
Multiclass Learning Approaches: A Theoretical Comparison with Implications

Amit Daniely
Department of Mathematics
The Hebrew University
Jerusalem, Israel

Sivan Sabato
Microsoft Research
1 Memorial Drive
Cambridge, MA 02142, USA

Shai Shalev-Shwartz
School of CS and Eng.
The Hebrew University
Jerusalem, Israel

Abstract

We theoretically analyze and compare the following five popular multiclass classification methods: One vs. All, All Pairs, Tree-based classifiers, Error Correcting Output Codes (ECOC) with randomly generated code matrices, and Multiclass SVM. In the first four methods, the classification is based on a reduction to binary classification. We consider the case where the binary classifier comes from a class of VC dimension d , and in particular from the class of halfspaces over \mathbb{R}^d . We analyze both the estimation error and the approximation error of these methods. Our analysis reveals interesting conclusions of practical relevance, regarding the success of the different approaches under various conditions. Our proof technique employs tools from VC theory to analyze the *approximation error* of hypothesis classes. This is in contrast to most previous uses of VC theory, which only deal with estimation error.

1 Introduction

In this work we consider multiclass prediction: The problem of classifying objects into one of several possible target classes. Applications include, for example, categorizing documents according to topic, and determining which object appears in a given image. We assume that objects (a.k.a. instances) are vectors in $\mathcal{X} = \mathbb{R}^d$ and the class labels come from the set $\mathcal{Y} = [k] = \{1, \dots, k\}$. Following the standard PAC model, the learner receives a training set of m examples, drawn i.i.d. from some unknown distribution, and should output a classifier which maps \mathcal{X} to \mathcal{Y} .

The centrality of the multiclass learning problem has spurred the development of various approaches for tackling the task. Perhaps the most straightforward approach is a reduction from multiclass classification to binary classification. For example, the One-vs-All (OvA) method is based on a reduction of the multiclass problem into k binary problems, each of which discriminates between one class to all the rest of the classes (e.g. [Rumelhart et al. \[1986\]](#)). A different reduction is the All-Pairs (AP) approach in which all pairs of classes are compared to each other [[Hastie and Tibshirani, 1998](#)]. These two approaches have been unified under the framework of Error Correction Output Codes (ECOC) [[Dietterich and Bakiri, 1995](#), [Allwein et al., 2000](#)]. A tree-based classifier (TC) is another reduction in which the prediction is obtained by traversing a binary tree, where at each node of the tree a binary classifier is used to decide on the rest of the path (see for example [Beygelzimer et al. \[2007\]](#)).

All of the above methods are based on reductions to binary classification. We pay special attention to the case where the underlying binary classifiers are linear separators (halfspaces). Formally, each $w \in \mathbb{R}^{d+1}$ defines the linear separator $h_w(x) = \text{sign}(\langle w, \bar{x} \rangle)$, where $\bar{x} = (x, 1) \in \mathbb{R}^{d+1}$ is the concatenation of the vector x and the scalar 1. While halfspaces are our primary focus, many of our results hold for any underlying binary hypothesis class of VC dimension $d + 1$.

Other, more direct approaches to multiclass classification over \mathbb{R}^d have also been proposed (e.g. Vapnik [1998], Weston and Watkins [1999], Crammer and Singer [2001]). In this paper we analyze the Multiclass SVM (MSVM) formulation of Crammer and Singer [2001], in which each hypothesis is of the form $h_W(x) = \operatorname{argmax}_{i \in [k]} (W\bar{x})_i$, where W is a $k \times (d+1)$ matrix and $(W\bar{x})_i$ is the i 'th element of the vector $W\bar{x} \in \mathbb{R}^k$.

We theoretically analyze the prediction performance of the aforementioned methods, namely, OvA, AP, ECOC, TC, and MSVM. The error of a multiclass predictor $h : \mathbb{R}^d \rightarrow [k]$ is defined to be the probability that $h(x) \neq y$, where (x, y) is sampled from the underlying distribution \mathcal{D} over $\mathbb{R}^d \times [k]$, namely, $\operatorname{Err}(h) = \mathbb{P}_{(x,y) \sim \mathcal{D}}[h(x) \neq y]$. Our main goal is to understand which method is preferable in terms of the error it will achieve, based on easy-to-verify properties of the problem at hand.

Our analysis pertains to the type of classifiers each method can potentially find, and does not depend on the specific training algorithm. More precisely, each method corresponds to a hypothesis class, \mathcal{H} , which contains the multiclass predictors that may be returned by the method. For example, the hypothesis class of MSVM is $\mathcal{H} = \{x \mapsto \operatorname{argmax}_{i \in [k]} (W\bar{x})_i : W \in \mathbb{R}^{k \times (d+1)}\}$.

A learning algorithm, A , receives a training set, $S = \{(x_i, y_i)\}_{i=1}^m$, sampled i.i.d. according to \mathcal{D} , and returns a multiclass predictor which we denote by $A(S) \in \mathcal{H}$. A learning algorithm is called an Empirical Risk Minimizer (ERM) if it returns a hypothesis in \mathcal{H} that minimizes the empirical error on the sample. We denote by h^* a hypothesis in \mathcal{H} with minimal error,¹ that is, $h^* \in \operatorname{argmin}_{h \in \mathcal{H}} \operatorname{Err}(h)$.

When analyzing the error of $A(S)$, it is convenient to decompose this error as a sum of *approximation error* and *estimation error*:

$$\operatorname{Err}(A(S)) = \underbrace{\operatorname{Err}(h^*)}_{\text{approximation}} + \underbrace{\operatorname{Err}(A(S)) - \operatorname{Err}(h^*)}_{\text{estimation}}. \quad (1)$$

- The **approximation error** is the minimum error achievable by a predictor in the hypothesis class, \mathcal{H} . The approximation error does not depend on the sample size, and is determined solely by the allowed hypothesis class².
- The **estimation error** of an algorithm is the difference between the approximation error, and the error of the classifier the algorithm chose based on the sample. This error exists both for statistical reasons, since the sample may not be large enough to determine the best hypothesis, and for algorithmic reasons, since the learning algorithm may not output the best possible hypothesis given the sample. For the ERM algorithm, the estimation error can be bounded from above by order of $\sqrt{C(\mathcal{H})/m}$ where $C(\mathcal{H})$ is a complexity measure of \mathcal{H} (analogous to the VC dimension) and m is the sample size. A similar term also bounds the estimation error from below *for any algorithm*. Thus $C(\mathcal{H})$ is an estimate of the best achievable estimation error for the class.

When studying the estimation error of different methods, we follow the standard distribution-free analysis. Namely, we will compare the algorithms based on the worst-case estimation error, where worst-case is over all possible distributions \mathcal{D} . Such an analysis can lead us to the following type of conclusion: If two hypothesis classes have roughly the same complexity, $C(\mathcal{H}_1) \approx C(\mathcal{H}_2)$, and the number of available training examples is significantly larger than this value of complexity, then for both hypothesis classes we are going to have a small estimation error. Hence, in this case the difference in prediction performance between the two methods will be dominated by the approximation error and by the success of the learning algorithm in approaching the best possible estimation error. In our discussion below we disregard possible differences in optimality which stem from algorithmic aspects and implementation details. A rigorous comparison of training heuristics would certainly be of interest and is left to future work.

For the approximation error we will provide even stronger results, by comparing the approximation error of classes for *any* distribution. We rely on the following definition.

¹For simplicity, we assume that the minimum is attainable.

²Note that, when comparing different hypothesis classes over the same distribution, the Bayes error is constant. Thus, in the definition of approximation error, we do not subtract the Bayes error.

Definition 1.1. Given two hypothesis classes, $\mathcal{H}, \mathcal{H}'$, we say that \mathcal{H} essentially contains \mathcal{H}' if for any distribution, the approximation error of \mathcal{H} is at most the approximation error of \mathcal{H}' . \mathcal{H} strictly contains \mathcal{H}' if, in addition, there is a distribution for which the approximation error of \mathcal{H} is strictly smaller than that of \mathcal{H}' .

Our main findings are as follows (see a full comparison in Table 1). The formal statements are given in Section 3.

- The estimation errors of OvA, MSVM, and TC are all roughly the same, in the sense that $C(\mathcal{H}) = \tilde{\Theta}(dk)$ for all of the corresponding hypothesis classes. The complexity of AP is $\tilde{\Theta}(dk^2)$. The complexity of ECOC with a code of length l and code-distance δ is at most $\tilde{O}(dl)$ and at least $d\delta/2$. It follows that for randomly generated codes, $C(\mathcal{H}) = \tilde{\Theta}(dl)$. Note that this analysis shows that a larger code-distance yields a larger estimation error and might therefore hurt performance. This contrasts with previous “reduction-based” analyses of ECOC, which concluded that a larger code distance improves performance.
- We prove that the hypothesis class of MSVM essentially contains the hypothesis classes of both OvA and TC. Moreover, these inclusions are strict. Since the estimation errors of these three methods are roughly the same, it follows that the MSVM method dominates both OvA and TC in terms of achievable prediction performance.
- In the TC method, one needs to associate each leaf of the tree to a label. If no prior knowledge on how to break the symmetry is known, it is suggested in [Beygelzimer et al. \[2007\]](#) to break symmetry by choosing a random permutation of the labels. We show that whenever $d \ll k$, for any distribution \mathcal{D} , with high probability over the choice of a random permutation, the approximation error of the resulting tree would be close to $1/2$. It follows that a random choice of a permutation is likely to yield a poor predictor.
- We show that if $d \ll k$, for any distribution \mathcal{D} , the approximation error of ECOC with a randomly generated code matrix is likely to be close to $1/2$.
- We show that the hypothesis class of AP essentially contains the hypothesis class of MSVM (hence also that of OvA and TC), and that there can be a substantial gap in the containment. Therefore, as expected, the relative performance of AP and MSVM depends on the well-known trade-off between estimation error and approximation error.

