

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Agnostic Active Learning Without Constraints

Anonymous Author(s)

Affiliation

Address

email

Abstract

We present and analyze an agnostic active learning algorithm that works without keeping a version space. This is unlike all previous approaches where a restricted set of candidate hypotheses is maintained throughout learning, and only hypotheses from this set are ever returned. By avoiding this version space approach, our algorithm sheds the computational burden and brittleness associated with maintaining version spaces, yet still allows for substantial improvements over supervised learning for classification.

1 Introduction

In active learning, a learner is given access to unlabeled data and is allowed to adaptively choose which ones to label. This learning model is motivated by applications in which the cost of labeling data is high relative to that of collecting the unlabeled data itself. Therefore, the hope is that the active learner only needs to query the labels of a small number of the unlabeled data, and otherwise perform as well as a fully supervised learner. In this work, we are interested in agnostic active learning algorithms for binary classification that are provably consistent, *i.e.* that converge to an optimal hypothesis in a given hypothesis class.

One technique that has proved theoretically profitable is to maintain a candidate set of hypotheses (sometimes called a version space), and to query the label of a point only if there is disagreement within this set about how to label the point. The criteria for membership in this candidate set needs to be carefully defined so that an optimal hypothesis is always included, but otherwise this set can be quickly whittled down as more labels are queried. This technique is perhaps most readily understood in the noise-free setting [1, 2], and it can be extended to noisy settings by using empirical confidence bounds [3, 4, 5, 6, 7].

The version space approach unfortunately has its share of significant drawbacks. The first is computational intractability: maintaining a version space and guaranteeing that *only* hypotheses from this set are returned is difficult for linear predictors and appears intractable for interesting nonlinear predictors such as neural nets and decision trees [1]. Another drawback of the approach is its brittleness: a single mishap (due to, say, modeling failures or computational approximations) might cause the learner to exclude the best hypothesis from the version space forever; this is an ungraceful failure mode that is not easy to correct. A third drawback is related to sample re-usability: if (labeled) data is collected using a version space-based active learning algorithm, and we later decide to use a different algorithm or hypothesis class, then the earlier data may not be freely re-used because its collection process is inherently biased.

Here, we develop a new strategy addressing all of the above problems given an oracle that returns an empirical risk minimizing (ERM) hypothesis. As this oracle matches our abstraction of many supervised learning algorithms, we believe active learning algorithms built in this way are immediately and widely applicable.

054 Our approach instantiates the importance weighted active learning framework of [5] using a rejection
 055 threshold similar to the algorithm of [4] which only accesses hypotheses via a supervised learning
 056 oracle. However, the oracle we require is simpler and avoids strict adherence to a candidate set
 057 of hypotheses. Moreover, our algorithm creates an importance weighted sample that allows for
 058 unbiased risk estimation, even for hypotheses from a class different from the one employed by the
 059 active learner. This is in sharp contrast to many previous algorithms (*e.g.*, [1, 3, 8, 4, 6, 7]) that create
 060 heavily biased data sets. We prove that our algorithm is always consistent and has an improved label
 061 complexity over passive learning in cases previously studied in the literature. We also describe a
 062 practical instantiation of our algorithm and report on some experimental results.

064 1.1 Related Work

065
 066 As already mentioned, our work is closely related to the previous works of [4] and [5], both of
 067 which in turn draw heavily on the work of [1] and [3]. The algorithm from [4] extends the selective
 068 sampling method of [1] to the agnostic setting using generalization bounds in a manner similar
 069 to that first suggested in [3]. It accesses hypotheses only through a special ERM oracle that can
 070 enforce an arbitrary number of example-based constraints; these constraints define a version space,
 071 and the algorithm only ever returns hypotheses from this space, which can be undesirable as we
 072 previously argued. Other previous algorithms with comparable performance guarantees also require
 073 similar example-based constraints (*e.g.*, [3, 5, 6, 7]). Our algorithm differs from these in that (i) it
 074 never restricts its attention to a version space when selecting a hypothesis to return, and (ii) it only
 075 requires an ERM oracle that enforces at most one example-based constraint, and this constraint is
 076 only used for selective sampling. Our label complexity bounds are comparable to those proved in [5]
 (though somewhat worse than those in [3, 4, 6, 7]).

077 The use of importance weights to correct for sampling bias is a standard technique for many machine
 078 learning problems (*e.g.*, [9, 10, 11]) including active learning [12, 13, 5]. Our algorithm is based
 079 on the importance weighted active learning (IWAL) framework introduced by [5]. In that work, a
 080 rejection threshold procedure called *loss-weighting* is rigorously analyzed and shown to yield im-
 081 proved label complexity bounds in certain cases. Loss-weighting is more general than our technique
 082 in that it extends beyond zero-one loss to a certain subclass of loss functions such as logistic loss. On
 083 the other hand, the loss-weighting rejection threshold requires optimizing over a restricted version
 084 space, which is computationally undesirable. Moreover, the label complexity bound given in [5]
 085 only applies to hypotheses selected from this version space, and not when selected from the entire
 086 hypothesis class (as the general IWAL framework suggests). We avoid these deficiencies using a
 087 new rejection threshold procedure and a more subtle martingale analysis.

088 Many of the previously mentioned algorithms are analyzed in the agnostic learning model, where
 089 no assumption is made about the noise distribution (see also [14]). In this setting, the label com-
 090 plexity of active learning algorithms cannot generally improve over supervised learners by more
 091 than a constant factor [15, 5]. However, under a parameterization of the noise distribution related to
 092 Tsybakov’s low-noise condition [16], active learning algorithms have been shown to have improved
 093 label complexity bounds over what is achievable in the purely agnostic setting [17, 8, 18, 6, 7]. We
 094 also consider this parameterization to obtain a tighter label complexity analysis.

096 2 Preliminaries

098 2.1 Learning Model

099
 100 Let \mathcal{D} be a distribution over $\mathcal{X} \times \mathcal{Y}$ where \mathcal{X} is the input space and $\mathcal{Y} = \{\pm 1\}$ are the labels. Let
 101 $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ be a pair of random variables with joint distribution \mathcal{D} . An active learner receives
 102 a sequence $(X_1, Y_1), (X_2, Y_2), \dots$ of i.i.d. copies of (X, Y) , with the label Y_i hidden unless it is
 103 explicitly queried. We use the shorthand $a_{1:k}$ to denote a sequence (a_1, a_2, \dots, a_k) (so $k = 0$
 104 correspond to the empty sequence).

105 Let \mathcal{H} be a set of hypotheses mapping from \mathcal{X} to \mathcal{Y} . For simplicity, we assume \mathcal{H} is finite but does
 106 not completely agree on any single $x \in \mathcal{X}$ (*i.e.*, $\forall x \in \mathcal{X}, \exists h, h' \in \mathcal{H}$ such that $h(x) \neq h'(x)$). This
 107 keeps the focus on the relevant aspects of active learning that differ from passive learning. The error
 of a hypothesis $h : \mathcal{X} \rightarrow \mathcal{Y}$ is $\text{err}(h) := \Pr(h(X) \neq Y)$. Let $h^* := \arg \min\{\text{err}(h) : h \in \mathcal{H}\}$ be

a hypothesis of minimum error in \mathcal{H} . The goal of the active learner is to return a hypothesis $h \in \mathcal{H}$ with error $\text{err}(h)$ not much more than $\text{err}(h^*)$, using as few label queries as possible.

2.2 Importance Weighted Active Learning

In the importance weighted active learning (IWAL) framework of [5], an active learner looks at the unlabeled data X_1, X_2, \dots one at a time. After each new point X_i , the learner determines a probability $P_i \in [0, 1]$. Then a coin with bias P_i is flipped, and the label Y_i is queried if and only if the coin comes up heads. The query probability P_i can depend on all previous unlabeled examples $X_{1:i-1}$, any previously queried labels, any past coin flips, and the current unlabeled point X_i .

Formally, an IWAL algorithm specifies a *rejection threshold* function $p : (\mathcal{X} \times \mathcal{Y} \times \{0, 1\})^* \times \mathcal{X} \rightarrow [0, 1]$ for determining these query probabilities. Let $Q_i \in \{0, 1\}$ be a random variable conditionally independent of the current label Y_i ,

$$Q_i \perp\!\!\!\perp Y_i \mid X_{1:i}, Y_{1:i-1}, Q_{1:i-1}$$

and with conditional expectation

$$\mathbb{E}[Q_i \mid Z_{1:i-1}, X_i] = P_i := p(Z_{1:i-1}, X_i).$$

where $Z_j := (X_j, Y_j, Q_j)$. That is, Q_i indicates if the label Y_i is queried (the outcome of the coin toss). Although the notation does not explicitly suggest this, the query probability $P_i = p(Z_{1:i-1}, X_i)$ is allowed to explicitly depend on a label Y_j ($j < i$) if and only if it has been queried ($Q_j = 1$).

2.3 Importance Weighted Estimators

We first review some standard facts about the importance weighting technique. For a function $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, define the *importance weighted estimator* of $\mathbb{E}[f(X, Y)]$ from $Z_{1:n} \in (\mathcal{X} \times \mathcal{Y} \times \{0, 1\})^n$ to be

$$\hat{f}(Z_{1:n}) := \frac{1}{n} \sum_{i=1}^n \frac{Q_i}{P_i} \cdot f(X_i, Y_i).$$

Note that this quantity depends on a label Y_i only if it has been queried (*i.e.*, only if $Q_i = 1$; it also depends on X_i only if $Q_i = 1$). Our rejection threshold will be based on a specialization of this estimator, specifically the *importance weighted empirical error* of a hypothesis h

$$\text{err}(h, Z_{1:n}) := \frac{1}{n} \sum_{i=1}^n \frac{Q_i}{P_i} \cdot \mathbb{1}[h(X_i) \neq Y_i].$$

In the notation of Algorithm 1, this is equivalent to

$$\text{err}(h, S_n) := \frac{1}{n} \sum_{(X_i, Y_i, 1/P_i) \in S_n} (1/P_i) \cdot \mathbb{1}[h(X_i) \neq Y_i] \quad (1)$$

where $S_n \subseteq \mathcal{X} \times \mathcal{Y} \times \mathbb{R}$ is the importance weighted sample collected by the algorithm.

