-|J|$  negative with  $h_{i,J \cup J'}$  where  $J'$  is any subset of  $\{d+1, \dots, 2d-1\}$  with  $|J'| = d-1-|J|$ ; also, we can realize the all-positive labeling of these  $d$  points with  $h_{d-1, \{1, \dots, d-1\}}$ . To show no set of  $d+1$  points can be shattered, note that if  $\{((i_1, j_1), -1), \dots, ((i_{d+1}, j_{d+1}), -1)\}$  is realizable, then  $i_1 = \dots = i_{d+1}$ ; but in this case,  $\{((i_1, j_1), 1), \dots, ((i_d, j_d), 1), ((i_{d+1}, j_{d+1}), -1)\}$  is not realizable, and hence no set of size  $d+1$  is shattered.

Next we argue that the dual Helly number of  $\mathbb{C}$  is indeed  $k_w$ . To see that it is at least  $k_w$ , note that  $\{((k_w+d-2, 1), -1), \dots, ((k_w+d-2, k_w), -1)\}$  is a hollow star set of size  $k_w$ , so that [Lemma 5](#) implies the dual Helly number is at least  $k_w$ . To see that it is also at most  $k_w$ , consider any unrealizable set  $S$ . If some  $x$  has  $\{(x, -1), (x, 1)\} \subseteq S$ , this is clearly an unrealizable subset of size  $2 \leq k_w$ . Otherwise if no such  $x$  exists, then note that since  $h_{d-1, \{1, \dots, d-1\}}$  is positive on all of  $\mathcal{X}$ , there must be some  $(i, j_-)$  with  $((i, j_-), -1) \in S$ . If there are in fact *two* points  $(i, j), (i', j')$  with  $i \neq i'$  and  $\{((i, j), -1), ((i', j'), -1)\} \subseteq S$ , then again this is an unrealizable subset of size  $2 \leq k_w$ . Otherwise, if every  $(i, j)$  with  $((i, j), -1) \in S$  has the *same*  $i$ , then it must be that either there exist  $j_1, \dots, j_d$  with  $\{((i, j_1), 1), \dots, ((i, j_d), 1)\} \subseteq S$ , in which case  $\{((i, j_1), 1), \dots, ((i, j_d), 1), ((i, j_-), -1)\}$  is an unrealizable subset of size  $d+1 \leq k_w$ , or else there exist  $j_1, \dots, j_{i-(d-2)}$  with  $\{((i, j_1), -1), \dots, ((i, j_{i-(d-2)}), -1)\} \subseteq S$ , in which case this is an unrealizable subset of size  $i-(d-2) \leq k_w$ . Since this covers all possible cases for the set  $S$ , we conclude that the dual Helly number is equal  $k_w$ .

Fix any  $\delta \in (0, 1/100)$ . For any  $\varepsilon \in (1/504, 1/8)$ , a lower bound  $\frac{cd}{\varepsilon} \text{Log}(\frac{k_w}{d} \wedge \frac{1}{\varepsilon})$  follows from the lower bound  $\frac{c'd}{\varepsilon}$  of [Ehrenfeucht, Haussler, Kearns, and Valiant \(1989\)](#) (for  $c$  a sufficiently small numerical constant). To address the remaining case, fix any  $\varepsilon \in (0, 1/504]$ . If  $\varepsilon \geq \frac{d-1}{4(k_w-1)}$ , let  $i_\varepsilon = \lfloor (d-1)/(4\varepsilon) \rfloor + d-1$ , and otherwise let  $i_\varepsilon = k_w + d - 2$ . We prove the lower bound via the probabilistic method. Let  $J^*$  be a subset of  $\{1, \dots, i_\varepsilon\}$  with  $|J^*| = d-1$  chosen uniformly at random (without replacement). Let  $\mathcal{P}(\{(i_\varepsilon, j)\}) = \frac{4\varepsilon}{d-1}$  for every  $j \in \{1, \dots, i_\varepsilon\} \setminus J^*$ , and let

$\mathcal{P}(\{(d-1, 1)\}) = 1 - (i_\varepsilon - (d-1)) \frac{4\varepsilon}{d-1}$ , and define the target concept  $f^* = h_{i_\varepsilon, J^*}$ : in particular,  $f^*((i_\varepsilon, j)) = 2\mathbb{1}[j \in J^*] - 1$ , and hence  $\mathcal{P}$  has zero mass on the set of all points  $(i_\varepsilon, j)$  where  $f^*((i_\varepsilon, j)) = 1$  and has mass  $\frac{4\varepsilon}{d-1}$  on every point  $(i_\varepsilon, j)$  where  $f^*((i_\varepsilon, j)) = -1$ ; any remaining probability mass is placed on  $(d-1, 1)$ , which is an uninformative point (since every  $h_{i, J}$  classifies it 1).