	TC	OvA	MSVM	AP	random ECOC
Estimation error	dk	dk	dk	dk^2	dl
Approximation error	\geq MSVM $\approx 1/2$ when $d \ll k$	\geq MSVM	\geq AP	smallest	incomparable $\approx 1/2$ when $d \ll k$
Testing run-time	$d \log(k)$	dk	dk	dk^2	dl

Table 1: Summary of comparison

The above findings suggest that in terms of performance, it may be wiser to choose MSVM over OvA and TC, and especially so when $d \ll k$. We note, however, that in some situations (e.g. $d = k$) the prediction success of these methods can be similar, while TC has the advantage of having a testing run-time of $d \log(k)$, compared to the testing run-time of dk for OvA and MSVM. In addition, TC and ECOC may be a good choice when there is additional prior knowledge on the distribution or on how to break symmetry between the different labels.

1.1 Related work

[Allwein et al. \[2000\]](#) analyzed the multiclass error of ECOC as a function of the binary error. The problem with such a “reduction-based” analysis is that such analysis becomes problematic if the underlying binary problems are very hard. Indeed, our analysis reveals that the underlying binary problems would be too hard if $d \ll k$ and the code is randomly generated. The experiments in [Allwein et al. \[2000\]](#) show that when using kernel-based SVM or AdaBoost as the underlying classifier, OvA is inferior to random ECOC. However, in their experiments, the number of classes is small relative to the dimension of the feature space, especially if working with kernels or with combinations of weak learners.

Crammer and Singer [2001] presented experiments demonstrating that MSVM outperforms OvA on several data sets. Rifkin and Klautau [2004] criticized the experiments of Crammer and Singer [2001], Allwein et al. [2000], and presented another set of experiments demonstrating that all methods perform roughly the same when the underlying binary classifier is very strong (SVM with a Gaussian kernel). As our analysis shows, it is not surprising that with enough data and powerful binary classifiers, all methods should perform well. However, in many practical applications, we will prefer not to employ kernels (either because of shortage of examples, which might lead to a large estimation error, or due to computational constraint), and in such cases we expect to see a large difference between the methods.

Beygelzimer et al. [2007] analyzed the *regret* of a specific training method for trees, called Filter Tree, as a function of the regret of the binary classifier. The regret is defined to be the difference between the learned classifier and the Bayes-optimal classifier for the problem. Here again we show that the regret values of the underlying binary classifiers are likely to be very large whenever $d \ll k$ and the leaves of the tree are associated to labels in a random way. Thus in this case the regret analysis is problematic. Several authors presented ways to learn better splits, which corresponds to learning the association of leaves to labels (see for example Bengio et al. [2011] and the references therein). Some of our negative results do not hold for such methods, as these do not randomly attach labels to tree leaves.

Daniely et al. [2011] analyzed the properties of multiclass learning with various ERM learners, and have also provided some bounds on the estimation error of multiclass SVM and of trees. In this paper we both improve these bounds, derive new bounds for other classes, and also analyze the approximation error of the classes.

2 Definitions and Preliminaries

We first formally define the hypothesis classes that we analyze in this paper.

Multiclass SVM (MSVM): For $W \in \mathbb{R}^{k \times (d+1)}$ define $h_W : \mathbb{R}^d \rightarrow [k]$ by $h_W(x) = \operatorname{argmax}_{i \in [k]} (W\bar{x})_i$ and let $\mathcal{L} = \{h_W : W \in \mathbb{R}^{k \times (d+1)}\}$. Though NP-hard in general, solving the ERM problem with respect to \mathcal{L} can be done efficiently in the realizable case (namely, whenever exists a hypothesis with zero empirical error on the sample).

Tree-based classifiers (TC): A tree-based multiclass classifier is a full binary tree whose leaves are associated with class labels and whose internal nodes are associated with binary classifiers. To classify an instance, we start with the root node and apply the binary classifier associated with it. If the prediction is 1 we traverse to the right child. Otherwise, we traverse to the left child. This process continues until we reach a leaf, and then we output the label associated with the leaf. Formally, a tree for k classes is a full binary tree T together with a bijection $\lambda : \operatorname{leaf}(T) \rightarrow [k]$, which associates a label to each of the leaves. We usually identify T with the pair (T, λ) . The set of internal nodes of T is denoted by $N(T)$. Let $\mathcal{H} \subset \{\pm 1\}^{\mathcal{X}}$ be a binary hypothesis class. Given a mapping $C : N(T) \rightarrow \mathcal{H}$, define a multiclass predictor, $h_C : \mathcal{X} \rightarrow [k]$, by setting $h_C(x) = \lambda(v)$ where v is the last node of the root-to-leaf path $v_1, \dots, v_m = v$ such that v_{i+1} is the left (resp. right) child of v_i if $C(v_i)(x) = -1$ (resp. $C(v_i)(x) = 1$). Let $\mathcal{H}_T = \{h_C \mid C : N(T) \rightarrow \mathcal{H}\}$. Also, let $\mathcal{H}_{\text{trees}} = \cup_{T \text{ is a tree for } k \text{ classes}} \mathcal{H}_T$. If \mathcal{H} is the class of linear separators over \mathbb{R}^d , then for any tree T the ERM problem with respect to \mathcal{H}_T can be solved efficiently in the realizable case. However, the ERM problem is NP-hard in the non-realizable case.

Error Correcting Output Codes (ECOC): An ECOC is a code $M \in \mathbb{R}^{k \times l}$ along with a bijection $\lambda : [k] \rightarrow [k]$. We sometimes identify λ with the identity function and M with (M, λ) ³. Given a code M , and the result of l binary classifiers represented by a vector $u \in \{-1, 1\}^l$, the code selects a label via $\tilde{M} : \{-1, 1\}^l \rightarrow [k]$, defined by $\tilde{M}(u) = \lambda \left(\operatorname{argmax}_{i \in [k]} \sum_{j=1}^l M_{ij} u_j \right)$. Given binary classifiers h_1, \dots, h_l for each column in the code matrix, the code assigns to the instance $x \in \mathcal{X}$ the label $\tilde{M}(h_1(x), \dots, h_l(x))$. Let $\mathcal{H} \subset \{\pm 1\}^{\mathcal{X}}$ be a binary hypothesis class. Denote by

³The use of λ here allows us to later consider codes with random association of rows to labels.

$\mathcal{H}_M \subseteq [k]^{\mathcal{X}}$ the hypotheses class $\mathcal{H}_M = \{h : \mathcal{X} \rightarrow [k] \mid \exists (h_1, \dots, h_l) \in \mathcal{H}^l \text{ s.t. } \forall x \in \mathcal{X}, h(x) = \bar{M}(h_1(x), \dots, h_l(x))\}$.

The *distance* of a binary code, denoted by $\delta(M)$ for $M \in \{\pm 1\}^{k \times l}$, is the minimal *hamming distance* between any two pairs of rows in the code matrix. Formally, the hamming distance between $u, v \in \{-1, +1\}^l$ is $\Delta_h(u, v) = |\{r : u[r] \neq v[r]\}|$, and $\delta(M) = \min_{1 \leq i < j \leq k} \Delta_h(M[i], M[j])$. The ECOC paradigm described in [Dietterich and Bakiri, 1995] proposes to choose a code with a large distance.

One vs. All (OvA) and All Pairs (AP): Let $\mathcal{H} \subset \{\pm 1\}^{\mathcal{X}}$ and $k \geq 2$. In the OvA method we train k binary problems, each of which discriminates between one class and the rest of the classes. In the AP approach all pairs of classes are compared to each other. This is formally defined as two ECOCs. Define $M^{\text{OvA}} \in \mathbb{R}^{k \times k}$ to be the matrix whose (i, j) element is 1 if $i = j$ and -1 if $i \neq j$. Then, the hypothesis class of OvA is $\mathcal{H}_{\text{OvA}} = \mathcal{H}_{M^{\text{OvA}}}$. For the AP method, let $M^{\text{AP}} \in \mathbb{R}^{k \times \binom{k}{2}}$ be such that for all $i \in [k]$ and $1 \leq j < l \leq k$, the coordinate corresponding to row i and column (j, l) is defined to be -1 if $i = j$, 1 if $i = l$, and 0 otherwise. Then, the hypothesis class of AP is $\mathcal{H}_{\text{AP}} = \mathcal{H}_{M^{\text{AP}}}$.

Our analysis of the estimation error is based on results that bound the sample complexity of multi-class learning. The *sample complexity* of an algorithm A is the function m_A defined as follows: For $\epsilon, \delta > 0$, $m_A(\epsilon, \delta)$ is the smallest integer such that for every $m \geq m_A(\epsilon, \delta)$ and every distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, with probability of $> 1 - \delta$ over the choice of an i.i.d. sample S of size m ,

$$\text{Err}(A(S_m)) \leq \min_{h \in \mathcal{H}} \text{Err}(h) + \epsilon. \quad (2)$$

The first term on the right-hand side is the approximation error of \mathcal{H} . Therefore, the sample complexity is the number of examples required to ensure that the estimation error of A is at most ϵ (with high probability). We denote the sample complexity of a class \mathcal{H} by $m_{\mathcal{H}}(\epsilon, \delta) = \inf_A m_A(\epsilon, \delta)$, where the infimum is taken over all learning algorithms.

To bound the sample complexity of a hypothesis class we rely on upper and lower bounds on the sample complexity in terms of two generalizations of the VC dimension for multiclass problems, called the *Graph dimension* and the *Natarajan dimension* and denoted $d_G(\mathcal{H})$ and $d_N(\mathcal{H})$. For completeness, these dimensions are formally defined in the appendix.

Theorem 2.1. *Daniely et al. [2011]* For every hypothesis class \mathcal{H} , and for every ERM rule,

$$\Omega\left(\frac{d_N(\mathcal{H}) + \ln(\frac{1}{\delta})}{\epsilon^2}\right) \leq m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\text{ERM}}(\epsilon, \delta) \leq O\left(\frac{\min\{d_N(\mathcal{H}) \ln(|\mathcal{Y}|), d_G(\mathcal{H})\} + \ln(\frac{1}{\delta})}{\epsilon^2}\right)$$

We note that the constants in the O, Ω notations are universal.

3 Main Results

In Section 3.1 we analyze the sample complexity of the different hypothesis classes. We provide lower bounds on the Natarajan dimensions of the various hypothesis classes, thus concluding, in light of Theorem 2.1, a lower bound on the sample complexity of *any* algorithm. We also provide upper bounds on the graph dimensions of these hypothesis classes, yielding, by the same theorem, an upper bound on the estimation error of ERM. In Section 3.2 we analyze the approximation error of the different hypothesis classes.