A basic property of these estimators is *unbiasedness*: $\mathbb{E}[\hat{f}(Z_{1:n})] = (1/n) \sum_{i=1}^n \mathbb{E}[\mathbb{E}[(Q_i/P_i) \cdot f(X_i, Y_i) \mid X_{1:i}, Y_{1:i-1}, Q_{1:i-1}]] = (1/n) \sum_{i=1}^n \mathbb{E}[(P_i/P_i) \cdot f(X_i, Y_i)] = \mathbb{E}[f(X, Y)]$. So, for example, the importance weighted empirical error of a hypothesis h is an unbiased estimator of its true error $\text{err}(h)$. This holds for *any* choice of the rejection threshold that guarantees $P_i > 0$.

3 A Deviation Bound for Importance Weighted Estimators

As mentioned before, the rejection threshold used by our algorithm is based on importance weighted error estimates $\text{err}(h, Z_{1:n})$. Even though these estimates are unbiased, they are only reliable when the variance is not too large. To get a handle on this, we need a deviation bound for importance weighted estimators. This is complicated by two factors that rules out straightforward applications of some standard bounds:

1. The importance weighted samples $(X_i, Y_i, 1/P_i)$ (or equivalently, the $Z_i = (X_i, Y_i, Q_i)$) are not i.i.d. This is because the query probability P_i (and thus the importance weight $1/P_i$) generally depends on $Z_{1:i-1}$ and X_i .
2. The effective range and variance of each term in the estimator are, themselves, random variables.

To address these issues, we develop a deviation bound using a martingale technique from [19].

Let $f : \mathcal{X} \times \mathcal{Y} \rightarrow [-1, 1]$ be a bounded function. Consider any rejection threshold function $p : (\mathcal{X} \times \mathcal{Y} \times \{0, 1\})^* \times \mathcal{X} \rightarrow (0, 1]$ for which $P_n = p(Z_{1:n-1}, X_n)$ is bounded below by some positive quantity (which may depend on n). Equivalently, the query probabilities P_n should have inverses $1/P_n$ bounded above by some deterministic quantity r_{max} (which, again, may depend on n). The *a priori* upper bound r_{max} on $1/P_n$ can be pessimistic, as the dependence on r_{max} in the final deviation bound will be very mild—it enters in as $\log \log r_{max}$. Our goal is to prove a bound on $|\hat{f}(Z_{1:n}) - \mathbb{E}[f(X, Y)]|$ that holds with high probability over the joint distribution of $Z_{1:n}$.

To start, we establish bounds on the range and variance of each term $W_i := (Q_i/P_i) \cdot f(X_i, Y_i)$ in the estimator, conditioned on $(X_{1:i}, Y_{1:i}, Q_{1:i-1})$. Let $\mathbb{E}_i[\cdot]$ denote $\mathbb{E}[\cdot | X_{1:i}, Y_{1:i}, Q_{1:i-1}]$. Note that $\mathbb{E}_i[W_i] = (\mathbb{E}_i[Q_i]/P_i) \cdot f(X_i, Y_i) = f(X_i, Y_i)$, so if $\mathbb{E}_i[W_i] = 0$, then $W_i = 0$. Therefore, the (conditional) range and variance are non-zero only if $\mathbb{E}_i[W_i] \neq 0$. For the range, we have $|W_i| = (Q_i/P_i) \cdot |f(X_i, Y_i)| \leq 1/P_i$, and for the variance, $\mathbb{E}_i[(W_i - \mathbb{E}_i[W_i])^2] \leq (\mathbb{E}_i[Q_i^2]/P_i^2) \cdot f(X_i, Y_i)^2 \leq 1/P_i$. These range and variance bounds indicate the form of the deviations we can expect, similar to that of other classical deviation bounds.

Theorem 1. *Pick any $t \geq 0$ and $n \geq 1$. Assume $1 \leq 1/P_i \leq r_{max}$ for all $1 \leq i \leq n$, and let $R_n := 1/\min(\{P_i : 1 \leq i \leq n \wedge f(X_i, Y_i) \neq 0\} \cup \{1\})$. With probability at least $1 - 2(3 + \log_2 r_{max})e^{-t/2}$,*

$$\left| \frac{1}{n} \sum_{i=1}^n \frac{Q_i}{P_i} \cdot f(X_i, Y_i) - \mathbb{E}[f(X, Y)] \right| \leq \sqrt{\frac{2R_n t}{n}} + \sqrt{\frac{2t}{n}} + \frac{R_n t}{3n}.$$

We defer all proofs to the appendices.

4 Algorithm

First, we state a deviation bound for the importance weighted error of hypotheses in a finite hypothesis class \mathcal{H} that holds for all $n \geq 1$. It is a simple consequence of Theorem 1 and union bounds; the form of the bound motivates certain algorithmic choices to be described below.

Lemma 1. *Pick any $\delta \in (0, 1)$. For all $n \geq 1$, let*

$$\varepsilon_n := \frac{16 \log(2(3 + n \log_2 n)n(n+1)|\mathcal{H}|/\delta)}{n} = O\left(\frac{\log(n|\mathcal{H}|/\delta)}{n}\right). \quad (3)$$

Let $(Z_1, Z_2, \dots) \in (\mathcal{X} \times \mathcal{Y} \times \{0, 1\})^$ be the sequence of random variables specified in Section 2.2 using a rejection threshold $p : (\mathcal{X} \times \mathcal{Y} \times \{0, 1\})^* \times \mathcal{X} \rightarrow [0, 1]$ that satisfies $p(z_{1:n}, x) \geq 1/n^n$ for all $(z_{1:n}, x) \in (\mathcal{X} \times \mathcal{Y} \times \{0, 1\})^n \times \mathcal{X}$ and all $n \geq 1$.*

The following holds with probability at least $1 - \delta$. For all $n \geq 1$ and all $h \in \mathcal{H}$,

$$|(\text{err}(h, Z_{1:n}) - \text{err}(h^*, Z_{1:n})) - (\text{err}(h) - \text{err}(h^*))| \leq \sqrt{\frac{\varepsilon_n}{P_{min,n}(h)}} + \frac{\varepsilon_n}{P_{min,n}(h)} \quad (4)$$

where $P_{min,n}(h) = \min\{P_i : 1 \leq i \leq n \wedge h(X_i) \neq h^(X_i)\} \cup \{1\}$.*

We let $C_0 = O(\log(|\mathcal{H}|/\delta)) \geq 2$ be a quantity such that ε_n (as defined in Eq. (3)) is bounded as $\varepsilon_n \leq C_0 \cdot \log(n+1)/n$. The following absolute constants are used in the description of the rejection threshold and the subsequent analysis: $c_1 := 5 + 2\sqrt{2}$, $c_2 := 5$, $c_3 := ((c_1 + \sqrt{2})/(c_1 - 2))^2$, $c_4 := (c_1 + \sqrt{c_3})^2$, $c_5 := c_2 + c_3$.

Our proposed algorithm is shown in Figure 1. The rejection threshold (Step 2) is based on the deviation bound from Lemma 1. First, the importance weighted error minimizing hypothesis h_k and

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

Algorithm 1
Notes: see Eq. (1) for the definition of err (importance weighted error), and Section 4 for the definitions of C_0 , c_1 , and c_2 .
Initialize: $S_0 := \emptyset$.
For $k = 1, 2, \dots, n$:
1. Obtain unlabeled data point X_k .
2. Let
 $h_k := \arg \min\{\text{err}(h, S_{k-1}) : h \in \mathcal{H}\}$, and
 $h'_k := \arg \min\{\text{err}(h, S_{k-1}) : h \in \mathcal{H} \wedge h(X_k) \neq h_k(X_k)\}$.
Let $G_k := \text{err}(h'_k, S_{k-1}) - \text{err}(h_k, S_{k-1})$, and

$$P_k := \begin{cases} 1 & \text{if } G_k \leq \sqrt{\frac{C_0 \log k}{k-1}} + \frac{C_0 \log k}{k-1} \\ s & \text{otherwise} \end{cases} \quad \left(= \min \left\{ 1, O \left(\frac{1}{G_k^2} + \frac{1}{G_k} \right) \cdot \frac{C_0 \log k}{k-1} \right\} \right)$$
where $s \in (0, 1)$ is the positive solution to the equation

$$G_k = \left(\frac{c_1}{\sqrt{s}} - c_1 + 1 \right) \cdot \sqrt{\frac{C_0 \log k}{k-1}} + \left(\frac{c_2}{s} - c_2 + 1 \right) \cdot \frac{C_0 \log k}{k-1}. \quad (2)$$
3. Toss a biased coin with $\Pr(\text{heads}) = P_k$.
If heads, then query Y_k , and let $S_k := S_{k-1} \cup \{(X_k, Y_k, 1/P_k)\}$.
Else, let $S_k := S_{k-1}$.
Return: $h_{n+1} := \arg \min\{\text{err}(h, S_n) : h \in \mathcal{H}\}$.

