

Learning with a Drifting Target Concept

Steve Hanneke¹✉, Varun Kanade², and Liu Yang³

¹ Princeton, NJ, USA

`steve.hanneke@gmail.com`

² Département d'informatique, École normale supérieure, Paris, France

`varun.kanade@ens.fr`

³ IBM T.J. Watson Research Center, Yorktown Heights, NY, USA

`yangli@us.ibm.com`

Abstract. We study the problem of learning in the presence of a drifting target concept. Specifically, we provide bounds on the error rate at a given time, given a learner with access to a history of independent samples labeled according to a target concept that can change on each round. One of our main contributions is a refinement of the best previous results for polynomial-time algorithms for the space of linear separators under a uniform distribution. We also provide general results for an algorithm capable of adapting to a variable rate of drift of the target concept. Some of the results also describe an active learning variant of this setting, and provide bounds on the number of queries for the labels of points in the sequence sufficient to obtain the stated bounds on the error rates.

1 Introduction

Much of the work on statistical learning has focused on learning settings in which the concept to be learned is static over time. However, there are many application areas where this is not the case. For instance, in the problem of face recognition, the concept to be learned actually changes over time as each individual's facial features evolve over time. In this work, we study the problem of learning with a drifting target concept. Specifically, we consider a statistical learning setting, in which data arrive i.i.d. in a stream, and for each data point, the learner is required to predict a label for the data point at that time. We are then interested in obtaining low error rates for these predictions. The target labels are generated from a function known to reside in a given concept space, and at each time t the target function is allowed to change by at most some distance Δ_t : that is, the probability the new target function disagrees with the previous target function on a random sample is at most Δ_t .

This framework has previously been studied in a number of articles. The classic works of [5, 6, 15, 16] and [7] together provide a general analysis of a very-much related setting. Though the objectives in these works are specified slightly differently, the results established there are easily translated into our present framework, and we summarize many of the relevant results in Section 3.

While the results in these classic works are general, the best guarantees on the error rates are only known for methods having no guarantees of computational efficiency. In a more recent effort, the work of [8] studies this problem in the specific context of learning a homogeneous linear separator, when all the Δ_t values are identical. They propose a polynomial-time algorithm (based on the modified Perceptron algorithm of [9]), and prove a bound on the number of mistakes it makes as a function of the number of samples, when the data distribution satisfies a certain condition called “ λ -good” (which generalizes a useful property of the uniform distribution on a sphere). However, their result is again worse than that obtainable by the known computationally-inefficient methods.

Thus, the natural question is whether there exists a polynomial-time algorithm achieving roughly the same guarantees on the error rates known for the inefficient methods. In the present work, we resolve this question in the case of learning homogeneous linear separators under the uniform distribution, by proposing a polynomial-time algorithm that indeed achieves roughly the same bounds on the error rates known for the inefficient methods in the literature. This represents the main technical contribution of this work.

We also study the interesting problem of *adaptivity* of an algorithm to the sequence of Δ_t values, in the setting where Δ_t may itself vary over time. Since the values Δ_t might typically not be accessible in practice, it seems important to have learning methods having no explicit dependence on the sequence Δ_t . We propose such a method below, and prove that it achieves roughly the same bounds on the error rates known for methods in the literature which require direct access to the Δ_t values. Also in the context of variable Δ_t sequences, we discuss conditions on the sequence Δ_t necessary and sufficient for there to exist a learning method guaranteeing *sublinear* growth of the number of mistakes.

We additionally study an *active learning* extension to this framework, in which, at each time, after making its prediction, the algorithm may decide whether or not to request access to the label assigned to the data point at that time. In addition to guarantees on the error rates (for *all* times, including those for which the label was not observed), we are also interested in bounding the number of labels we expect the algorithm to request, as a function of the number of samples encountered thus far.

2 Definitions and Notation

Formally, in this setting, there is a fixed distribution \mathcal{P} over the instance space \mathcal{X} , and there is a sequence of independent \mathcal{P} -distributed unlabeled data X_1, X_2, \dots . There is also a concept space \mathbb{C} , and a sequence of target functions $\mathbf{h}^* = \{h_1^*, h_2^*, \dots\}$ in \mathbb{C} . Each t has an associated target label $Y_t = h_t^*(X_t)$. In this context, a (passive) learning algorithm is required, on each round t , to produce a classifier \hat{h}_t based on the observations $(X_1, Y_1), \dots, (X_{t-1}, Y_{t-1})$, and we denote by $\hat{Y}_t = \hat{h}_t(X_t)$ the corresponding prediction by the algorithm for the label of X_t . For any classifier h , we define $\text{er}_t(h) = \mathcal{P}(x : h(x) \neq h_t^*(x))$.