3.1 Sample Complexity

Together with Theorem 2.1, the following theorems estimate, up to logarithmic factors, the sample complexity of the classes under consideration. We note that these theorems support the rule of thumb that the Natarajan and Graph dimensions are of the same order of the number of parameters. The first theorem shows that the sample complexity of MSVM depends on $\tilde{\Theta}(dk)$.

Theorem 3.1. $d(k-1) \leq d_N(\mathcal{L}) \leq d_G(\mathcal{L}) \leq O(dk \log(dk))$.

Next, we analyze the sample complexities of TC and ECOC. These methods rely on an underlying hypothesis class of binary classifiers. While our main focus is the case in which the binary hypothesis class is halfspaces over \mathbb{R}^d , the upper bounds on the sample complexity we derive below holds for any binary hypothesis class of VC dimension $d+1$.

Theorem 3.2. *For every binary hypothesis class of VC dimension $d + 1$, and for any tree T , $d_G(\mathcal{H}_T) \leq d_G(\mathcal{H}_{\text{trees}}) \leq O(dk \log(dk))$. If the underlying hypothesis class is halfspaces over \mathbb{R}^d , then also*

$$d(k-1) \leq d_N(\mathcal{H}_T) \leq d_G(\mathcal{H}_T) \leq d_G(\mathcal{H}_{\text{trees}}) \leq O(dk \log(dk)).$$

Theorems 3.1 and 3.2 improve results from Daniely et al. [2011] where it was shown that $\lfloor \frac{d}{2} \rfloor \lfloor \frac{k}{2} \rfloor \leq d_N(\mathcal{L}) \leq O(dk \log(dk))$, and for every tree $d_G(\mathcal{H}_T) \leq O(dk \log(dk))$. Further it was shown that if \mathcal{H} is the set of halfspaces over \mathbb{R}^d , then $\Omega\left(\frac{dk}{\log(k)}\right) \leq d_N(\mathcal{H}_T)$.

We next turn to results for ECOC, and its special cases OvA and AP.

Theorem 3.3. *For every $M \in \mathbb{R}^{k \times l}$ and every binary hypothesis class of VC dimension d , $d_G(\mathcal{H}_M) \leq O(dl \log(dl))$. Moreover, if $M \in \{\pm 1\}^{k \times l}$ and the underlying hypothesis class is halfspaces over \mathbb{R}^d , then*

$$d \cdot \delta(M)/2 \leq d_N(\mathcal{H}_M) \leq d_G(\mathcal{H}_M) \leq O(dl \log(dl)).$$

We note if the code has a large distance, which is the case, for instance, in random codes, then $\delta(M) = \Omega(l)$. In this case, the bound is tight up to logarithmic factors.

Theorem 3.4. *For any binary hypothesis class of VC dimension d , $d_G(\mathcal{H}_{\text{OvA}}) \leq O(dk \log(dk))$ and $d_G(\mathcal{H}_{\text{AP}}) \leq O(dk^2 \log(dk))$. If the underlying hypothesis class is halfspaces over \mathbb{R}^d we also have:*

$$\begin{aligned} d(k-1) &\leq d_N(\mathcal{H}_{\text{OvA}}) \leq d_G(\mathcal{H}_{\text{OvA}}) \leq O(dk \log(dk)) \quad \text{and} \\ d\binom{k-1}{2} &\leq d_N(\mathcal{H}_{\text{AP}}) \leq d_G(\mathcal{H}_{\text{AP}}) \leq O(dk^2 \log(dk)). \end{aligned}$$

3.2 Approximation error

We first show that the class \mathcal{L} essentially contains \mathcal{H}_{OvA} and \mathcal{H}_T for any tree T , assuming, of course, that \mathcal{H} is the class of halfspaces in \mathbb{R}^d . We find this result quite surprising, since the sample complexity of all of these classes is of the same order.

Theorem 3.5. *\mathcal{L} essentially contains $\mathcal{H}_{\text{trees}}$ and \mathcal{H}_{OvA} . These inclusions are strict for $d \geq 2$ and $k \geq 3$.*

One might suggest that a small increase in the dimension would perhaps allow us to embed \mathcal{L} in \mathcal{H}_T for some tree T or for OvA. The next result shows that this is not the case.

Theorem 3.6. *Any embedding into a higher dimension that allows \mathcal{H}_{OvA} or \mathcal{H}_T (for some tree T for k classes) to essentially contain \mathcal{L} , necessarily embeds into a dimension of at least $\tilde{\Omega}(dk)$.*

The next theorem shows that the approximation error of AP is better than that of MSVM (and hence also better than OvA and TC). This is expected as the sample complexity of AP is considerably higher, and therefore we face the usual trade-off between approximation and estimation error.

Theorem 3.7. *\mathcal{H}_{AP} essentially contains \mathcal{L} . Moreover, there is a constant $k^* > 0$, independent of d , such that the inclusion is strict for all $k \geq k^*$.*

For a random ECOC of length $o(k)$, it is easy to see that it does not contain MSVM, as MSVM has higher complexity. It is also not contained in MSVM, as it generates non-convex regions of labels.

We next derive absolute lower bounds on the approximation errors of ECOC and TC when $d \ll k$. Recall that both methods are built upon binary classifiers that should predict $h(x) = 1$ if the label of x is in L , for some $L \subset [k]$, and should predict $h(x) = -1$ if the label of x is not in L . As the following lemma shows, when the partition of $[k]$ into the two sets L and $[k] \setminus L$ is arbitrary and balanced, and $k \gg d$, such binary classifiers will almost always perform very poorly.

Lemma 3.8. *There exists a constant $C > 0$ for which the following holds. Let $\mathcal{H} \subseteq \{\pm 1\}^{\mathcal{X}}$ be any hypothesis class of VC-dimension d , let $\mu \in (0, 1/2]$, and let \mathcal{D} be any distribution over $\mathcal{X} \times [k]$ such that $\forall i \mathbb{P}_{(x,y) \sim \mathcal{D}}(y = i) \leq \frac{10}{k}$. Let $\phi : [k] \rightarrow \{\pm 1\}$ be a randomly chosen function which is sampled according to one of the following rules: (1) For each $i \in [k]$, each coordinate $\phi(i)$ is chosen independently from the other coordinates and $\mathbb{P}(\phi(i) = -1) = \mu$; or (2) ϕ is chosen uniformly among all functions satisfying $|\{i \in [k] : \phi(i) = -1\}| = \mu k$.*

Let \mathcal{D}_ϕ be the distribution over $\mathcal{X} \times \{\pm 1\}$ obtained by drawing (x, y) according to \mathcal{D} and replacing it with $(x, \phi(y))$. Then, for any $\nu > 0$, if $k \geq C \cdot \left(\frac{d + \ln(\frac{1}{\delta})}{\nu^2}\right)$, then with probability of at least $1 - \delta$ over the choice of ϕ , the approximation error of \mathcal{H} with respect to \mathcal{D}_ϕ will be at least $\mu - \nu$.

As the corollaries below show, Lemma 3.8 entails that when $k \gg d$, both random ECOCs with a small code length, and balanced trees with a random labeling of the leaves, are expected to perform very poorly.

Corollary 3.9. *There is a constant $C > 0$ for which the following holds. Let (T, λ) be a tree for k classes such that $\lambda : \text{leaf}(T) \rightarrow [k]$ is chosen uniformly at random. Denote by k_L and k_R the number of leaves of the left and right sub-trees (respectively) that descend from root, and let $\mu = \min\{\frac{k_1}{k}, \frac{k_2}{k}\}$. Let $\mathcal{H} \subseteq \{\pm 1\}^{\mathcal{X}}$ be a hypothesis class of VC-dimension d , let $\nu > 0$, and let \mathcal{D} be any distribution over $\mathcal{X} \times [k]$ such that $\forall i \mathbb{P}_{(x,y) \sim \mathcal{D}}(y = i) \leq \frac{10}{k}$. Then, for $k \geq C \cdot \left(\frac{d + \ln(\frac{1}{\delta})}{\nu^2}\right)$, with probability of at least $1 - \delta$ over the choice of λ , the approximation error of \mathcal{H}_T with respect to \mathcal{D} is at least $\mu - \nu$.*

Corollary 3.10. *There is a constant $C > 0$ for which the following holds. Let (M, λ) be an ECOC where $M \in \mathbb{R}^{k \times l}$, and assume that the bijection $\lambda : [k] \rightarrow [k]$ is chosen uniformly at random. Let $\mathcal{H} \subseteq \{\pm 1\}^{\mathcal{X}}$ be a hypothesis class of VC-dimension d , let $\nu > 0$, and let \mathcal{D} be any distribution over $\mathcal{X} \times [k]$ such that $\forall i \mathbb{P}_{(x,y) \sim \mathcal{D}}(y = i) \leq \frac{10}{k}$. Then, for $k \geq C \cdot \left(\frac{dl \log(dl) + \ln(\frac{1}{\delta})}{\nu^2}\right)$, with probability of at least $1 - \delta$ over the choice of λ , the approximation error of \mathcal{H}_M with respect to \mathcal{D} is at least $1/2 - \nu$.*

Note that the first corollary holds even if only the top level of the binary tree is balanced and splits the labels randomly to the left and the right sub-trees. The second corollary holds even if the code itself is not random (nor does it have to be binary), and only the association of rows with labels is random. In particular, if the length of the code is $O(\log(k))$, as suggested in [Allwein et al. \[2000\]](#), and the number of classes is $\tilde{\Omega}(d)$, then the code is expected to perform poorly.

For an ECOC with a matrix of length $\Omega(k)$ and $d = o(k)$, we do not have such a negative result as stated in Corollary 3.10. Nonetheless, Lemma 3.8 implies that the prediction of the binary classifiers when $d = o(k)$ is just slightly better than a random guess, thus it seems to indicate that the ECOC method will still perform poorly. Moreover, most current theoretical analyses of ECOC estimate the error of the learned multiclass hypothesis in terms of the average error of the binary classifiers. Alas, when the number of classes is large, Lemma 3.8 shows that this average will be close to $\frac{1}{2}$.