Figure 1: Algorithm for importance weighted active learning with an error minimization oracle.

the “alternative” hypothesis h'_k are found. Note that both optimizations are over the entire hypothesis class \mathcal{H} (with h'_k only being required to disagree with h_k on x_k)—this is a key aspect where our algorithm differs from previous approaches. The difference in importance weighted errors G_k of the two hypotheses is then computed. If $G_k \leq \sqrt{(C_0 \log k)/(k-1)} + (C_0 \log k)/(k-1)$, then the query probability P_k is set to 1. Otherwise, P_k is set to the positive solution s to the quadratic equation in Eq. (2). The functional form of P_k is roughly

$$\min \left\{ 1, O \left(\frac{1}{G_k^2} + \frac{1}{G_k} \right) \cdot \frac{C_0 \log k}{k-1} \right\}.$$

It can be checked that $P_k \in (0, 1]$ and that P_k is non-increasing with G_k . It is also useful to note that $(\log k)/(k-1)$ is monotonically decreasing with $k \geq 1$ (we use the convention $\log(1)/0 = \infty$).

In order to apply Lemma 1 with our rejection threshold, we need to establish the (very crude) bound $P_k \geq 1/k^k$ for all k .

Lemma 2. *The rejection threshold of Algorithm 1 satisfies $p(z_{1:n-1}, x) \geq 1/n^n$ for all $n \geq 1$ and all $(z_{1:n-1}, x) \in (\mathcal{X} \times \mathcal{Y} \times \{0, 1\})^{n-1} \times \mathcal{X}$.*

Note that this is a worst-case bound; our analysis shows that the probabilities P_k are more like $1/\text{poly}(k)$ in the typical case.

5 Analysis

5.1 Correctness

We first prove a consistency guarantee for Algorithm 1 that bounds the generalization error of the importance weighted empirical error minimizer. The proof actually establishes a lower bound on the query probabilities $P_i \geq 1/2$ for X_i such that $h_n(X_i) \neq h^*(X_i)$. This offers an intuitive characterization of the weighting landscape induced by the importance weights $1/P_i$.

Theorem 2. *The following holds with probability at least $1 - \delta$. For any $n \geq 1$,*

$$0 \leq \text{err}(h_n) - \text{err}(h^*) \leq \text{err}(h_n, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1}) + \sqrt{\frac{2C_0 \log n}{n-1}} + \frac{2C_0 \log n}{n-1}.$$

This implies, for all $n \geq 1$,

$$\text{err}(h_n) \leq \text{err}(h^*) + \sqrt{\frac{2C_0 \log n}{n-1}} + \frac{2C_0 \log n}{n-1}.$$

Therefore, the final hypothesis returned by Algorithm 1 after seeing n unlabeled data has roughly the same error bound as a hypothesis returned by a standard passive learner with n labeled data. A variant of this result under certain noise conditions is given in the appendix.

5.2 Label Complexity Analysis

We now bound the number of labels requested by Algorithm 1 after n iterations. The following lemma bounds the probability of querying the label Y_n ; this is subsequently used to establish the final bound on the expected number of labels queried. The key to the proof is in relating empirical error differences and their deviations to the probability of querying a label. This is mediated through the *disagreement coefficient*, a quantity first used by [14] for analyzing the label complexity of the A^2 algorithm of [3]. The disagreement coefficient $\theta := \theta(h^*, \mathcal{H}, \mathcal{D})$ is defined as

$$\theta(h^*, \mathcal{H}, \mathcal{D}) := \sup \left\{ \frac{\Pr(X \in \text{DIS}(h^*, r))}{r} : r > 0 \right\}$$

where

$$\text{DIS}(h^*, r) := \{x \in \mathcal{X} : \exists h' \in \mathcal{H} \text{ such that } \Pr(h^*(X) \neq h'(X)) \leq r \text{ and } h^*(x) \neq h'(x)\}$$

(the disagreement region around h^* at radius r). This quantity is bounded for many learning problems studied in the literature; see [14, 6, 20, 21] for more discussion. Note that the supremum can instead be taken over $r > \epsilon$ if the target excess error is ϵ , which allows for a more detailed analysis.

Lemma 3. *Assume the bounds from Eq. (4) holds for all $h \in \mathcal{H}$ and $n \geq 1$. For any $n \geq 1$,*

$$\mathbb{E}[Q_n] \leq \theta \cdot 2 \text{err}(h^*) + O \left(\theta \cdot \sqrt{\frac{C_0 \log n}{n-1}} + \theta \cdot \frac{C_0 \log^2 n}{n-1} \right).$$

Theorem 3. *With probability at least $1 - \delta$, the expected number of labels queried by Algorithm 1 after n iterations is at most*

$$1 + \theta \cdot 2 \text{err}(h^*) \cdot (n-1) + O \left(\theta \cdot \sqrt{C_0 n \log n} + \theta \cdot C_0 \log^3 n \right).$$

Proof. Follows from assuming Y_1 is always queried; applying Lemmas 1, 2, 3, and linearity of expectation. \square

The bound is dominated by a linear term scaled by $\text{err}(h^*)$, plus a sublinear term. The linear term $\text{err}(h^*) \cdot n$ is unavoidable in the worst case, as evident from label complexity lower bounds [15, 5]. When $\text{err}(h^*)$ is negligible (e.g., the data is separable) and θ is bounded (as is the case for many problems studied in the literature [14]), then the bound represents a polynomial label complexity improvement over supervised learning, similar to that achieved by the version space algorithm from [5].

5.3 Analysis under Low Noise Conditions

Some recent work on active learning has focused on improved label complexity under certain noise conditions [17, 8, 18, 6, 7]. Specifically, it is assumed that there exists constants $\kappa > 0$ and $0 < \alpha \leq 1$ such that

$$\Pr(h(X) \neq h^*(X)) \leq \kappa \cdot (\text{err}(h) - \text{err}(h^*))^\alpha \quad (5)$$

for all $h \in \mathcal{H}$. This is related to Tsybakov's low noise condition [16]. Essentially, this condition requires that low error hypotheses not be too far from the optimal hypothesis h^* under the disagreement metric $\Pr(h^*(X) \neq h(X))$. Under this condition, Lemma 3 can be improved, which in turn yields the following theorem.

Theorem 4. Assume that for some value of $\kappa > 0$ and $0 < \alpha \leq 1$, the condition in Eq. (5) holds for all $h \in \mathcal{H}$. There is a constant $c_\alpha > 0$ depending only on α such that the following holds. With probability at least $1 - \delta$, the expected number of labels queried by Algorithm 1 after n iterations is at most

$$\theta \cdot \kappa \cdot c_\alpha \cdot (C_0 \log n)^{\alpha/2} \cdot n^{1-\alpha/2}.$$

Note that the bound is sublinear in n for all $0 < \alpha \leq 1$, which implies label complexity improvements whenever θ is bounded (an improved analogue of Theorem 2 under these conditions can be established using similar techniques). The previous algorithms of [6, 7] obtain even better rates under these noise conditions using specialized data dependent generalization bounds, but these algorithms also required optimizations over restricted version spaces, even for the bound computation.

6 Experiments

Although agnostic learning is typically intractable in the worst case, empirical risk minimization can serve as a useful abstraction for many practical supervised learning algorithms in non-worst case scenarios. With this in mind, we conducted a preliminary experimental evaluation of Algorithm 1, implemented using a popular algorithm for learning decision trees in place of the required ERM oracle. Specifically, we use the J48 algorithm from Weka v3.6.2 (with default parameters) to select the hypothesis h_k in each round k ; to produce the “alternative” hypothesis h'_k , we just modify the decision tree h_k by changing the label of the node used for predicting on x_k . Both of these procedures are clearly heuristic, but they are similar in spirit to the required optimizations. We set $C_0 = 8$ and $c_1 = c_2 = 1$ —these can be regarded as tuning parameters, with C_0 controlling the aggressiveness of the rejection threshold. We did not perform parameter tuning with active learning although the importance weighting approach developed here could potentially be used for that. Rather, the goal of these experiments is to assess the compatibility of Algorithm 1 with an existing, practical supervised learning procedure.

6.1 Data Sets

We constructed two binary classification tasks using MNIST and KDDCUP99 data sets. For MNIST, we randomly chose 4000 training 3s and 5s for training (using the 3s as the positive class), and used all of the 1902 testing 3s and 5s for testing. For KDDCUP99, we randomly chose 5000 examples for training, and another 5000 for testing. In both cases, we reduced the dimension of the data to 25 using PCA.

To demonstrate the versatility of our algorithm, we also conducted a multi-class classification experiment using the entire MNIST data set (all ten digits, so 60000 training data and 10000 testing data). This required modifying how h'_k is selected: we force $h'_k(x_k) \neq h_k(x_k)$ by changing the label of the prediction node for x_k to the next best label. We used PCA to reduce the dimension to 40.

6.2 Results

We examined the test error as a function of (i) the number of unlabeled data seen, and (ii) the number of labels queried. We compared the performance of the active learner described above to a passive learner (one that queries every label, so (i) and (ii) are the same) using J48 with default parameters.

In all three cases, the test errors as a function of the number of unlabeled data were roughly the same for both the active and passive learners. This agrees with the consistency guarantee from Theorem 2. We note that this is a basic property *not* satisfied by many active learning algorithms (this issue is discussed further in [22]).