Fix any sample size  $n \in \mathbb{N}$  with  $n < \frac{d-1}{32\varepsilon} \ln\left(\frac{1}{d-1} \min\{\lfloor \frac{d-1}{4\varepsilon} \rfloor, k_w - 1\}\right)$ , fix any proper learning algorithm  $\mathbb{A}'$ , and let  $\hat{h}$  be the classifier returned by running  $\mathbb{A}'$  on a conditionally i.i.d. (given  $J^*$ ) training set  $D_n$  of size  $n$  (with each  $(X, Y) \in D_n$  having  $X \sim \mathcal{P}$  and  $Y = f^*(X)$  given  $J^*$ ). Let  $\mathcal{Z}_{i_\varepsilon} = \{((i_\varepsilon, j), -1) : j \in \{1, \dots, i_\varepsilon\} \setminus J^*\}$  and  $\hat{n} = |D_n \cap \mathcal{Z}_{i_\varepsilon}|$ , and note that we have  $\mathbb{E}[\hat{n}|J^*] = \frac{4\varepsilon}{d-1}(i_\varepsilon - (d-1))n$ . Thus, by a Chernoff bound and the law of total probability, with probability at least  $1/2$ , it holds that  $\hat{n} \leq 1 + 2\varepsilon \mathbb{E}[\hat{n}|J^*] = 1 + \frac{8\varepsilon}{d-1}(i_\varepsilon - (d-1))n$  (see [Motwani and Raghavan, 2010](#)). Combining this with the constraint on  $n$ , on this event we have  $\hat{n} < 1 + \frac{i_\varepsilon - (d-1)}{4} \ln\left(\frac{1}{d-1} \min\{\lfloor \frac{d-1}{4\varepsilon} \rfloor, k_w - 1\}\right) \leq \frac{i_\varepsilon - (d-1)}{2} \ln\left(\frac{i_\varepsilon - (d-1)}{d-1}\right)$  (using the fact that  $\varepsilon \leq 1/504$ ). Furthermore, note that the samples in  $D_n \cap \mathcal{Z}_{i_\varepsilon}$  are conditionally i.i.d.  $\text{Uniform}(\mathcal{Z}_{i_\varepsilon})$  given  $\hat{n}$ . Also note that the assumptions that  $k_w \geq 126d$  and  $\varepsilon \leq 1/504$  imply  $i_\varepsilon - (d-1) \geq 126(d-1)$ , so that  $1 + \sqrt{\frac{2}{d-1}} \leq \frac{1}{2} \ln \frac{i_\varepsilon - (d-1)}{d-1}$ . Therefore, denoting by  $\hat{n}_1$  the number of *distinct* elements of  $\mathcal{Z}_{i_\varepsilon}$  present in  $D_n$ , we have that, on the event that  $\hat{n} < \frac{i_\varepsilon - (d-1)}{2} \ln\left(\frac{i_\varepsilon - (d-1)}{d-1}\right)$ , Lemma 19 implies  $\Pr(\hat{n}_1 < i_\varepsilon - 2(d-1) | \hat{n}, J^*) > \frac{1}{2}$ .