Finally, let us briefly discuss the tightness of Lemma 3.8. Let $x_1, \dots, x_{d+1} \in \mathbb{R}^d$ be affinely independent and let \mathcal{D} be the distribution over $\mathbb{R}^d \times [d+1]$ defined by $\mathbb{P}_{(x,y) \sim \mathcal{D}}((x, y) = (x_i, i)) = \frac{1}{d+1}$. It is not hard to see that for every $\phi : [d+1] \rightarrow \{\pm 1\}$, the approximation error of the class of half-spaces with respect to \mathcal{D}_ϕ is zero. Thus, in order to ensure a large approximation error *for every distribution*, the number of classes must be at least linear in the dimension, so in this sense, the lemma is tight. Yet, this example is very simple, since each class is concentrated on a single point and the points are linearly independent. It is possible that in real-world distributions, a large approximation error will be exhibited even when $k < d$.

We note that the phenomenon of a large approximation error, described in Corollaries 3.9 and 3.10, does not reproduce in the classes \mathcal{L} , \mathcal{H}_{OVA} and \mathcal{H}_{AP} , since these classes are symmetric.

4 Proof Techniques

Due to lack of space, the proofs for all the results stated above are provided in the appendix. In this section we give a brief description of our main proof techniques.

Most of our proofs for the estimation error results, stated in Section 3.1, are based on a similar method which we now describe. Let $L : \{\pm 1\}^l \rightarrow [k]$ be a multiclass-to-binary reduction (e.g., a tree), and for $\mathcal{H} \subseteq \{\pm 1\}^{\mathcal{X}}$, denote $L(\mathcal{H}) = \{x \mapsto L(h_1(x), \dots, h_l(x)) \mid h_1, \dots, h_l \in \mathcal{H}\}$. Our upper bounds for $d_G(L(\mathcal{H}))$ are mostly based on the following simple lemma.

Lemma 4.1. *If $\text{VC}(\mathcal{H}) = d$ then $d_G(L(\mathcal{H})) = O(ld \ln(ld))$.*

The technique for the lower bound on $d_N(L(\mathcal{W}))$ when \mathcal{W} is the class of halfspaces in \mathbb{R}^d is more involved, and quite general. We consider a binary hypothesis class $\mathcal{G} \subseteq \{\pm 1\}^{[d] \times [l]}$ which consists of functions having an arbitrary behaviour over $[d] \times \{i\}$, and a very uniform behaviour on other inputs (such as mapping all other inputs to a constant). We show that $L(\mathcal{G})$ N -shatters the set $[d] \times [l]$. Since \mathcal{G} is quite simple, this is usually not very hard to show. Finally, we show that the class of halfspaces is richer than \mathcal{G} , in the sense that the inputs to \mathcal{G} can be mapped to points in \mathbb{R}^d such that the functions of \mathcal{G} can be mapped to halfspaces. We conclude that $d_N(L(\mathcal{W})) \geq d_N(L(\mathcal{G}))$.

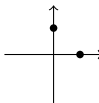
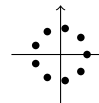
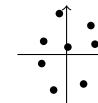
To prove the approximation error lower bounds stated in Section 3.2, we use the techniques of VC theory in an unconventional way. The idea of this proof is as follows: Using a uniform convergence argument based on the VC dimension of the binary hypothesis class, we show that there exists a small labeled sample S whose approximation error for the hypothesis class is close to the approximation error for the distribution, for all possible label mappings. This allows us to restrict our attention to a finite set of hypotheses, by their restriction to the sample. For these hypotheses, we show that with high probability over the choice of label mapping, the approximation error on the sample is high. A union bound on the finite set of possible hypotheses shows that the approximation error on the distribution will be high, with high probability over the choice of the label mapping.

5 Implications

The first immediate implication of our results is that whenever the number of examples in the training set is $\tilde{\Omega}(dk)$, MSVM should be preferred to OvA and TC. This is certainly true if the hypothesis class of MSVM, \mathcal{L} , has a zero approximation error (the realizable case), since the ERM is then solvable with respect to \mathcal{L} . Note that since the inclusions given in Theorem 3.5 are strict, there are cases where the data is realizable with MSVM but not with \mathcal{H}_{OvA} or with respect to any tree.

In the non-realizable case, implementing the ERM is intractable for all of these methods. Nonetheless, for each method there are reasonable heuristics to approximate the ERM, which should work well when the approximation error is small. Therefore, we believe that MSVM should be the method of choice in this case as well due to its lower approximation error. However, variations in the optimality of algorithms for different hypothesis classes should also be taken into account in this analysis. We leave this detailed analysis of specific training heuristics for future work. Our analysis also implies that it is highly unrecommended to use TC with a randomly selected λ or ECOC with a random code whenever $k > d$. Finally, when the number of examples is much larger than dk^2 , the analysis implies that it is better to choose the AP approach.

To conclude this section, we illustrate the relative performance of MSVM, OvA, TC, and ECOC, by considering the simplistic case where $d = 2$, and each class is concentrated on a single point in \mathbb{R}^2 . In the leftmost graph below, there are two classes in \mathbb{R}^2 , and the approximation error of all algorithms is zero. In the middle graph, there are 9 classes ordered on the unit circle of \mathbb{R}^2 . Here, both MSVM and OvA have a zero approximation error, but the error of TC and of ECOC with a random code will most likely be large. In the rightmost graph, we chose random points in \mathbb{R}^2 . MSVM still has a zero approximation error. However, OvA cannot learn the binary problem of distinguishing between the middle point and the rest of the points and hence has a larger approximation error.

			
MSVM	✓	✓	✓
OvA	✓	✓	✗
TC/ECOC	✓	✗	✗

Acknowledgements: Shai Shalev-Shwartz was supported by the John S. Cohen Senior Lectureship in Computer Science. Amit Daniely is a recipient of the Google Europe Fellowship in Learning Theory, and this research is supported in part by this Google Fellowship.

References

- E. L. Allwein, R.E. Schapire, and Y. Singer. Reducing multiclass to binary: A unifying approach for margin classifiers. *Journal of Machine Learning Research*, 1:113–141, 2000.
- S. Ben-David, N. Cesa-Bianchi, D. Haussler, and P. Long. Characterizations of learnability for classes of $\{0, \dots, n\}$ -valued functions. *Journal of Computer and System Sciences*, 50:74–86, 1995.
- S. Bengio, J. Weston, and D. Grangier. Label embedding trees for large multi-class tasks. In *NIPS*, 2011.
- A. Beygelzimer, J. Langford, and P. Ravikumar. Multiclass classification with filter trees. *Preprint, June*, 2007.
- K. Crammer and Y. Singer. On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of Machine Learning Research*, 2:265–292, 2001.
- A. Daniely, S. Sabato, S. Ben-David, and S. Shalev-Shwartz. Multiclass learnability and the erm principle. In *COLT*, 2011.
- T. G. Dietterich and G. Bakiri. Solving multiclass learning problems via error-correcting output codes. *Journal of Artificial Intelligence Research*, 2:263–286, January 1995.
- Trevor Hastie and Robert Tibshirani. Classification by pairwise coupling. *The Annals of Statistics*, 26(1):451–471, 1998.
- Ryan Rifkin and Aldebaro Klautau. In defense of one-vs-all classification. *Journal of Machine Learning Research*, 5:101–141, 2004.
- David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning internal representations by error propagation. In David E. Rumelhart and James L. McClelland, editors, *Parallel Distributed Processing – Explorations in the Microstructure of Cognition*, chapter 8, pages 318–362. MIT Press, 1986.
- G. Takacs. *Convex polyhedron learning and its applications*. PhD thesis, Budapest University of Technology and Economics, 2009.
- V. N. Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- J. Weston and C. Watkins. Support vector machines for multi-class pattern recognition. In *Proceedings of the Seventh European Symposium on Artificial Neural Networks*, April 1999.

A Proofs

A.1 Notation and Definitions

Throughout the proofs, we fix $d, k \geq 2$. We denote by $\mathcal{W} = \mathcal{W}^d = \{h_w : w \in \mathbb{R}^{d+1}\}$ the class of linear separators (with bias) over \mathbb{R}^d . We assume the following "tie breaking" conventions:

- For $f : [k] \rightarrow \mathbb{R}$, $\operatorname{argmax}_{i \in [k]} f(i)$ is the *minimal* number $i_0 \in [k]$ for which $f(i_0) = \max_{i \in [k]} f(i)$;
- $\operatorname{sign}(0) = 1$.

Given a hypotheses class $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$, denote its restriction to $A \subseteq \mathcal{X}$ by $\mathcal{H}|_A = \{f|_A : f \in \mathcal{H}\}$. Let $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ be a hypothesis class and let $\phi : \mathcal{Y} \rightarrow \mathcal{Y}'$, $\iota : \mathcal{X} \rightarrow \mathcal{X}'$ be functions. Denote $\phi \circ \mathcal{H} = \{\phi \circ h : h \in \mathcal{H}\}$ and $\mathcal{H} \circ \iota = \{h \circ \iota : h \in \mathcal{H}\}$.

Given $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ and a distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, denote the approximation error by $\operatorname{Err}_{\mathcal{D}}^*(\mathcal{H}) = \inf_{h \in \mathcal{H}} \operatorname{Err}_{\mathcal{D}}(h)$. Recall that by definition 1.1, \mathcal{H} essentially contains $\mathcal{H}' \subseteq \mathcal{Y}^{\mathcal{X}}$ if and only if $\operatorname{Err}_{\mathcal{D}}^*(\mathcal{H}) \leq \operatorname{Err}_{\mathcal{D}}^*(\mathcal{H}')$ for every distribution \mathcal{D} . For a binary hypothesis class \mathcal{H} , denote its VC dimension by $\operatorname{VC}(\mathcal{H})$.

Let $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ be a hypothesis class and let $S \subseteq \mathcal{X}$. We say that \mathcal{H} *G-shatters* S if there exists an $f : S \rightarrow \mathcal{Y}$ such that for every $T \subseteq S$ there is a $g \in \mathcal{H}$ such that

$$\forall x \in T, g(x) = f(x), \text{ and } \forall x \in S \setminus T, g(x) \neq f(x).$$

We say that \mathcal{H} *N-shatters* S if there exist $f_1, f_2 : S \rightarrow \mathcal{Y}$ such that $\forall y \in S, f_1(y) \neq f_2(y)$, and for every $T \subseteq S$ there is a $g \in \mathcal{H}$ such that

$$\forall x \in T, g(x) = f_1(x), \text{ and } \forall x \in S \setminus T, g(x) = f_2(x).$$

The *graph dimension* of \mathcal{H} , denoted $d_G(\mathcal{H})$, is the maximal cardinality of a set that is G-shattered by \mathcal{H} . The *Natarajan dimension* of \mathcal{H} , denoted $d_N(\mathcal{H})$, is the maximal cardinality of a set that is N-shattered by \mathcal{H} . Both of these dimensions coincide with the VC-dimension for $|\mathcal{Y}| = 2$. Note also that we always have $d_N(\mathcal{H}) \leq d_G(\mathcal{H})$. As shown in Ben-David et al. [1995], it also holds that $d_G(\mathcal{H}) \leq 4.67 \log_2(|\mathcal{Y}|) d_N(\mathcal{H})$.