In terms of test error as a function of the number of labels queried (Figure 2), the active learner had minimal improvement over the passive learner on the binary MNIST task, but a substantial improvement over the passive learner on the KDDCUP99 task (even at small numbers of label queries). For the multi-class MNIST task, the active learner had a moderate improvement over the passive learner. Note that KDDCUP99 is far less noisy (more separable) than MNIST 3s vs 5s task, so the results are in line with the label complexity behavior suggested by Theorem 3, which states that the label complexity improvement may scale with the error of the optimal hypothesis. Also,

378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404
 405
 406
 407
 408
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431

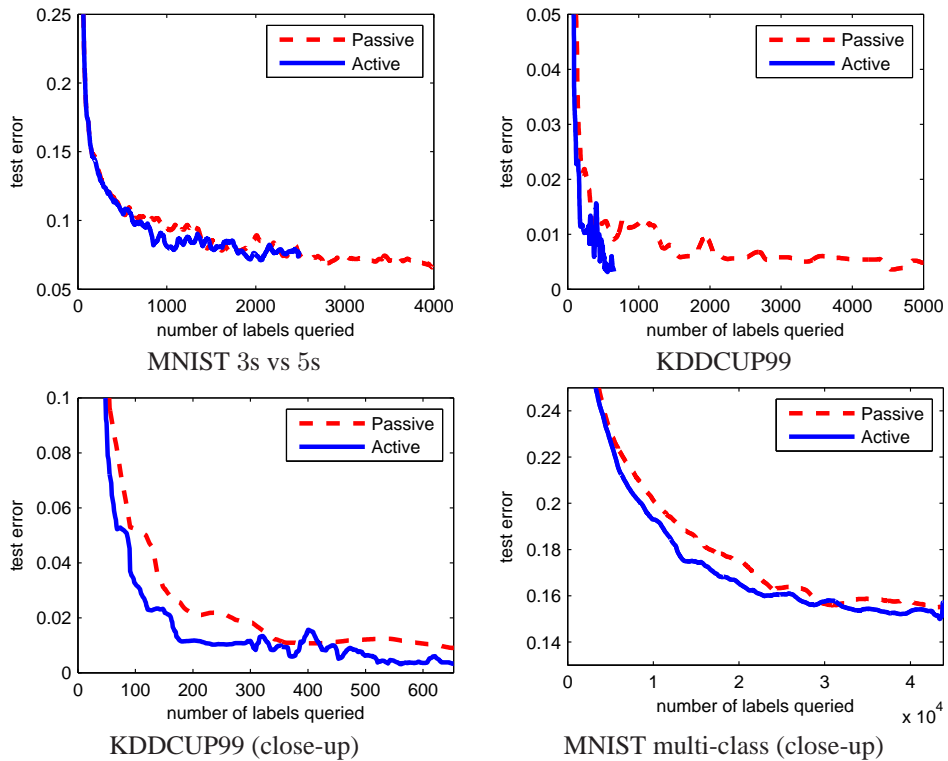


Figure 2: Test errors as a function of the number of labels queried.

the results from MNIST tasks suggest that the active learner may require an initial random sampling phase during which it is equivalent to the passive learner, and the advantage manifests itself after this phase. This again is consistent with the analysis (also see [14]), as the disagreement coefficient can be large at initial scales, yet much smaller as the number of (unlabeled) data increases and the scale becomes finer.

7 Conclusion

This paper provides a new active learning algorithm based on error minimization oracles, a departure from the version space approach adopted by previous works. The algorithm we introduce here motivates computationally tractable and effective methods for active learning with many classifier training algorithms. The overall algorithmic template applies to any training algorithm that (i) operates by approximate error minimization and (ii) for which the cost of switching a class prediction (as measured by example errors) can be estimated. Furthermore, although these properties might only hold in an approximate or heuristic sense, the created active learning algorithm will be “safe” in the sense that it will eventually converge to the same solution as a passive supervised learning algorithm. Consequently, we believe this approach can be widely used to reduce the cost of labeling in situations where labeling is expensive.

Recent theoretical work on active learning has focused on improving rates of convergence. However, in some applications, it may be desirable to improve performance at much smaller sample sizes, perhaps even at the cost of improved rates as long as consistency is ensured. Importance sampling and weighting techniques like those analyzed in this work may be useful for developing more aggressive strategies with such properties.

References

[1] D. Cohn, L. Atlas, and R. Ladner. Improving generalization with active learning. *Machine Learning*, 15(2):201–221, 1994.

- 432 [2] S. Dasgupta. Coarse sample complexity bounds for active learning. In *Advances in Neural Information*
433 *Processing Systems 18*, 2005.
- 434 [3] M.-F. Balcan, A. Beygelzimer, and J. Langford. Agnostic active learning. In *Twenty-Third International*
435 *Conference on Machine Learning*, 2006.
- 436 [4] S. Dasgupta, D. Hsu, and C. Monteleoni. A general agnostic active learning algorithm. In *Advances in*
437 *Neural Information Processing Systems 20*, 2007.
- 438 [5] A. Beygelzimer, S. Dasgupta, and J. Langford. Importance weighted active learning. In *Twenty-Sixth*
439 *International Conference on Machine Learning*, 2009.
- 440 [6] S. Hanneke. Adaptive rates of convergence in active learning. In *Twenty-Second Annual Conference on*
441 *Learning Theory*, 2009.
- 442 [7] V. Koltchinskii. Rademacher complexities and bounding the excess risk in active learning. Manuscript,
443 2009.
- 444 [8] M.-F. Balcan, A. Broder, and T. Zhang. Margin based active learning. In *Twentieth Annual Conference*
445 *on Learning Theory*, 2007.
- 446 [9] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- 447 [10] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem.
448 *SIAM Journal of Computing*, 32:48–77, 2002.
- 449 [11] M. Sugiyama, M. Krauledat, and K.-R. Müller. Covariate shift adaptation by importance weighted cross
450 validation. *Journal of Machine Learning Research*, 8:985–1005, 2007.
- 451 [12] M. Sugiyama. Active learning for misspecified models. In *Advances in Neural Information Processing*
452 *Systems 18*, 2005.
- 453 [13] F. Bach. Active learning for misspecified generalized linear models. In *Advances in Neural Information*
454 *Processing Systems 19*, 2006.
- 455 [14] S. Hanneke. A bound on the label complexity of agnostic active learning. In *Twenty-Fourth International*
456 *Conference on Machine Learning*, 2007.
- 457 [15] M. Kääriäinen. Active learning in the non-realizable case. In *Seventeenth International Conference on*
458 *Algorithmic Learning Theory*, 2006.
- 459 [16] A. B. Tsybakov. Optimal aggregation of classifiers in statistical learning. *Annals of Statistics*, 32(1):135–
460 166, 2004.
- 461 [17] R. Castro and R. Nowak. Upper and lower bounds for active learning. In *Allerton Conference on Com-*
462 *munication, Control and Computing*, 2006.
- 463 [18] R. Castro and R. Nowak. Minimax bounds for active learning. In *Twentieth Annual Conference on*
464 *Learning Theory*, 2007.
- 465 [19] T. Zhang. Data dependent concentration bounds for sequential prediction algorithms. In *Eighteenth*
466 *Annual Conference on Learning Theory*, 2005.
- 467 [20] E. Friedman. Active learning for smooth problems. In *Twenty-Second Annual Conference on Learning*
468 *Theory*, 2009.
- 469 [21] L. Wang. Sufficient conditions for agnostic active learnable. In *Advances in Neural Information Process-*
470 *ing Systems 22*, 2009.
- 471 [22] S. Dasgupta and D. Hsu. Hierarchical sampling for active learning. In *Twenty-Fifth International Confer-*
472 *ence on Machine Learning*, 2008.
- 473
- 474
- 475
- 476
- 477
- 478
- 479
- 480
- 481
- 482
- 483
- 484
- 485

A Proof of Deviation Bound for Importance Weighted Estimators

The techniques here are mostly developed in [19]; for completeness, we detail the proofs for our particular application. The first two lemmas establish a basic bound in terms of conditional moment generating functions.

Lemma 4. For all $n \geq 1$ and all functionals $\Xi_i := \xi_i(Z_{1:i})$,

$$\mathbb{E} \left[\exp \left(\sum_{i=1}^n \Xi_i - \sum_{i=1}^n \ln \mathbb{E}_i[\exp(\Xi_i)] \right) \right] = 1.$$

Proof. A straightforward induction on n . □

Lemma 5. For all $t \geq 0$, $\lambda \in \mathbb{R}$, $n \geq 1$, and functionals $\Xi_i := \xi_i(Z_{1:i})$,

$$\Pr \left(\lambda \sum_{i=1}^n \Xi_i - \sum_{i=1}^n \ln \mathbb{E}_i[\exp(\lambda \Xi_i)] \geq t \right) \leq e^{-t}.$$

Proof. The claim follows by Markov's inequality and Lemma 4 (replacing Ξ_i with $\lambda \Xi_i$). □

In order to specialize Lemma 5 for our purposes, we first analyze the conditional moment generating function of $W_i - \mathbb{E}_i[W_i]$.

Lemma 6. If $0 < \lambda < 3P_i$, then

$$\ln \mathbb{E}_i[\exp(\lambda(W_i - \mathbb{E}_i[W_i]))] \leq \frac{1}{P_i} \cdot \frac{\lambda^2}{2(1 - \lambda/(3P_i))}.$$

If $\mathbb{E}_i[W_i] = 0$, then

$$\ln \mathbb{E}_i[\exp(\lambda(W_i - \mathbb{E}_i[W_i]))] = 0.$$

Proof. Let $g(x) := (\exp(x) - x - 1)/x^2$ for $x \neq 0$, so $\exp(x) = 1 + x + x^2 \cdot g(x)$. Note that $g(x)$ is non-decreasing. Thus,

$$\begin{aligned} & \mathbb{E}_i[\exp(\lambda(W_i - \mathbb{E}_i[W_i]))] \\ &= \mathbb{E}_i[1 + \lambda(W_i - \mathbb{E}_i[W_i]) + \lambda^2(W_i - \mathbb{E}_i[W_i])^2 \cdot g(\lambda(W_i - \mathbb{E}_i[W_i]))] \\ &= 1 + \lambda^2 \cdot \mathbb{E}_i[(W_i - \mathbb{E}_i[W_i])^2 \cdot g(\lambda(W_i - \mathbb{E}_i[W_i]))] \\ &\leq 1 + \lambda^2 \cdot \mathbb{E}_i[(W_i - \mathbb{E}_i[W_i])^2 \cdot g(\lambda/P_i)] \\ &= 1 + \lambda^2 \cdot \mathbb{E}_i[(W_i - \mathbb{E}_i[W_i])^2] \cdot g(\lambda/P_i) \\ &\leq 1 + (\lambda^2/P_i) \cdot g(\lambda/P_i) \end{aligned}$$

where the first inequality follows from the range bound $|W_i| \leq 1/P_i$ and the second follows from variance bound $\mathbb{E}_i[(W_i - \mathbb{E}_i[W_i])^2] \leq 1/P_i$. Now the first claim follows from the definition of $g(x)$, the facts $\exp(x) - x - 1 \leq x^2/(2(1 - x/3))$ for $0 \leq x < 3$ and $\ln(1 + x) \leq x$.