Since  $\mathbb{A}'$  is a proper learning algorithm, it must be that  $\hat{h} = h_{\hat{i}, \hat{J}}$  for some  $\hat{i} \in \{d-1, \dots, k_w + d-2\}$  and  $\hat{J} \subseteq \{1, \dots, \hat{i}\}$  with  $|\hat{J}| = d-1$ . If  $\hat{i} \neq i_\varepsilon$ , then  $\text{er}(\hat{h}) = (i_\varepsilon - (d-1)) \frac{4\varepsilon}{d-1} > \varepsilon$ . Otherwise, suppose  $\hat{i} = i_\varepsilon$ . Then  $\text{er}(\hat{h}) = |\hat{J} \setminus J^*| \frac{4\varepsilon}{d-1}$ . Note that, conditioned on  $D_n$ , the variable  $J^*$  has conditional distribution uniform on the subsets of the (size  $i_\varepsilon - \hat{n}_1$ ) set  $\{j \in \{1, \dots, i_\varepsilon\} : ((i_\varepsilon, j), -1) \notin D_n\}$  of size  $d-1$ . In particular, on the event  $\hat{i} = i_\varepsilon$ , we have  $\mathbb{E}[|\hat{J} \setminus J^*| | D_n, \hat{J}, \hat{i}] \geq (d-1) \frac{i_\varepsilon - \hat{n}_1 - (d-1)}{i_\varepsilon - \hat{n}_1}$ . On the event that  $\hat{n}_1 < i_\varepsilon - 2(d-1)$ , this implies  $\mathbb{E}[|\hat{J} \setminus J^*| | D_n, \hat{J}, \hat{i}] > \frac{d-1}{2}$ . Therefore, a Chernoff bound (for sampling without replacement; see [Hoeffding, 1963](#)) implies that, on the events that  $\hat{n}_1 < i_\varepsilon - 2(d-1)$  and  $\hat{i} = i_\varepsilon$ , we have  $\Pr\left(|\hat{J} \setminus J^*| \leq \frac{d-1}{4} \mid D_n, \hat{J}, \hat{i}\right) \leq \exp\left\{-\frac{d-1}{16}\right\} \leq e^{-1/16} < 1 - \frac{1}{17}$ . In particular, note that if  $\hat{i} = i_\varepsilon$  and  $|\hat{J} \setminus J^*| > \frac{d-1}{4}$ , then  $\text{er}(\hat{h}) > \varepsilon$ . Altogether, we have that

$$\begin{aligned}
 \Pr(\text{er}(\hat{h}) > \varepsilon) &\geq \Pr(\hat{i} \neq i_\varepsilon) + \mathbb{E}\left[\Pr\left(|\hat{J} \setminus J^*| > \frac{d-1}{4} \mid D_n, \hat{J}, \hat{i}\right) \mathbb{1}[\hat{n}_1 < i_\varepsilon - 2(d-1)] \mathbb{1}[\hat{i} = i_\varepsilon]\right] \\
 &\geq \frac{1}{17} \mathbb{E}\left[\Pr\left(\hat{n}_1 < i_\varepsilon - 2(d-1) \mid \hat{n}, J^*\right) \mathbb{1}\left[\hat{n} < \frac{i_\varepsilon - (d-1)}{2} \ln\left(\frac{i_\varepsilon - (d-1)}{d-1}\right)\right]\right] \\
 &\geq \frac{1}{34} \Pr\left(\hat{n} < \frac{i_\varepsilon - (d-1)}{2} \ln\left(\frac{i_\varepsilon - (d-1)}{d-1}\right)\right) \geq \frac{1}{68} > \delta.
 \end{aligned}$$

In particular, this implies that there exists a non-random choice of  $f^* \in \mathbb{C}$  and  $\mathcal{P}$  such that, with probability strictly greater than  $\delta$ , it holds that  $\text{er}(\hat{h}) > \varepsilon$ . The claimed lower bound on  $\mathcal{M}_{\text{prop}}(\varepsilon, \delta)$  follows by simplifying the expression of the constraint on  $n$  above (which is lower-bounded by the expression in the theorem, for a sufficiently small choice of the numerical constant  $c$ ).

For the final claim in the theorem, it is clear that we can extend the above construction to an infinite space by allowing all  $i \in \mathbb{N}$  with  $i \geq d-1$ , in which case  $k_o = k_w = \infty$  (since

there exist hollow star sets of unbounded sizes, following the same argument given above), and the  $d\text{Log}\left(\frac{\kappa_w}{d} \wedge \frac{1}{\varepsilon}\right)$  term simplifies to  $d\text{Log}\left(\frac{1}{\varepsilon}\right)$ .  $\blacksquare$

## Appendix D. Proof of Theorem 15

The essence of the proof of this result is in fact very simple, relying only on one technical construction: a set system on the data indices. Specifically, for any  $m \in \mathbb{N}$ , consider a family  $\mathcal{I}_m$  of subsets of  $\{1, \dots, m\}$  satisfying the following two properties, for some  $T_m \in \{1, \dots, m\}$ :

- (i) each  $I \in \mathcal{I}_m$  has size  $|I| \leq m - T_m$ ,
- (ii) for every  $i_1, i_2, \dots, i_\ell \in \{1, \dots, m\}$  there exists  $I \in \mathcal{I}_m$  such that  $\{i_1, i_2, \dots, i_\ell\} \subseteq I$ .