Proof of Lemma 4.1. Let $A \subseteq \mathcal{X}$ be a G-shattered set with $|A| = d_G(L(\mathcal{H}))$. By Sauer's Lemma, $2^{|A|} \leq |\mathcal{H}|_A|^l \leq |A|^{dl}$, thus $d_G(L(\mathcal{H})) = |A| = O(ld \log(ld))$. \square

A.2 Multiclass SVM

Proof of Theorem 3.1. The lower bound follows from Theorems 3.5 and 3.2. To upper bound $d_G := d_G(\mathcal{L})$, let $S = \{x_1, \dots, x_{d_G}\} \subseteq \mathbb{R}^d$ be a set which is G-shattered by \mathcal{L} , and let $f : S \rightarrow [k]$ be a function that witnesses the shattering. For $x \in \mathbb{R}^d$ and $j \in [k]$, denote

$$\phi(x, j) = (0, \dots, 0, x[1], \dots, x[d], 1, 0, \dots, 0) \in \mathbb{R}^{(d+1)k},$$

where $x[1]$ is in the $(d+1)(j-1)$ coordinate. For every $(i, j) \in [d_G] \times [k]$, define $z_{i,j} = \phi(x_i, f(x_i)) - \phi(x_i, j)$. Denote $Z = \{z_{i,j} \mid (i, j) \in [d_G] \times [k]\}$. Since $\operatorname{VC}(\mathcal{W}^{(d+1)k}) = (d+1)k+1$, by Sauer's lemma,

$$|\mathcal{W}^{(d+1)k}|_Z \leq |Z|^{(d+1)k+1} = (d_G k)^{(d+1)k+1}.$$

We now show that there is a one-to-one mapping from subsets of S to $\mathcal{W}^{(d+1)k}|_Z$, thus concluding an upper bound on the size of S . For any $T \subseteq S$, choose $W(T) \in \mathbb{R}^{k \times (d+1)}(\mathbb{R})$ such that

$$\{x \in S \mid h_{W(T)}(x) = f(x)\} = T.$$

Such a $W(T)$ exists because of the G-shattering of S by \mathcal{L} using the witness f . Define the vector $w(T) \in \mathbb{R}^{k(d+1)}$ which is the concatenation of the rows of $W(T)$, that is $w(T) = (W(T)_{(1,1)}, \dots, W(T)_{(1,d+1)}, \dots, W(T)_{(k,1)}, \dots, W(T)_{(k,d+1)})$.

Now, suppose that $T_1 \neq T_2$ for $T_1, T_2 \subseteq S$. We now show that $w(T_1)|_Z \neq w(T_2)|_Z$. Suppose w.l.o.g. that there is some $x_i \in T_1 \setminus T_2$. Thus, $f(x_i) = h_{W(T_1)}(x_i) \neq h_{W(T_2)}(x_i) =: j$. It

follows that the inner product of x_i with row $f(x_i)$ of $W(T_1)$ is greater than the inner product of x_i with row j of $W(T_1)$, while for $W(T_2)$, the situation is reversed. Therefore, $\text{sign}(\langle w(T_1), z_{i,j} \rangle) \neq \text{sign}(\langle w(T_2), z_{i,j} \rangle)$, so $w(T_1)$ and $w(T_2)$ induce different labelings of Z . It follows that the number of subsets of S is bounded by the size of $\mathcal{W}^{(d+1)k}|_Z$, thus $2^{d_G} \leq (kd_G)^{(d+1)k+1}$. We conclude that $d_G \leq O(dk \log(dk))$. \square

A.3 Simple classes that can be represented by the class of linear separators

In this section we define two fairly simple hypothesis classes, and show that the class of linear separators is richer than them. We will later use this observation to prove lower bounds on the Natarajan dimension of various multiclass hypothesis classes.

Let $l \geq 2$. For $f \in \{-1, 1\}^{[d]}$, $i \in [l]$, $j \in \{-1, 1\}$ define $f^{i,j} : [d] \times [l] \rightarrow \{-1, 1\}$ by

$$f^{i,j}(u, v) = \begin{cases} f(u) & v = i \\ j & v \neq i, \end{cases}$$

And define the hypothesis class \mathcal{F}^l as

$$\mathcal{F}^l = \{f^{i,j} : f \in \{\pm 1\}^{[d]}, i \in [l], j \in \{-1, 1\}\}.$$

For $g \in \{-1, 1\}^{[d]}$, $i \in [l]$, $j \in \{\pm 1\}$ define $g^{i,j} : [d] \times [l] \rightarrow \{-1, 1\}$ by

$$g^{i,j}(u, v) = \begin{cases} h(u) & v = i \\ j & v > i \\ -j & v < i, \end{cases}$$

And define the hypothesis class \mathcal{G}^l as

$$\mathcal{G}^l = \{g^{i,j} : g \in \{-1, 1\}^{[d]}, i \in [l], j \in \{\pm 1\}\}.$$

Let $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$, $\mathcal{H}' \subset \mathcal{Y}^{\mathcal{X}'}$ be two hypotheses classes. We say that \mathcal{H} is *richer* than \mathcal{H}' if there is a mapping $\iota : \mathcal{X}' \rightarrow \mathcal{X}$ such that $\mathcal{H}' = \mathcal{H} \circ \iota$. It is clear that if \mathcal{H} is richer than \mathcal{H}' then $d_N(\mathcal{H}') \leq d_N(\mathcal{H})$ and $d_G(\mathcal{H}') \leq d_G(\mathcal{H})$. Thus, the notion of richness can be used to establish lower and upper bounds on the Natarajan and Graph dimension, respectively. The following lemma shows that \mathcal{W} is richer than \mathcal{F}^l and \mathcal{G}^l for every l . This will allow us to use the classes \mathcal{F}^l , \mathcal{G}^l instead of \mathcal{W} when bounding from below the dimension of an ECOC or TC hypothesis class in which the binary classifiers are from \mathcal{W} .

Lemma A.1. *For any integer $l \geq 2$, \mathcal{W} is richer than \mathcal{F}^l and \mathcal{G}^l .*

Proof. We shall first prove that \mathcal{W} is richer than \mathcal{F}^l . Choose l unit vectors $e_1, \dots, e_l \in \mathbb{R}^d$. For every $i \in [l]$, choose d affinely independent vectors such that

$$x_{1,i}, \dots, x_{d,i} \in \{x \in \mathbb{R}^d : \langle x, e_i \rangle = 1, \forall i' \neq i, \langle x, e_{i'} \rangle < 1\}.$$

This can be done by choosing d affinely independent vectors in $\{x \in \mathbb{R}^d : \langle x, e_i \rangle = 1\}$ that are very close to e_i . Define $\iota(m, i) = x_{m,i}$. Now fix $i \in [l]$ and $j \in \{-1, +1\}$, and let $f^{i,j} \in \mathcal{F}^l$. We must show that $f^{i,j} = h \circ \iota$ for some $h \in \mathcal{W}$. We will show that there exists an affine map $\Lambda : \mathbb{R}^d \rightarrow \mathbb{R}$ for which $f^{i,j} = \text{sign} \circ \Lambda \circ \iota$. This suffices, since \mathcal{W} is exactly the set of all functions of the form $\text{sign} \circ \Lambda$ where Λ is an affine map. Define $M = \{x \in \mathbb{R}^d : \langle x, e_i \rangle = 1\}$, and let $A : M \rightarrow \mathbb{R}$ be the affine map defined by

$$\forall m \in [d], A(x_{m,i}) = f(m, i).$$

Let $P : \mathbb{R}^d \rightarrow M$ be the orthogonal projection of \mathbb{R}^d on M . For $\alpha \in \mathbb{R}$, define an affine map $\Lambda_\alpha : \mathbb{R}^d \rightarrow \mathbb{R}$ by

$$\Lambda_\alpha(x) = A(P(x)) + \alpha \cdot \langle x - e_i, e_i \rangle.$$

Note that, $\forall m \in [d]$, $\Lambda_\alpha(x_{m,i}) = f(m, i)$. Moreover, for every $i' \neq i$ and $m \in [d]$ we have $\langle x_{m,i'} - e_i, e_i \rangle < 0$. Thus, by choosing $|\alpha|$ sufficiently large and choosing $\text{sign}(\alpha)$ depending on j , we can make sure that $f^{i,j} = \text{sign} \circ \Lambda_\alpha \circ \iota$.

The proof that \mathcal{W} is richer than \mathcal{G}^l is similar and simpler. Let $e_1, \dots, e_d \in \mathbb{R}^{d-1}$ be affinely independent. Define

$$\iota(m, i) = (e_m, i) \in \mathbb{R}^{d-1} \times \mathbb{R} \cong \mathbb{R}^d,$$

Given $g^{i,j} \in \mathcal{G}^{d,l}$, let $A : \mathbb{R}^{d-1} \times \{i\} \rightarrow \mathbb{R}$ be the affine map defined by $A(e_m, i) = g^{i,j}(m, i)$ and let $P : \mathbb{R}^d \rightarrow \mathbb{R}^{d-1} \times \{i\}$ be the orthogonal projection. Define $\Lambda : \mathbb{R}^d \rightarrow \mathbb{R}$ by

$$\Lambda(x, y) = A(P(x, y)) + j \cdot 10 \cdot (y - i).$$

It is easy to check that $\text{sign} \circ \Lambda \circ \iota = g^{i,j}$. \square

Note A.2. From Lemma A.1 it follows that $\text{VC}(\mathcal{F}^l), \text{VC}(\mathcal{G}^l) \leq d + 1$. On the other hand, both \mathcal{F}^l and \mathcal{G}^l shatter $([d] \times \{1\}) \cup \{(1, 2)\}$. Thus, $\text{VC}(\mathcal{F}^l) = \text{VC}(\mathcal{G}^l) = d + 1$

A.4 Trees

Proof of Theorem 3.2. We first prove the upper bound. Let $A \subseteq \mathcal{X}$ be a G -shattered set with $|A| = d_G(\mathcal{H}_{\text{trees}})$. By Sauer's Lemma, and since the number of trees is bounded by k^k , we have $2^{|A|} \leq k^k \cdot |\mathcal{H}|_{|A|}^k \leq k^k \cdot |A|^{dk}$, thus $d_G(\mathcal{H}_{\text{trees}}) = |A| = O(dk \log(dk))$.