The second claim is immediate from the definition of W_i and the fact $\mathbb{E}_i[W_i] = f(X_i, Y_i)$. □

We now combine Lemma 6 and Lemma 5 to bound the deviation of the importance weighted estimator $\hat{f}(Z_{1:n})$ from $(1/n) \sum_{i=1}^n \mathbb{E}_i[W_i]$.

Lemma 7. Pick any $t \geq 0$, $n \geq 1$, and $p_{\min} > 0$, and let E be the (joint) event

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n W_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_i[W_i] \geq \sqrt{\frac{1}{p_{\min}} \cdot \frac{2t}{n}} + \frac{1}{p_{\min}} \cdot \frac{t}{3n} \\ & \text{and } \min\{P_i : 1 \leq i \leq n \wedge \mathbb{E}_i[W_i] \neq 0\} \geq p_{\min}. \end{aligned}$$

Then $\Pr(E) \leq e^{-t}$.

540 *Proof.* With foresight, let

$$541 \lambda := 3p_{\min} \cdot \frac{\sqrt{\frac{1}{3p_{\min}} \cdot \frac{2t}{3n}}}{1 + \sqrt{\frac{1}{3p_{\min}} \cdot \frac{2t}{3n}}}. \quad 542$$

543 Note that $0 < \lambda < 3p_{\min}$. By Lemma 6 and the choice of λ , we have that if $\min\{P_i : 1 \leq i \leq n \wedge \mathbb{E}_i[W_i] \neq 0\} \geq p_{\min}$, then

$$544 \frac{1}{n\lambda} \cdot \sum_{i=1}^n \ln \mathbb{E}_i[\exp(\lambda(W_i - \mathbb{E}_i[W_i]))] \leq \frac{1}{p_{\min}} \cdot \frac{\lambda}{2(1 - \lambda/(3p_{\min}))} = \sqrt{\frac{1}{p_{\min}} \cdot \frac{t}{2n}} \quad (6)$$

545 and

$$546 \frac{t}{n\lambda} = \sqrt{\frac{1}{p_{\min}} \cdot \frac{t}{2n}} + \frac{1}{p_{\min}} \cdot \frac{t}{3n}. \quad (7)$$

547 Let E' be the event that

$$548 \frac{1}{n} \cdot \sum_{i=1}^n (W_i - \mathbb{E}_i[W_i]) - \frac{1}{n\lambda} \cdot \sum_{i=1}^n \ln \mathbb{E}_i[\exp(\lambda(W_i - \mathbb{E}_i[W_i]))] \geq \frac{t}{n\lambda}$$

549 and let E'' be the event $\min\{P_i : 1 \leq i \leq n \wedge \mathbb{E}_i[W_i] \neq 0\} \geq p_{\min}$. Together, Eq. (6) and Eq. (7) imply $E \subseteq E' \cap E''$. And of course, $E' \cap E'' \subseteq E'$, so $\Pr(E) \leq \Pr(E' \cap E'') \leq \Pr(E') \leq e^{-t}$ by Lemma 5. \square

550 To do away with the joint event in Lemma 7, we use the standard trick of taking a union bound over a geometric sequence of possible values for p_{\min} .

551 **Lemma 8.** Pick any $t \geq 0$ and $n \geq 1$. Assume $1 \leq 1/P_i \leq r_{\max}$ for all $1 \leq i \leq n$, and let $R_n := 1/\min\{P_i : 1 \leq i \leq n \wedge \mathbb{E}_i[W_i] \neq 0\} \cup \{1\}$. We have

$$552 \Pr \left(\left| \frac{1}{n} \sum_{i=1}^n W_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_i[W_i] \right| \geq \sqrt{\frac{2R_n t}{n}} + \frac{R_n t}{3n} \right) \leq 2(2 + \log_2 r_{\max})e^{-t/2}.$$

553 *Proof.* The assumption on P_i implies $1 \leq R_n \leq r_{\max}$. Let $r_j := 2^j$ for $-1 \leq j \leq m := \lceil \log_2 r_{\max} \rceil$. Then

$$554 \Pr \left(\frac{1}{n} \sum_{i=1}^n W_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_i[W_i] \geq \sqrt{\frac{2R_n t}{n}} + \frac{R_n t}{3n} \right)$$

$$555 = \sum_{j=0}^m \Pr \left(\frac{1}{n} \sum_{i=1}^n W_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_i[W_i] \geq \sqrt{\frac{2R_n t}{n}} + \frac{R_n t}{3n} \wedge r_{j-1} < R_n \leq r_j \right)$$

$$556 \leq \sum_{j=0}^m \Pr \left(\frac{1}{n} \sum_{i=1}^n W_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_i[W_i] \geq \sqrt{\frac{2r_{j-1} t}{n}} + \frac{r_{j-1} t}{3n} \wedge R_n \leq r_j \right)$$

$$557 = \sum_{j=0}^m \Pr \left(\frac{1}{n} \sum_{i=1}^n W_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_i[W_i] \geq \sqrt{\frac{2r_j(t/2)}{n}} + \frac{r_j(t/2)}{3n} \wedge R_n \leq r_j \right)$$

$$558 \leq (2 + \log_2 r_{\max})e^{-t/2}$$

559 where the last inequality follows from Lemma 7. Replacing W_i with $-W_i$ bounds the probability of deviations in the other direction in exactly the same way. The claim then follows by the union bound. \square

560 *Proof of Theorem 1.* By Hoeffding's inequality and the fact $|f(X_i, Y_i)| \leq 1$, we have

$$561 \Pr \left(\left| \frac{1}{n} \sum_{i=1}^n f(X_i, Y_i) - \mathbb{E}[f(X, Y)] \right| \geq \sqrt{\frac{2t}{n}} \right) \leq 2e^{-t/2}.$$

562 Since $\mathbb{E}_i[W_i] = f(X_i, Y_i)$, the claim follows by combining this and Lemma 8 with the triangle inequality and the union bound. \square

B Remaining Proofs

In this section, we use the notation $\varepsilon_k := C_0 \log(k+1)/k$.

B.1 Proof of Lemma 2

By induction on n . Trivial for $n = 1$ (since $p(\text{empty sequence}, x) = 1$ for all $x \in \mathcal{X}$), so now fix any $n \geq 2$ and assume as the inductive hypothesis $p_{n-1} = p(z_{1:n-2}, x) \geq 1/(n-1)^{n-1}$ for all $(z_{1:n-2}, x) \in (\mathcal{X} \times \mathcal{Y}) \times \{0, 1\}^{n-2} \times \mathcal{X}$. Fix any $(z_{1:n-1}, x) \in (\mathcal{X} \times \mathcal{Y}) \times \{0, 1\}^{n-1} \times \mathcal{X}$, and consider the error difference $g_n := \text{err}(h'_n, z_{1:n-1}) - \text{err}(h_n, z_{1:n-1})$ used to determine $p_n := p(z_{1:n-1}, x)$. We only have to consider the case $g_n > \sqrt{\varepsilon_{n-1}} + \varepsilon_{n-1}$. By the inductive hypothesis and triangle inequality, we have $g_n \leq 2(n-1)^{n-1}$. Solving the quadratic in Eq. (2) implies

$$\begin{aligned}
\sqrt{p_n} &= \frac{c_1 \cdot \sqrt{\varepsilon_{n-1}} + \sqrt{c_1^2 \cdot \varepsilon_{n-1} + 4 \cdot (g_n + (c_1 - 1) \cdot \sqrt{\varepsilon_{n-1}} + (c_2 - 1) \cdot \varepsilon_{n-1}) \cdot c_2 \cdot \varepsilon_{n-1}}}{2(g_n + (c_1 - 1) \cdot \sqrt{\varepsilon_{n-1}} + (c_2 - 1) \cdot \varepsilon_{n-1})} \\
&> \frac{\sqrt{4 \cdot (g_n + (c_1 - 1) \cdot \sqrt{\varepsilon_{n-1}} + (c_2 - 1) \cdot \varepsilon_{n-1}) \cdot c_2 \cdot \varepsilon_{n-1}}}{2(g_n + (c_1 - 1) \cdot \sqrt{\varepsilon_{n-1}} + (c_2 - 1) \cdot \varepsilon_{n-1})} \quad (\text{dropping terms}) \\
&= \sqrt{\frac{c_2 \cdot \varepsilon_{n-1}}{g_n + (c_1 - 1) \cdot \sqrt{\varepsilon_{n-1}} + (c_2 - 1) \cdot \varepsilon_{n-1}}} \\
&\geq \sqrt{\frac{c_2 \cdot \varepsilon_{n-1}}{g_n + (c_1 - 1) \cdot \sqrt{\varepsilon_{n-1}} + (c_1 - 1) \cdot \varepsilon_{n-1}}} \quad (\text{since } c_2 \leq c_1) \\
&\geq \sqrt{\frac{c_2 \cdot \varepsilon_{n-1}}{c_1 \cdot g_n}} \quad (\text{since } g_n > \sqrt{\varepsilon_{n-1}} + \varepsilon_{n-1}) \\
&= \sqrt{\frac{c_2 \cdot C_0 \log n}{c_1 \cdot (n-1) \cdot g_n}} \\
&\geq \sqrt{\frac{c_2 \cdot C_0 \log n}{2c_1 \cdot (n-1) \cdot (n-1)^{n-1}}} \quad (\text{inductive hypothesis}) \\
&> \sqrt{\frac{1}{e(n-1)^n}} \quad (\text{since } C_0 \geq 2, n \geq 2, \text{ and } (c_2 \cdot C_0 \log 2)/(2c_1) > 1/e) \\
&\geq \sqrt{\frac{1}{n^n}} \quad (\text{since } (n/(n-1))^n \geq e)
\end{aligned}$$