Let  $(\kappa, \rho)$  be a stable compression scheme of size  $\ell$ . Fix any distribution  $\mathcal{P}$ , any  $f^* \in \mathbb{C}$ , and any  $\delta \in (0, 1)$ , and let  $S = (X_{1:m}, f^*(X_{1:m}))$  be such that  $X_{1:m} \sim \mathcal{P}^m$ . Given any family  $\mathcal{I}_m$  satisfying (i) and (ii), we will establish that with probability at least  $1 - \delta$ ,

$$\text{er}(\rho(\kappa(S))) \leq \frac{1}{T_m} \left( \ln(|\mathcal{I}_m|) + \ln\left(\frac{1}{\delta}\right) \right). \quad (17)$$

As a simple example of such a family  $\mathcal{I}_m$  that yields Theorem 15, consider any partition of  $\{1, \dots, m\}$  into disjoint blocks  $I_1, \dots, I_{2\ell}$ , each of size either  $\lceil m/(2\ell) \rceil$  or  $\lfloor m/(2\ell) \rfloor$ . Then we can define

$$\mathcal{I}_m = \left\{ \bigcup \{I_j : j \in \mathcal{J}\} : \mathcal{J} \subseteq \{1, \dots, 2\ell\}, |\mathcal{J}| = \ell \right\}. \quad (18)$$

This clearly satisfies the above properties, with  $T_m = \ell \lfloor m/(2\ell) \rfloor$ , and has size  $|\mathcal{I}_m| = \binom{2\ell}{\ell} < 4^\ell$ , and hence plugging this into (17) yields the bound stated in Theorem 15. We now finish the proof of Theorem 15 by establishing the bound (17).

Fix any family  $\mathcal{I}_m$  satisfying (i) and (ii) above. For the set  $S$  as introduced above, for any  $I \subseteq \{1, \dots, m\}$  define  $S_I = \{(X_i, f^*(X_i)) : i \in I\}$ . For any  $I \in \mathcal{I}_m$ , since  $S_{(1:m)\setminus I}$  is independent of  $S_I$ , and property (i) implies  $|S_{(1:m)\setminus I}| \geq T_m$ , we have

$$\Pr\left(\rho(\kappa(S_I)) \text{ is correct on } S_{(1:m)\setminus I} \text{ and } \text{er}(\rho(\kappa(S_I))) > \varepsilon\right) \leq (1 - \varepsilon)^{T_m}.$$

However, by property (ii) there must exist at least one  $I^* \in \mathcal{I}$  with  $\kappa(S) \subseteq S_{I^*}$ , which means  $\rho(\kappa(S_{I^*})) = \rho(\kappa(S))$  (by the stability property). Thus, since  $\rho(\kappa(S))$  is correct on all of  $S$  (because  $(\kappa, \rho)$  is a valid compression scheme), including  $S_{(1:m)\setminus I^*}$ , we have for this (data-dependent) choice of  $I^*$  that  $\rho(\kappa(S_{I^*}))$  is correct on  $S_{(1:m)\setminus I^*}$ . Therefore, by basic inequalities and a union bound,

$$\begin{aligned} \Pr(\text{er}(\rho(\kappa(S))) > \varepsilon) &= \Pr(\text{er}(\rho(\kappa(S_{I^*}))) > \varepsilon) \\ &\leq \Pr(\exists I \in \mathcal{I}_m : \rho(\kappa(S_I)) \text{ is correct on } S_{(1:m)\setminus I} \text{ and } \text{er}(\rho(\kappa(S_I))) > \varepsilon) \\ &\leq |\mathcal{I}_m| (1 - \varepsilon)^{T_m} \leq |\mathcal{I}_m| e^{-\varepsilon T_m}. \end{aligned} \quad (19)$$

In particular, for any  $\delta \in (0, 1)$ , choosing  $\varepsilon$  equal the expression on the right hand side of (17) makes the rightmost expression in (19) equal  $\delta$ , which therefore completes the proof of the abstract bound (17). The bound in Theorem 15 follows by plugging in the family  $\mathcal{I}_m$  from (18), which has  $|\mathcal{I}_m| < 4^\ell$  and  $T_m = \ell \lfloor m/(2\ell) \rfloor > (m - 2\ell)/2$ .  $\blacksquare$