To prove the lower bound, by Lemma A.1, it is enough to show that $d_N(\mathcal{G}_T^l) \geq d \cdot (k - 1)$ for some l . We will take $l = |N(T)| = k - 1$. Linearly order $N(T)$ such that for every node v , the nodes in the left sub-tree emanating from v are smaller than the nodes in the corresponding right sub-tree. We will identify $[l]$ with $N(T)$ by an order-preserving map, thus $\mathcal{G}^l \subset \{-1, 1\}^{[d] \times N(T)}$. We also identify the labels with the leaves.

Define $g_1 : [d] \times N(T) \rightarrow \text{leaf}(T)$ by setting $g_1(i, v)$ to be the leaf obtained by starting from the node v , going right once and then going left until reaching a leaf. Similarly, define $g_2 : [d] \times N(T) \rightarrow \text{leaf}(T)$ by setting $g_2(i, v)$ to be the leaf obtained by starting from the node v , going left once and then going right until reaching a leaf.

We shall show that g_1, g_2 witness the N -shattering of $[d] \times N(T)$ by \mathcal{G}_T^l . Given $S \subset [d] \times N(T)$ define $C : N(T) \rightarrow \mathcal{G}^l$ by

$$C(v)(i, u) = \begin{cases} -1 & u < v \\ 1 & u > v \\ 1 & u = v, (i, u) \in S \\ -1 & u = v, (i, u) \notin S. \end{cases}$$

It is not hard to check that $\forall (i, u) \in S, h_C(i, u) = g_1(i, u)$, and $\forall (i, u) \notin S, h_C(i, u) = g_2(i, u)$. \square

Note A.3. Define $\tilde{\mathcal{G}}^l = \{g^{i,1} : g \in \{-1, 1\}^{[d]}, i \in [l]\}$. The proof shows that $d_N(\tilde{\mathcal{G}}_T^l) \geq d \cdot (k - 1)$. Since $\text{VC}(\tilde{\mathcal{G}}^l) = d$, we obtain a simpler proof of Theorem 23 from Daniely et al. [2011], which states that for every tree T there exists a class \mathcal{H} of VC dimension d for which $d_N(\mathcal{H}_T) \geq d(k - 1)$.

A.5 ECOC, One vs. All and All Pairs

To prove the results for ECOC and its special cases, we first prove a more general theorem, based on the notion of a sensitive vector for a given code. Fix a code $M \in \mathbb{R}^{k \times l}(\mathbb{R})$. We say that a binary vector $u \in \{\pm 1\}^l$ is q -sensitive for M if there are q indices $j \in [l]$ for which $\tilde{M}(u) \neq \tilde{M}(u \oplus e_j)$. Here, $u \oplus e_j := (u[1], \dots, -u[j], \dots, u[l])$.

Theorem A.4. *If there exists a q -sensitive vector for a code $M \in \mathbb{R}^{k \times l}(\mathbb{R})$ then $d_N(\mathcal{W}_M) \geq d \cdot q$.*

Proof. By Lemma A.1, it suffices to show that $d_N(\mathcal{F}_M^l) \geq d \cdot q$. Let $u \in \{\pm 1\}^l$ be a q -sensitive vector. Assume w.l.o.g. that the sensitive coordinates are $1, \dots, q$. We shall show that $[d] \times [q]$ is N -shattered by \mathcal{F}_M^l . Define $g_1, g_2 : [d] \times [q] \rightarrow [k]$ by

$$g_1(x, y) = \tilde{M}(u), \quad g_2(x, y) = \tilde{M}(u \oplus e_y)$$

Let $T \subset [d] \times [q]$. Define $h_1, \dots, h_l \in \mathcal{F}^l$ as follows. For every $j > q$, define $h_j \equiv u[j]$. For $j \leq q$ define

$$h_j(x, y) = \begin{cases} u[j] & y \neq j \\ u[j] & y = j, (x, y) \in T \\ -u[j] & y = j, (x, y) \in [d] \times [q] \setminus T. \end{cases}$$

For $h = (h_1, \dots, h_l)$, it is not hard to check that

$$\begin{aligned} \forall (x, y) \in T, \quad \tilde{M}(h_1(x, y), \dots, h_l(x, y)) &= g_1(x, y), \text{ and} \\ \forall (x, y) \in [d] \times [q] \setminus T, \quad \tilde{M}(h_1(x, y), \dots, h_l(x, y)) &= g_2(x, y). \end{aligned}$$

□

The following lemma shows that a code with a large distance is also highly sensitive. In fact, we prove a stronger claim: the sensitivity is actually at least as large as the distance between any row and the row closest to it in Hamming distance. Formally, we consider $\Delta(M) = \max_i \min_{j \neq i} \Delta_h(M[i], M[j]) \geq \delta(M)$.

Lemma A.5. *For any binary code $M \in \mathbb{R}^{k \times l}(\pm 1)$, there is a q -sensitive vector for M , where $q \geq \frac{1}{2}\Delta(M) \geq \frac{1}{2}\delta(M)$.*

Proof. Let i_1 the row in M such that its hamming distance to the row closest to it is $\Delta(M)$. Denote by i_2 the index of the closest row (if there is more than one such row, choose one of them arbitrarily). We have $\Delta_h(M[i_1], M[i_2]) = \Delta(M)$. In addition, $\forall i \neq i_1, i_2, \Delta_h(M[i_1], M[i]) \geq \Delta(M)$. Assume w.l.o.g. that the indices in which rows i_1 and i_2 differ are $1, \dots, \Delta(M)$. Consider first the case that $i_1 < i_2$. Define $u \in \{\pm 1\}^{[l]}$ by

$$u[j] = \begin{cases} M_{(i_1, j)} & j \leq \lceil \frac{\Delta}{2} \rceil \\ M_{(i_2, j)} & \text{otherwise.} \end{cases}$$

It is not hard to check that for every $1 \leq j \leq \lceil \frac{\Delta}{2} \rceil$, $i_1 = \tilde{M}(u)$ and $\tilde{M}(u \oplus e_j) = i_2$, thus u is $\lceil \frac{\Delta}{2} \rceil$ -sensitive. If $i_1 > i_2$, the proof is similar except that u is defined as

$$u[j] = \begin{cases} M_{(i_2, j)} & j \leq \lceil \frac{\Delta}{2} \rceil \\ M_{(i_1, j)} & \text{otherwise.} \end{cases}$$

□

Proof of Theorem 3.3. The upper bound follows from Lemma 4.1. The lower bound follows from Theorem A.4 and Lemma A.5. □

Proof of Theorem 3.4. The upper bounds follow from Theorem 3.3. To show that $d_N(\mathcal{W}_{\text{OVA}}) \geq (k-1)d$, we note that the all-negative vector $u = (-1, \dots, -1)$ of length k is $(k-1)$ -sensitive for the code M^{OVA} , and apply Theorem A.4.

To show that $d_N(\mathcal{W}_{\text{AP}}) \geq d \binom{k-1}{2}$, assume for simplicity that k is odd (a similar analysis can be given when k is even). Define $u \in \{\pm 1\}^{\binom{k}{2}}$ by

$$\forall i < j, u[i, j] = \begin{cases} 1 & j - i \leq \frac{k-1}{2} \\ -1 & \text{otherwise.} \end{cases}$$

For every $n \in [k]$, we have $\sum_{1 \leq i < j \leq k} u[i, j] \cdot M_{n, (i, j)}^{\text{AP}} = 0$, as the summation counts the number of pairs (i, j) such that $n \in \{i, j\}$ and $M_{n, (i, j)}^{\text{AP}}$ agrees with $u[i, j]$. Thus, $\tilde{M}^{\text{AP}}(u) = 1$, by our tie-breaking assumptions. Moreover, it follows that for every $1 < i < j \leq k$, we have $\tilde{M}^{\text{AP}}(u \oplus e_{(i, j)}) \in \{i, j\}$, since flipping entry $[i, j]$ of u increases $(M^{\text{AP}}u)_j$ or $(M^{\text{AP}}u)_i$ by 1 and does not increase the rest of the coordinates of the vector $M^{\text{AP}}u$. This shows that u is $\binom{k-1}{2}$ -sensitive. □

A.6 Approximation

Proof of Theorem 3.5. We first show that for any tree for k classes T , \mathcal{L} essentially contains \mathcal{W}_T . It follows that \mathcal{L} essentially contains $\mathcal{W}_{\text{trees}}$ as well. Let \mathcal{D} a distribution over \mathbb{R}^d , let $C : N(T) \rightarrow \mathcal{W}$ be a mapping associating nodes in T to binary classifiers in \mathcal{W} , and let $\epsilon > 0$. We will show that there exists a matrix $W \in \mathbb{R}^{k \times (d+1)}$ such that $\Pr_{x \sim \mathcal{D}}[h_W(x) \neq h_C(x)] < \epsilon$.