as required. \square

B.2 Proof of Theorem 2

We condition on the $1 - \delta$ probability event that the deviation bounds from Lemma 1 hold (also using Lemma 2). The proof now proceeds by induction on n . The claim is trivially true for $n = 1$. Now pick any $n \geq 2$ and assume as the (strong) inductive hypothesis that

$$0 \leq \text{err}(h_k) - \text{err}(h^*) \leq \text{err}(h_k, Z_{1:k-1}) - \text{err}(h^*, Z_{1:k-1}) + \sqrt{2\varepsilon_{k-1}} + 2\varepsilon_{k-1} \quad (8)$$

for all $1 \leq k \leq n-1$. We need to show Eq. (8) holds for $k = n$.

Let $P_{\min} := \min\{P_i : 1 \leq i \leq n-1 \wedge h_n(X_i) \neq h^*(X_i)\} \cup \{1\}$. If $P_{\min} \geq 1/2$, then Eq. (4) implies that Eq. (8) holds for $k = n$ as needed. So assume for sake of contradiction that $P_{\min} < 1/2$, and let $n_0 := \max\{i \leq n-1 : P_i = P_{\min} \wedge h_n(X_i) \neq h^*(X_i)\}$. By definition of P_{n_0} , we have

$$\text{err}(h'_{n_0}, Z_{1:n_0-1}) - \text{err}(h_{n_0}, Z_{1:n_0-1}) = \left(\frac{c_1}{\sqrt{P_{\min}}} - c_1 + 1 \right) \sqrt{\varepsilon_{n_0-1}} + \left(\frac{c_2}{P_{\min}} - c_2 + 1 \right) \varepsilon_{n_0-1}.$$

Using this fact together with the inductive hypothesis, we have

$$\begin{aligned}
& \text{err}(h'_{n_0}, Z_{1:n_0-1}) - \text{err}(h^*, Z_{1:n_0-1}) \\
&= \text{err}(h'_{n_0}, Z_{1:n_0-1}) - \text{err}(h_{n_0}, Z_{1:n_0-1}) + \text{err}(h_{n_0}, Z_{1:n_0-1}) - \text{err}(h^*, Z_{1:n_0-1}) \\
&\geq \left(\frac{c_1}{\sqrt{P_{\min}}} - c_1 + 1 \right) \cdot \sqrt{\varepsilon_{n_0-1}} + \left(\frac{c_2}{P_{\min}} - c_2 + 1 \right) \cdot \varepsilon_{n_0-1} - \sqrt{2\varepsilon_{n_0-1}} - 2\varepsilon_{n_0-1} \\
&= \left(\frac{c_1}{\sqrt{P_{\min}}} - c_1 + 1 - \sqrt{2} \right) \cdot \sqrt{\varepsilon_{n_0-1}} + \left(\frac{c_2}{P_{\min}} - c_2 - 1 \right) \cdot \varepsilon_{n_0-1} \quad . \quad (9)
\end{aligned}$$

We use the assumption $P_{\min} < 1/2$ to lower bound the righthand side to get the inequality

$$\text{err}(h'_{n_0}, Z_{1:n_0-1}) - \text{err}(h^*, Z_{1:n_0-1}) > (c_1 - 1) \cdot (\sqrt{2} - 1) \cdot \sqrt{\varepsilon_{n_0-1}} + (c_2 - 1) \cdot \varepsilon_{n_0-1} > 0.$$

which implies $\text{err}(h'_{n_0}, Z_{1:n_0-1}) > \text{err}(h^*, Z_{1:n_0-1})$. Since h'_{n_0} minimizes $\text{err}(h, Z_{1:n_0-1})$ among hypotheses $h \in \mathcal{H}$ that disagree with h_{n_0} on X_{n_0} , it must be that h^* agrees with h_{n_0} on X_{n_0} . By transitivity and the definition of n_0 , we conclude that $h_n(X_{n_0}) = h'_{n_0}(X_{n_0})$; so $\text{err}(h_n, Z_{1:n_0-1}) \geq \text{err}(h'_{n_0}, Z_{1:n_0-1})$. Then

$$\begin{aligned}
& \text{err}(h_n, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1}) \\
&\geq \text{err}(h_n) - \text{err}(h^*) - \sqrt{\frac{1}{P_{\min}} \cdot \varepsilon_{n-1}} - \frac{1}{P_{\min}} \cdot \varepsilon_{n-1} \\
&\geq \text{err}(h_n, Z_{1:n_0-1}) - \text{err}(h^*, Z_{1:n_0-1}) - 2 \cdot \sqrt{\frac{1}{P_{\min}} \cdot \varepsilon_{n_0-1}} - 2 \cdot \frac{1}{P_{\min}} \cdot \varepsilon_{n_0-1} \\
&\geq \left(\frac{c_1 - 2}{\sqrt{P_{\min}}} - c_1 + 1 - \sqrt{2} \right) \cdot \sqrt{\varepsilon_{n_0-1}} + \left(\frac{c_2 - 2}{P_{\min}} - c_2 - 1 \right) \cdot \varepsilon_{n_0-1} \\
&> \left((c_1 - 1) \cdot (\sqrt{2} - 1) - 2\sqrt{2} \right) \cdot \sqrt{\varepsilon_{n_0-1}} + (c_2 - 5) \cdot \varepsilon_{n_0-1}
\end{aligned}$$

where Eq. (4) is used in the first two inequalities, Eq. (9) and the fact $\text{err}(h_n, Z_{1:n_0-1}) \geq \text{err}(h'_{n_0}, Z_{1:n_0-1})$ are used in the third inequality, and the fact $P_{\min} < 1/2$ is used in the last inequality. This final quantity is non-negative, so we have the contradiction $\text{err}(h_n, Z_{1:n-1}) > \text{err}(h^*, Z_{1:n-1})$. \square

B.3 Proof of Lemma 3

First, we establish a property of the query probabilities that relates error deviations (via P_{\min}) to empirical error differences (via P_n). Both quantities play essential roles in bounding the label complexity through the disagreement metric structure around h^* .

Lemma 9. *Assume the bounds from Eq. (4) hold for all $h \in \mathcal{H}$ and $n \geq 1$. For any $n \geq 1$, we have $P_n \leq c_3 \cdot P_{\min}$, where $P_{\min} := \min(\{P_i : 1 \leq i \leq n-1 \wedge h(X_i) \neq h^*(X_i)\} \cup \{1\})$ and*

$$h := \begin{cases} h_n & \text{if } h_n \text{ disagrees with } h^* \text{ on } X_n \\ h'_n & \text{if } h'_n \text{ disagrees with } h^* \text{ on } X_n. \end{cases} \quad (10)$$

Proof. We can assume $P_{\min} < 1/c_3$, since otherwise the claim is trivial. Pick any $n_0 \leq n-1$ such that $h(X_{n_0}) \neq h^*(X_{n_0})$ and $P_{n_0} = P_{\min}$ (such an n_0 is guaranteed to exist given the above assumption). We now proceed as in the proof of Theorem 2. We first show a lower bound on $\text{err}(h, Z_{1:n_0-1}) - \text{err}(h^*, Z_{1:n_0-1})$. Note that

$$\begin{aligned}
& \text{err}(h'_{n_0}, Z_{1:n_0-1}) - \text{err}(h^*, Z_{1:n_0-1}) \\
&= \text{err}(h'_{n_0}, Z_{1:n_0-1}) - \text{err}(h_{n_0}, Z_{1:n_0-1}) + \text{err}(h_{n_0}, Z_{1:n_0-1}) - \text{err}(h^*, Z_{1:n_0-1}) \\
&\geq \left(\frac{c_1}{\sqrt{P_{\min}}} - c_1 + 1 \right) \cdot \sqrt{\varepsilon_{n_0-1}} + \left(\frac{c_2}{P_{\min}} - c_2 + 1 \right) \cdot \varepsilon_{n_0-1} - \sqrt{2\varepsilon_{n_0-1}} - 2\varepsilon_{n_0-1} \\
&= \left(\frac{c_1}{\sqrt{P_{\min}}} - c_1 + 1 - \sqrt{2} \right) \cdot \sqrt{\varepsilon_{n_0-1}} + \left(\frac{c_2}{P_{\min}} - c_2 - 1 \right) \cdot \varepsilon_{n_0-1} \quad (11)
\end{aligned}$$

702 where the inequality follows from Theorem 2. The righthand side is positive, so h^* must disagree
703 with h'_{n_0} on X_{n_0} . By transitivity (recalling that $h(X_{n_0}) \neq h^*(X_{n_0})$), h must agree with h'_{n_0}
704 on X_{n_0} . Therefore $\text{err}(h, Z_{1:n_0-1}) - \text{err}(h'_{n_0}, Z_{1:n_0-1}) \geq 0$, so the inequality in Eq. (11) holds with
705 h in place of h'_{n_0} on the lefthand side.