For every $v \in N(T)$, denote by $w(v) \in \mathbb{R}^{d+1}$ the linear separator such that $C(v) = h_{w(v)}$. For every $w \in \mathbb{R}^{d+1}$ define $\tilde{w} = w + (0, \dots, 0, \gamma)$. Recall that for $x \in \mathbb{R}^d$, $\bar{x} \in \mathbb{R}^{d+1}$ is simply the concatenation $(x, 1)$. Choose $r > 0$ large enough so that $\Pr_{x \sim \mathcal{D}}[\|\bar{x}\| > r] < \epsilon/2$ and $\forall v \in N(T)$, $\|\tilde{w}(v)\| < r$. Choose $\gamma > 0$ small enough so that

$$\Pr_{x \sim \mathcal{D}}[\exists v \in N(T), \langle \tilde{w}(v), \bar{x} \rangle \in (-\gamma, \gamma)] = \Pr_{x \sim \mathcal{D}}[\exists v \in N(T), \langle w(v), \bar{x} \rangle \in (-2\gamma, 0)] < \epsilon/2.$$

Let $a = 2r^2/\gamma + 1$. For $i \in [k]$, let $v_{i,1}, \dots, v_{i,m_i}$ be the path from the root to the leaf associated with label i . For each $1 \leq j < m_i$ define $b_{i,j} = 1$ if $v_{i,j+1}$ is the right son of $v_{i,j}$, and $b_{i,j} = -1$ otherwise. Now, define $W \in \mathbb{R}^{k \times (d+1)}$ to be the matrix whose i 'th row is $w_i = \sum_{j=1}^{m_i-1} a^{-j} \cdot b_{i,j} \tilde{w}(v_{i,j})$.

To prove that $\Pr_{x \sim \mathcal{D}}[h_W(x) \neq h_C(x)] < \epsilon$, it suffices to show that $h_W(x) = h_C(x)$ for every $x \in \mathbb{R}^d$ satisfying $\|\bar{x}\| < r$ and $\forall v \in N(T)$, $\langle \tilde{w}(v), \bar{x} \rangle \notin (-\gamma, \gamma)$, since the probability mass of the rest of the vectors is less than ϵ . Let $x \in \mathbb{R}^d$ be a vector that satisfies these assumptions. Denote $i_1 = h_C(x)$. It suffices to show that for all $i_2 \in [k] \setminus \{i_1\}$, $\langle w_{i_1}, \bar{x} \rangle > \langle w_{i_2}, \bar{x} \rangle$, since this would imply that $h_W(x) = i_1$ as well.

Indeed, fix $i_2 \neq i_1$, and let j_0 be the length of the joint prefix of the two root-to-leaf paths that match the labels i_1 and i_2 . In other words, $\forall j \leq j_0$, $v_{i_1,j} = v_{i_2,j}$ and $v_{i_1,j_0+1} \neq v_{i_2,j_0+1}$. Note that

$$\langle \bar{x}, (b_{i_1,j_0} - b_{i_2,j_0}) \tilde{w}(v_{i_1,j_0}) \rangle = \langle \bar{x}, 2b_{i_1,j_0} \tilde{w}(v_{i_1,j_0}) \rangle = 2|\langle \bar{x}, \tilde{w}(v_{i_1,j_0}) \rangle| \geq 2\gamma.$$

The last equality holds because b_{i_1,j_0} and $\langle \bar{x}, w(v_{i_1,j_0}) \rangle$ have the same sign by definition of $b_{i,j}$. We have

$$\begin{aligned} \langle w_{i_1}, \bar{x} \rangle - \langle w_{i_2}, \bar{x} \rangle &= \langle \bar{x}, \sum_{j=1}^{m_{i_1}-1} a^{-j} b_{i_1,j} \tilde{w}(v_{i_1,j}) - \sum_{j=1}^{m_{i_2}-1} a^{-j} b_{i_2,j} \tilde{w}(v_{i_2,j}) \rangle \\ &= \langle \bar{x}, a^{-j_0} (b_{i_1,j_0} - b_{i_2,j_0}) \tilde{w}(v_{i_1,j_0}) \rangle + \langle \bar{x}, \sum_{j=j_0+1}^{m_{i_1}-1} a^{-j} b_{i_1,j} \tilde{w}(v_{i_1,j}) - \sum_{j=j_0+1}^{m_{i_2}-1} a^{-j} b_{i_2,j} \tilde{w}(v_{i_2,j}) \rangle \\ &\geq \langle \bar{x}, a^{-j_0} (b_{i_1,j_0} - b_{i_2,j_0}) \tilde{w}(v_{i_1,j_0}) \rangle - \sum_{j=j_0+1}^{\infty} a^{-j} 2r^2 \\ &\geq 2a^{-j_0} \left(\gamma - \frac{r^2}{a-1} \right) > 0. \end{aligned}$$

Since this holds for all $i_2 \neq i_1$, it follows that $h_W(x) = i_1$. Thus, we have proved that \mathcal{L} essentially contains $\mathcal{W}_{\text{trees}}$.

Next, we show that \mathcal{L} strictly contains $\mathcal{W}_{\text{trees}}$, by showing a distribution over labeled examples such that the approximation error using \mathcal{L} is strictly smaller than the approximation error using $\mathcal{W}_{\text{trees}}$. Assume w.l.o.g. that $d = 2$ and $k = 3$: even if they are larger we can always restrict the support of the distribution to a subspace of dimension 2 and to only three of the labels. Consider the distribution \mathcal{D} over $\mathbb{R}^2 \times [3]$ such that its marginal over \mathbb{R}^2 is uniform in the unit circle, and $\Pr_{(X,Y) \sim \mathcal{D}}[Y = i \mid X = x] = \mathbb{I}[x \in D_i]$, where D_1, D_2, D_3 be subsets sectors of equal angle of the unit circle (see Figure 1):

Clearly, by taking the rows of W to point to the middle of each sector (dashed arrows in the illustration), we get $\text{Err}_{\mathcal{D}}^*(\mathcal{L}) = 0$. In contrast, no linear separator can split the three labels into two groups without error, thus $\text{Err}_{\mathcal{D}}^*(\mathcal{W}_{\text{trees}}) > 0$.

Finally, to see that \mathcal{L} essentially contains \mathcal{W}_{OVA} , we note that $\mathcal{W}_{\text{OVA}} = \mathcal{W}_T$ where T is a tree such that each of its internal nodes has a leaf corresponding to one of the labels as its left son. Thus \mathcal{W}_{OVA} is essentially contained in $\mathcal{W}_{\text{trees}}$. \square

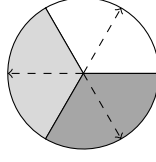


Figure 1: Illustration for the proof of Theorem 3.5

Proof of Theorem 3.7. It is easily seen that \mathcal{W}_{AP} contains \mathcal{L} : Let $W \in \mathbb{R}^{d+1 \times k}$, and denote its i 'th row by $W[i]$. For each column (i, j) of M^{AP} , define the binary classifier $h_{i,j} \in \mathcal{W}$ such that $\forall x \in \mathbb{R}^d$, $h_{i,j}(\bar{x}) = \text{sign}(\langle W[j] - W[i], \bar{x} \rangle)$. Then for all x , $h_W(x) = \tilde{M}^{AP}(h_{1,1}(x), \dots, h_{k-1,k}(x))$.

To show that the inclusion is strict, as in the proof of Theorem 3.5, we can and will assume that $d = 2$. Choose k^* to be the minimal number such that for every $k \geq k^*$, $d_N(\mathcal{W}_{AP}) > d_N(\mathcal{L})$: This number exists by Theorems 3.4 and 3.1 (note that though we chose k^* w.r.t. $d = 2$, the same k^* is valid for every d). For any $k \geq k^*$, it follows that there is a set $S \subseteq \mathbb{R}^2$ that is N -shattered by \mathcal{W}_{AP} but not by \mathcal{L} . Thus, there is a hypothesis $h \in \mathcal{W}_{AP}$ such that for every $g \in \mathcal{L}$, $g|_S \neq h|_S$. Define the distribution \mathcal{D} to be uniform over $\{(x, h(x)) : x \in S\}$. Then clearly $\text{Err}_{\mathcal{D}}^*(\mathcal{L}) > \text{Err}_{\mathcal{D}}^*(\mathcal{W}_{AP}) = 0$. \square

Next, we prove Theorem 3.6, which we restate more formally as follows. Note that the result on OvA is implied since there exists a tree that implements OvA.

Theorem A.6. (Restatement of Theorem 3.6) *If there exists an embedding $\iota : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ and a tree T such that $\mathcal{W}_T^{d'} \circ \iota$ essentially contains \mathcal{L} , then necessarily $d' \geq \tilde{\Omega}(dk)$.*

Proof. Assume that $i \in [k]$ is the class corresponding to the leaf with the least depth, l . Note that $l \leq \log_2(k)$. Let $\phi : [k] \rightarrow \{\pm 1\}$ be the function that is 1 on $\{i\}$ and -1 otherwise. It is not hard to see that $\phi \circ \mathcal{L}$ is the hypothesis class of convex polyhedra in \mathbb{R}^d having $k - 1$ faces. Thus,

$$\text{VC}(\phi \circ \mathcal{L}) \geq (k - 1)d, \quad (3)$$

[see e.g. Takacs, 2009]. On the other hand, $\phi \circ \mathcal{W}_T^{d'} \circ \iota$, is the class of convex polyhedra in $\mathbb{R}^{d'}$ having $l \leq \log_2(k)$ faces. Thus, by Lemma 4.1

$$\text{VC}(\phi \circ \mathcal{W}_T^{d'} \circ \iota) \leq \text{VC}(\phi \circ \mathcal{W}_T^{d'}) \leq O(ld' \log(ld')) \leq O(\log(k)d' \log(\log(k)d')) \quad (4)$$

By the assumption that $\mathcal{W}_T^{d'} \circ \iota$ essentially contains \mathcal{L} , $\text{VC}(\phi \circ \mathcal{L}) \leq \text{VC}(\phi \circ \mathcal{W}_T^{d'} \circ \iota)$. Combining with equations (3) and (4) it follows that $d(k - 1) = O(\log(k)d' \log(\log(k)d'))$. Thus, $d' = \tilde{\Omega}(dk)$. \square

To prove Lemma 3.8, we first state the classic VC-dimension theorem, which will be useful to us.

Theorem A.7 (Vapnik [1998]). *There exists a constant $C > 0$ such that for every hypothesis class $\mathcal{H} \subseteq \{\pm 1\}^{\mathcal{X}}$ of VC dimension d , a distribution \mathcal{D} over \mathcal{X} , $\epsilon, \delta > 0$ and $m \geq C \frac{d + \ln(\frac{1}{\delta})}{\epsilon^2}$ we have*

$$\Pr_{S \sim \mathcal{D}^m} \left[\text{Err}_{\mathcal{D}}^*(\mathcal{H}) \geq \inf_{h \in \mathcal{H}} \text{Err}_S(h) - \epsilon \right] \geq 1 - \delta.$$

We also use the following lemma, which proves a variant of Hoeffding's inequality.