706 Now $\text{err}(h, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1})$ is related to $\text{err}(h, Z_{1:n_0-1}) - \text{err}(h^*, Z_{1:n_0-1})$ through
707 $\text{err}(h) - \text{err}(h^*)$ using the deviation bound from Eq. (4) (as well as the fact $\varepsilon_{n_0-1} \geq \varepsilon_{n-1}$):

$$\begin{aligned}
709 & \text{err}(h, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1}) \\
710 & \geq \text{err}(h, Z_{1:n_0-1}) - \text{err}(h^*, Z_{1:n_0-1}) - 2 \cdot \sqrt{\frac{1}{P_{\min}} \cdot \varepsilon_{n_0-1}} - 2 \cdot \frac{1}{P_{\min}} \cdot \varepsilon_{n_0-1} \\
711 & \geq \left(\frac{c_1 - 2}{\sqrt{P_{\min}}} - c_1 + 1 - \sqrt{2} \right) \cdot \sqrt{\varepsilon_{n-1}} + \left(\frac{c_2 - 2}{P_{\min}} - c_2 - 1 \right) \cdot \varepsilon_{n-1} > 0. \quad (12)
\end{aligned}$$

715 If $h = h_n$, then $\text{err}(h, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1}) = \text{err}(h_n, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1}) \leq 0$ by the
716 minimality of $\text{err}(h_n, Z_{1:n-1})$; this contradicts Eq. (12). Therefore it must be that $h = h'_n$. In this
717 case,

$$\begin{aligned}
718 & \text{err}(h, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1}) \leq \text{err}(h'_n, Z_{1:n-1}) - \text{err}(h_n, Z_{1:n-1}) \\
719 & = \left(\frac{c_1}{\sqrt{P_n}} - c_1 + 1 \right) \cdot \sqrt{\varepsilon_{n-1}} + \left(\frac{c_2}{P_n} - c_2 + 1 \right) \cdot \varepsilon_{n-1} \quad (13)
\end{aligned}$$

722 where the inequality follows from the minimality of $\text{err}(h_n, Z_{1:n-1})$, and the subsequent step fol-
723 lows from the definition of P_n . Combining the lower bound in Eq. (12) and the upper bound in
724 Eq. (13) implies that

$$\frac{c_1}{\sqrt{P_n}} \cdot \sqrt{\varepsilon_{n-1}} + \frac{c_2}{P_n} \cdot \varepsilon_{n-1} \geq \left(\frac{c_1 - 2}{\sqrt{P_{\min}}} - \sqrt{2} \right) \cdot \sqrt{\varepsilon_{n-1}} + \left(\frac{c_2 - 2}{P_{\min}} - 2 \right) \cdot \varepsilon_{n-1}.$$

728 It is easily checked that this implies $P_n \leq c_3 \cdot P_{\min}$. □

730 *Proof of Lemma 3.* Define h as in Eq. (10). By Lemma 9, we have $\min(\{P_i : 1 \leq i \leq n-1 \wedge$
731 $h(X_i) \neq h^*(X_i)\} \cup \{1\}) \geq P_n/c_3$. We first show that

$$\begin{aligned}
732 & \text{err}(h) - \text{err}(h^*) \leq \text{err}(h, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1}) + \sqrt{\frac{c_3}{P_n} \cdot \varepsilon_{n-1}} + \frac{c_3}{P_n} \cdot \varepsilon_{n-1} \\
733 & \leq \sqrt{\frac{c_4}{P_n}} \cdot \sqrt{\varepsilon_{n-1}} + \frac{c_5}{P_n} \cdot \varepsilon_{n-1}. \quad (14)
\end{aligned}$$

738 The first inequality follows from Eq. (4) and Lemma 9. For the second inequality, we consider two
739 cases depending on h . If $h = h'_n$, then we bound $\text{err}(h, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1})$ from above by
740 $\text{err}(h'_n, Z_{1:n-1}) - \text{err}(h_n, Z_{1:n-1})$ (by definition of h and minimality of $\text{err}(h_n, Z_{1:n-1})$), and then
741 simplify

$$\begin{aligned}
742 & \text{err}(h'_n, Z_{1:n-1}) - \text{err}(h_n, Z_{1:n-1}) + \sqrt{\frac{c_3}{P_n} \cdot \varepsilon_{n-1}} + \frac{c_3}{P_n} \cdot \varepsilon_{n-1} \\
743 & \leq \left(\frac{c_1 + \sqrt{c_3}}{\sqrt{P_n}} - c_1 + 1 \right) \cdot \sqrt{\varepsilon_{n-1}} + \left(\frac{c_2 + c_3}{P_n} - c_2 + 1 \right) \cdot \varepsilon_{n-1} \leq \sqrt{\frac{c_4}{P_n}} \cdot \sqrt{\varepsilon_{n-1}} + \frac{c_5}{P_n} \cdot \varepsilon_{n-1}
\end{aligned}$$

748 using the definition of P_n and the facts $c_1 \geq 1$ and $c_2 \geq 1$. If instead $h = h_n$, then we use the facts
749 $\text{err}(h, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1}) = \text{err}(h_n, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1}) \leq 0$ and $c_3 \leq \min\{c_4, c_5\}$.

750 If $\text{err}(h) - \text{err}(h^*) = \gamma > 0$, then solving the quadratic inequality in Eq. (14) for P_n gives the
751 bound

$$P_n \leq \min \left\{ 1, \frac{3}{2} \cdot \left(\frac{c_4}{\gamma^2} + \frac{c_5}{\gamma} \right) \cdot \varepsilon_{n-1} \right\}.$$

754 If $\text{err}(h) - \text{err}(h^*) \leq \bar{\gamma}$, then by the triangle inequality we have

$$\Pr(h^*(X) \neq h(X)) \leq \text{err}(h^*) + \text{err}(h) \leq 2\text{err}(h^*) + \bar{\gamma}$$

756 which in turn implies $X_n \in \text{DIS}(h^*, 2 \text{err}(h^*) + \bar{\gamma})$. Note that $\Pr(X_n \in \text{DIS}(h^*, 2 \text{err}(h^*) + \bar{\gamma})) \leq$
 757 $\theta \cdot (2 \text{err}(h^*) + \bar{\gamma})$ by definition of θ , so $\Pr(\text{err}(h) - \text{err}(h^*) \leq \bar{\gamma}) \leq \theta \cdot (2 \text{err}(h^*) + \bar{\gamma})$.

758 Let $f(\gamma) := \partial \Pr(\text{err}(h) - \text{err}(h^*) \leq \gamma) / \partial \gamma$ be the probability density (mass) function of the
 759 error difference $\text{err}(h) - \text{err}(h^*)$; note that this error difference is a function of $(Z_{1:n-1}, X_n)$. We
 760 compute the expected value of Q_n by conditioning on $\text{err}(h) - \text{err}(h^*)$ and integrating (an upper
 761 bound on) $\mathbb{E}[Q_n | \text{err}(h) - \text{err}(h^*) = \gamma]$ with respect to $f(\gamma)$.

762 Let $\gamma_0 > 0$ be the positive solution to $1.5(c_4/\gamma^2 + c_5/\gamma)\varepsilon_{n-1} = 1$. It can be checked that $\gamma_0 >$
 763 $\sqrt{1.5c_4\varepsilon_{n-1}}$. We have

$$\begin{aligned}
 764 \mathbb{E}[Q_n] &= \mathbb{E}[\mathbb{E}[Q_n | Z_{1:n-1}, X_n]] \quad (\text{the outer expectation is over } (Z_{1:n-1}, X_n)) \\
 765 &= \int_0^1 \left(\frac{\partial}{\partial \gamma} \Pr(\text{err}(h) - \text{err}(h^*) \leq \gamma) \right) \cdot \mathbb{E}[Q_n | \text{err}(h) - \text{err}(h^*) = \gamma] \cdot d\gamma \\
 766 &\leq \int_0^1 \left(\frac{\partial}{\partial \gamma} \Pr(\text{err}(h) - \text{err}(h^*) \leq \gamma) \right) \cdot \min \left\{ 1, \frac{3}{2} \cdot \left(\frac{c_4}{\gamma^2} + \frac{c_5}{\gamma} \right) \cdot \varepsilon_{n-1} \right\} \cdot d\gamma \\
 767 &\leq \frac{3}{2} \cdot (c_4 + c_5) \cdot \varepsilon_{n-1} \cdot \Pr(\text{err}(h) - \text{err}(h^*) \leq 1) \\
 768 &\quad - \int_0^1 \left(\frac{\partial}{\partial \gamma} \min \left\{ 1, \frac{3}{2} \cdot \left(\frac{c_4}{\gamma^2} + \frac{c_5}{\gamma} \right) \cdot \varepsilon_{n-1} \right\} \right) \cdot \Pr(\text{err}(h) - \text{err}(h^*) \leq \gamma) \cdot d\gamma \\
 769 &\leq \frac{3}{2} \cdot (c_4 + c_5) \cdot \varepsilon_{n-1} + \int_{\gamma_0}^1 \frac{3}{2} \cdot \left(\frac{2c_4}{\gamma^3} + \frac{c_5}{\gamma^2} \right) \cdot \varepsilon_{n-1} \cdot \theta \cdot (2 \text{err}(h^*) + \gamma) \cdot d\gamma \\
 770 &= \frac{3}{2} \cdot (c_4 + c_5) \cdot \varepsilon_{n-1} + \theta \cdot 2 \text{err}(h^*) \cdot \frac{3}{2} \cdot \left(c_4 \left(\frac{1}{\gamma_0^2} - 1 \right) + c_5 \left(\frac{1}{\gamma_0} - 1 \right) \right) \cdot \varepsilon_{n-1} \\
 771 &\quad + \theta \cdot \frac{3}{2} \cdot \left(2c_4 \left(\frac{1}{\gamma_0} - 1 \right) + c_5 \ln \frac{1}{\gamma_0} \right) \cdot \varepsilon_{n-1} \\
 772 &\leq \frac{3}{2} \cdot (c_4 + c_5) \cdot \varepsilon_{n-1} + \theta \cdot 2 \text{err}(h^*) + \theta \cdot \sqrt{6c_4\varepsilon_{n-1}} + \theta \cdot \frac{3c_5}{4} \cdot \varepsilon_{n-1} \cdot \ln \frac{1}{1.5c_4\varepsilon_{n-1}}
 \end{aligned}$$