Lemma A.8. *Let $\beta_1, \dots, \beta_k \geq 0$ and let $\gamma_1, \dots, \gamma_k \in \mathbb{R}$, such that $\forall i$, $|\gamma_i| \leq \beta_i$. Fix an integer $j \in \{1, \dots, \lfloor \frac{k}{2} \rfloor\}$ and let $\mu = j/k$. Let $(X_1, \dots, X_k) \in \{\pm 1\}^k$ be a random vector sampled uniformly from the set $\{(x_1, \dots, x_k) : \sum_{i=1}^k \frac{x_i + 1}{2} = \mu k\}$. Define $Y_i = \beta_i + X_i \gamma_i$ and denote $\alpha_i = \beta_i + |\gamma_i|$. Assume that $\sum_{i=1}^k \alpha_i = 1$. Then*

$$\Pr \left[\sum_{i=1}^k Y_i \leq \mu - \epsilon \right] \leq 2 \exp \left(- \frac{\epsilon^2}{2 \sum_{i=1}^k \alpha_i^2} \right).$$

Proof. First, since $\mu < \frac{1}{2}$, it suffices to prove the claim for the case $\forall i, \gamma_i \geq 0$ since this is the “harder” case. Let $Z_1, \dots, Z_k \in \{\pm 1\}$ be independent random variables such that $\Pr[Z_i = 1] = \mu - \frac{\epsilon}{2}$. Denote $W_i = \beta_i + Z_i \gamma_i$. Further denote $\bar{W} = \sum_{i=1}^k W_i$ and $\bar{Z} = \sum_{i=1}^k \frac{Z_i + 1}{2}$.

Note that for every $j_0 \leq j = \mu k$, given that $\bar{Z} = j_0$, \bar{W} can be described as follows: We start with the value $\sum_{i=1}^k \beta_i - \gamma_i$ and then choose j_0 indices uniformly from $[k]$. For each chosen index i , the value of \bar{W} is increased by $2\gamma_i$. $\sum_{i=1}^k Y_i$ can be described in the same way, except that that $j \geq j_0$ indices are chosen. Thus, $\Pr \left[\sum_{i=1}^k Y_i \leq \mu - \epsilon \right] \leq \Pr \left[\bar{W} \leq \mu - \epsilon \mid \bar{Z} = j_0 \right]$. Thus, we have

$$\begin{aligned} \Pr \left[\sum_{i=1}^k Y_i \leq \mu - \epsilon \right] &\leq \Pr \left[\bar{W} \leq \mu - \epsilon \mid \bar{Z} \leq \mu k \right] \\ &\leq \Pr \left[\bar{W} \leq \mu - \epsilon \right] / \Pr \left[\bar{Z} \leq \mu k \right] \\ &\leq 2 \Pr \left[\bar{W} \leq \mu - \epsilon \right] \\ &\leq 2 \exp \left(-\frac{\epsilon^2}{2 \sum_{i=1}^k \alpha_i^2} \right). \end{aligned}$$

The last inequality follows from Hoeffding’s inequality and noting that

$$E[W_i] = \beta_i + (2(\mu - \frac{\epsilon}{2}) - 1)\gamma_i = (\mu - \frac{\epsilon}{2})(\beta_i + \gamma_i) + (1 - \mu + \frac{\epsilon}{2})(\beta_i - \gamma_i) \geq (\mu - \frac{\epsilon}{2})\alpha_i.$$

So that $\sum_{i=1}^k E[W_i] \geq (\mu - \frac{\epsilon}{2}) \sum_{i=1}^k \alpha_i = \mu - \frac{\epsilon}{2}$. \square

Proof of Lemma 3.8. The idea of this proof is as follows: Using a uniform convergence argument based on the VC dimension of the binary hypothesis class, we show that there exists a labeled sample S such that $|S| \approx \frac{d+k}{\nu^2}$, and for all possible mappings ϕ , the approximation error of the hypothesis class on the sample is close to the approximation error on the distribution \mathcal{D}_ϕ . This allows us to restrict our attention to a finite set of hypotheses, based on their restriction to the sample. For these hypotheses, we show that with high probability over the choice of ϕ , the approximation error on the sample is high. Using a union bound on the possible hypotheses, we conclude that the approximation error on the distribution will be high, with high probability over the choice of ϕ .

For $i \in [k]$, denote $p_i = \Pr_{x \sim \mathcal{D}}[f(x) = i]$. Let $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \subseteq \mathcal{X} \times [k]$ be an i.i.d. sample drawn according to \mathcal{D} where $m = \lceil C \frac{d+(k+2)\ln(2)}{(\nu/2)^2} \rceil$, for the constant from C from Theorem A.7. Given S , denote $S_\phi = \{(x_1, \phi(y_1)), \dots, (x_m, \phi(y_m))\} \subseteq \mathcal{X} \times \{\pm 1\}$. For $i \in [k]$, let $\hat{p}_i = \frac{|\{j: y_j = i\}|}{m}$.

For any fixed $\phi : [k] \rightarrow \{\pm 1\}$, with probability $> 1 - 2^{-(k+2)}$ over the choice of S we have, by Theorem A.7, that $\text{Err}_{\mathcal{D}_\phi}^*(\mathcal{H}) > \inf_{h \in \mathcal{H}} \text{Err}_{S_\phi}(h) - \nu$. Since $|\{\pm 1\}^{[k]}| = 2^k$, w.p. $> 1 - \frac{1}{4}$,

$$\forall \phi \in \{\pm 1\}^{[k]}, \quad \text{Err}_{\mathcal{D}_\phi}^*(\mathcal{H}) > \inf_{h \in \mathcal{H}} \text{Err}_{S_\phi}(h) - \frac{\nu}{2}. \quad (5)$$

Moreover, we have

$$E \left[\sum_{i=1}^k \hat{p}_i^2 \right] = \frac{1}{m^2} \sum_{i=1}^k \left(\binom{m}{2} p_i^2 + m p_i \right) \leq k \cdot \left(\frac{m(m-1)}{2m^2} \frac{100}{k^2} + \frac{10}{mk} \right) \leq \frac{60}{k}.$$

Thus, by Markov’s inequality, w.p. $\geq \frac{1}{2}$ we have

$$\sum_{i=1}^k \hat{p}_i^2 < \frac{120}{k}. \quad (6)$$

Thus, with probability at least $1 - \frac{1}{4} - \frac{1}{2} > 0$, both (6) and (5) holds. In particular, there exists a sample S for which both (6) and (5) hold. Let us fix such an $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$.

Assume now that $\phi \in \{\pm 1\}^{[k]}$ is sampled according to the first condition. Denote

$$Y_i = |\{j : h(x_j) \neq \phi(y_j) \text{ and } y_j = i\}| / m.$$

For a fixed $h \in \mathcal{H}$ we have

$$\Pr_{\phi} \left[\text{Err}_{S_{\phi}}(h) < \mu - \frac{\nu}{2} \right] = \Pr_{\phi} \left[\sum_{i=1}^k Y_i < \mu - \frac{\nu}{2} \right]$$

We note that Y_i are independent random variables with $E[Y_i] \geq \mu \hat{p}_i$ and $0 \leq Y_i \leq \hat{p}_i$. Thus, by Hoeffding's inequality,

$$\Pr_{\phi} \left[\text{Err}_{S_{\phi}}(h) < \mu - \frac{\nu}{2} \right] \leq \exp \left(-\frac{\nu^2}{2 \sum_{i=1}^k \hat{p}_i^2} \right) \leq \exp \left(-\frac{\nu^2 k}{240} \right).$$

By Sauer's lemma, $|\mathcal{H}|_{\{x_1, \dots, x_m\}} \leq \left(\frac{em}{d}\right)^d$. Thus, with probability $\geq 1 - \left(\frac{em}{d}\right)^d \exp\left(-\frac{\nu^2 k}{240}\right)$ over the choice of ϕ , $\inf_{h \in \mathcal{H}} \text{Err}_{S_{\phi}}(h) \geq \mu - \frac{\nu}{2}$ and by (5) also

$$\text{Err}_{\mathcal{D}_{\phi}}^*(\mathcal{H}) \geq \frac{1}{2} - \nu. \quad (7)$$

Finally, since $m = O\left(\frac{k+d}{\nu^2}\right)$, if $k = \Omega\left(\frac{d \ln(1/\nu) + \ln(1/\delta)}{\nu^2}\right)$ then Eq. (7) holds w.p $> 1 - \delta$, concluding the proof for the case when the first condition holds. If the second condition holds, the proof is very similar, with the sole difference that Lemma A.8 is used instead of Hoeffding's inequality. \square

Proof of Corollary 3.9. The Corollary follows from Lemma 3.8, by noting that $\text{Err}_{\mathcal{D}}^*(\mathcal{H}_T) \geq \text{Err}_{\mathcal{D}_{\phi}}^*(\mathcal{H})$, where $\phi : [k] \rightarrow \{\pm 1\}$ is defined as $\phi(i) = 1$ if and only if $\lambda^{-1}(i)$ is in the right subtree emanating from the root of T . \square

Proof of Corollary 3.10. Let $\phi : [k] \rightarrow \{\pm 1\}$ be the function that is -1 on $\left[\left\lfloor \frac{k}{2} \right\rfloor\right]$ and 1 otherwise. By Lemma 4.1, applied to $L(\mathcal{H}) = \phi \circ \mathcal{H}_{(M, \text{Id})}$, $\text{VC}(\phi \circ \mathcal{H}_{(M, \text{Id})}) = O(d \log(d))$, so that, by Lemma 3.8 (applied to a random choice of λ instead of ϕ), $\text{Err}_{\mathcal{D}_{\phi \circ \lambda}}^*(\phi \circ \mathcal{H}_{(M, \text{Id})}) \geq \frac{1}{2} - \nu$ with probability $> 1 - \delta$ over the choice of λ . The proof follows as we note that for every λ , $\text{Err}_{\mathcal{D}}^*(\mathcal{H}_{(M, \lambda^{-1})}) = \text{Err}_{\mathcal{D}_{\lambda}}^*(\mathcal{H}_{(M, \text{Id})}) \geq \text{Err}_{\mathcal{D}_{\phi \circ \lambda}}^*(\phi \circ \mathcal{H}_{(M, \text{Id})})$. \square