773 where the first inequality uses the bound on $\mathbb{E}[Q_n | \text{err}(h) - \text{err}(h^*) = \gamma]$; the second inequality
 774 uses integration-by-parts; the third inequality uses the fact that the integrand from the previous line
 775 is 0 for $0 \leq \gamma \leq \gamma_0$, as well as the bound on $\Pr(\text{err}(h) - \text{err}(h^*) \leq \gamma)$; and the fourth inequality
 776 uses the definition of γ_0 . \square

777 B.4 Proof of Theorem 4

778 The theorem is a simple consequence of the following analogue of Lemma 3.

779 **Lemma 10.** *Assume that for some value of $\kappa > 0$ and $0 < \alpha \leq 1$, the condition in Eq. (5) holds*
 780 *for all $h \in \mathcal{H}$. Assume the bounds from Eq. (4) holds for all $h \in \mathcal{H}$ and $n \geq 1$. There is a constant*
 781 *$c_\alpha > 0$ such that the following holds. For any $n \geq 1$,*

$$782 \mathbb{E}[Q_n] \leq \theta \cdot \kappa \cdot c_\alpha \cdot \left(\frac{C_0 \log n}{n-1} \right)^{\alpha/2}.$$

783 *Proof.* For the most part, the proof is the same as that of Lemma 3. The key difference is to use the
 784 noise condition in Eq. (5) to directly bound $\Pr(h(X) \neq h^*(X)) \leq \kappa \cdot (\text{err}(h) - \text{err}(h^*))^\alpha$, which
 785 in turn implies the bound $\Pr(\text{err}(h) - \text{err}(h^*) \leq \gamma) \leq \theta \kappa \gamma^\alpha$. As before, let $\gamma_0 > \sqrt{1.5c_4\varepsilon_{n-1}}$ be
 786 the solution to $1.5(c_4/\gamma^2 + c_5/\gamma)\varepsilon_{n-1} = 1$. First consider the case $\alpha < 1$. Then, the expectation of
 787 Q_n can be bounded as

$$\begin{aligned}
 788 \mathbb{E}[Q_n] &\leq \frac{3}{2} \cdot (c_4 + c_5) \cdot \varepsilon_{n-1} + \int_{\gamma_0}^1 \frac{3}{2} \cdot \left(\frac{2c_4}{\gamma^3} + \frac{c_5}{\gamma^2} \right) \cdot \varepsilon_{n-1} \cdot \Pr(\text{err}(h) - \text{err}(h^*) \leq \gamma) \cdot d\gamma \\
 789 &\leq \frac{3}{2} \cdot (c_4 + c_5) \cdot \varepsilon_{n-1} + \int_{\gamma_0}^1 \frac{3}{2} \cdot \left(\frac{2c_4}{\gamma^3} + \frac{c_5}{\gamma^2} \right) \cdot \varepsilon_{n-1} \cdot \theta \kappa \gamma^\alpha \cdot d\gamma \\
 790 &\leq \frac{3}{2} \cdot (c_4 + c_5) \cdot \varepsilon_{n-1} + \theta \cdot \kappa \cdot \frac{3}{2} \cdot \left(\frac{2c_4}{2-\alpha} \cdot \frac{1}{\gamma_0^{2-\alpha}} + \frac{c_5}{1-\alpha} \cdot \frac{1}{\gamma_0^{1-\alpha}} \right) \cdot \varepsilon_{n-1}.
 \end{aligned}$$

810 The case $\alpha = 1$ is handled similarly. □
 811

812 B.5 Analogue of Theorem 2 under Low Noise Conditions

813 We first state a variant of Lemma 1 that takes into account the probability of disagreement between
 814 a hypothesis h and the optimal hypothesis h^* .

815 **Lemma 11.** *There exists an absolute constant $c > 0$ such that the following holds. Pick any*
 816 *$\delta \in (0, 1)$. For all $n \geq 1$, let*

$$817 \varepsilon_n := \frac{c \cdot \log((n+1)|\mathcal{H}|/\delta)}{n}.$$

818 *Let $(Z_1, Z_2, \dots) \in (\mathcal{X} \times \mathcal{Y} \times \{0, 1\})^n$ be the sequence of random variables specified in Section 2.2*
 819 *using a rejection threshold $p : (\mathcal{X} \times \mathcal{Y} \times \{0, 1\})^n \times \mathcal{X} \rightarrow [0, 1]$ that satisfies $p(z_{1:n}, x) \geq 1/n^n$ for*
 820 *all $(z_{1:n}, x) \in (\mathcal{X} \times \mathcal{Y} \times \{0, 1\})^n \times \mathcal{X}$ and all $n \geq 1$.*

821 *The following holds with probability at least $1 - \delta$. For all $n \geq 1$ and all $h \in \mathcal{H}$,*

$$822 |(\text{err}(h, Z_{1:n}) - \text{err}(h^*, Z_{1:n})) - (\text{err}(h) - \text{err}(h^*))| \leq \sqrt{\frac{\Pr(h(X) \neq h^*(X))}{P_{\min, n}(h)}} \cdot \varepsilon_n + \frac{\varepsilon_n}{P_{\min, n}(h)}$$

823 where $P_{\min, n}(h) = \min\{P_i : 1 \leq i \leq n \wedge h(X_i) \neq h^*(X_i)\} \cup \{1\}$.

824 *Proof sketch.* The proof of this lemma follows along the same lines as that of Lemma 1. A key
 825 difference comes in Lemma 7: the joint event is modified to also conjoin with

$$826 \frac{1}{n} \sum_{i=1}^n \mathbb{1}(\mathbb{E}_i[f(X_i, Y_i)] \leq 0) \leq a$$

827 for some fixed $a > 0$. In the proof, the parameter λ should be chosen as

$$828 \lambda := 3p_{\min} \cdot \frac{\sqrt{\frac{1}{3p_{\min}} \cdot \frac{2at}{3n}}}{a + \sqrt{\frac{1}{3p_{\min}} \cdot \frac{2at}{3n}}}.$$

829 Lemma 8 is modified to also take a union bound over a sequence of possible values for a (in fact,
 830 only $n + 1$ different values need to be considered). Finally, instead of combining with Hoeffding's
 831 inequality, we use Bernstein's inequality (or a multiplicative form of Chernoff's bound) so the re-
 832 sulting bound (an analogue of Theorem 1) involves an empirical average inside the square-root term:
 833 with probability at least $1 - O(n \cdot \log_2 r_{\max})e^{-t/2}$,

$$834 \left| \frac{1}{n} \sum_{i=1}^n \frac{Q_i}{P_i} \cdot f(X_i, Y_i) - \mathbb{E}[f(X, Y)] \right| \leq O\left(\sqrt{\frac{R_n A_n t}{n}} + \frac{R_n t}{3n}\right)$$

835 where

$$836 A_n := \frac{1}{n} \sum_{i=1}^n \mathbb{1}(f(X_i, Y_i) \neq 0).$$

837 Finally, we apply this deviation bound to obtain uniform error bounds over all hypotheses \mathcal{H} (a
 838 few extra steps are required to replace the empirical quantity A_n in the bound with a distributional
 839 quantity). □

840 Using the previous lemma, a modified version of Theorem 2 follows from essentially the same proof.
 841 We note that the quantity $C_1 := O(\log(|\mathcal{H}|/\delta))$ used here may differ from C_0 by constant factors.

842 **Lemma 12.** *The following holds with probability at least $1 - \delta$. For any $n \geq 1$,*

$$843 0 \leq \text{err}(h_n) - \text{err}(h^*) \leq \text{err}(h_n, Z_{1:n-1}) - \text{err}(h^*, Z_{1:n-1})$$

$$844 + \sqrt{\frac{2 \Pr(h_n(X) \neq h^*(X)) C_1 \log n}{n-1}} + \frac{2C_1 \log n}{n-1}.$$

845 *This implies, for all $n \geq 1$,*

$$846 \text{err}(h_n) \leq \text{err}(h^*) + \sqrt{\frac{2 \Pr(h_n(X) \neq h^*(X)) C_1 \log n}{n-1}} + \frac{2C_0 \log n}{n-1}.$$

864 Finally, using the noise condition to bound $\Pr(h_n(X) \neq h^*(X)) \leq \kappa \cdot (\text{err}(h_n) - \text{err}(h^*))^\alpha$, we
865 obtain the final error bound.

866 **Theorem 5.** *The following holds with probability at least $1 - \delta$. For any $n \geq 1$,*
867

$$868 \text{err}(h_n) \leq \text{err}(h^*) + c_\kappa \cdot \left(\frac{C_1 \log n}{n-1} \right)^{\frac{1}{2-\alpha}}$$

869 where c_κ is a constant that depends only on κ .
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917