We also say the algorithm makes a “mistake” on instance X_t if $\hat{Y}_t \neq Y_t$; thus, $\text{er}_t(\hat{h}_t) = \mathbb{P}(\hat{Y}_t \neq Y_t | (X_1, Y_1), \dots, (X_{t-1}, Y_{t-1}))$.

For notational convenience, we will suppose the h_t^* sequence is chosen independently from the X_t sequence (i.e., h_t^* is chosen prior to the “draw” of $X_1, X_2, \dots \sim \mathcal{P}$), and is not random. In each results, we will suppose \mathbf{h}^* is chosen from some set S of sequences in \mathbb{C} . In particular, we are interested in describing the sequence \mathbf{h}^* in terms of the magnitudes of *changes* in h_t^* from one time to the next. Specifically, for any sequence $\Delta = \{\Delta_t\}_{t=2}^\infty$ in $[0, 1]$, we denote by S_Δ the set of all sequences \mathbf{h}^* in \mathbb{C} such that, $\forall t \in \mathbb{N}$, $\mathcal{P}(x : h_t(x) \neq h_{t+1}(x)) \leq \Delta_{t+1}$. Throughout this article, we denote by d the VC dimension of \mathbb{C} [18], and we suppose $1 \leq d < \infty$. Also, $\forall x \in \mathbb{R}$, define $\text{Log}(x) = \ln(\max\{x, e\})$.

3 Background: (ϵ, S) -Tracking Algorithms

As mentioned, the classic literature on learning with a drifting target concept is expressed in terms of a slightly different model. In order to relate those results to our present setting, we first introduce the classic setting. Specifically, we consider a model introduced by [15], presented here in a more-general form inspired by [5]. For a set S of sequences $\{h_t\}_{t=1}^\infty$ in \mathbb{C} , and a value $\epsilon > 0$, an algorithm \mathcal{A} is said to be (ϵ, S) -tracking if $\exists t_\epsilon \in \mathbb{N}$ such that, for any choice of $\mathbf{h}^* \in S$, $\forall T \geq t_\epsilon$, the prediction \hat{Y}_T produced by \mathcal{A} at time T satisfies $\mathbb{P}(\hat{Y}_T \neq Y_T) \leq \epsilon$. Note that the value of this probability may be influenced by $\{X_t\}_{t=1}^T$, $\{h_t^*\}_{t=1}^T$, and any internal randomness of the algorithm \mathcal{A} .

The focus of the results expressed in this classical model is determining sufficient conditions on the set S for there to exist an (ϵ, S) -tracking algorithm, along with bounds on the sufficient size of t_ϵ . These conditions on S typically take the form of an assumption on the drift rate, expressed in terms of ϵ . Below, we summarize several of the strongest known results for this setting.

Bounded Drift Rate: The simplest, and perhaps most elegant, results for (ϵ, S) -tracking algorithms is for the set S of sequences with a bounded drift rate. Specifically, for any $\Delta \in [0, 1]$, define $S_\Delta = S_\Delta$, where Δ is such that $\Delta_{t+1} = \Delta$ for every $t \in \mathbb{N}$. The study of this problem was initiated in the original work of [15]. The best known general results are due to [16]: namely, that for some $\Delta_\epsilon = \Theta(\epsilon^2/d)$, for every $\epsilon \in (0, 1]$, there exists an (ϵ, S_Δ) -tracking algorithm for all values of $\Delta \leq \Delta_\epsilon$.¹ This refined an earlier result of [15] by a logarithmic factor. [16] further argued that this result can be achieved with $t_\epsilon = \Theta(d/\epsilon)$. The algorithm itself involves a beautiful modification of the one-inclusion graph prediction strategy of [14]; since its specification is somewhat involved, we refer the interested reader to the original work of [16] for the details.

¹ In fact, [16] also allowed the distribution \mathcal{P} to vary gradually over time. For simplicity, we will only discuss the case of fixed \mathcal{P} .

Varying Drift Rates (Nonadaptive Algorithm): In addition to the concrete bounds for the case $\mathbf{h}^* \in S_\Delta$, [15] additionally present an elegant general result. Specifically, they argue that, for any $\epsilon > 0$, and any $m = \Omega\left(\frac{d}{\epsilon} \text{Log} \frac{1}{\epsilon}\right)$, if $\sum_{i=1}^m \mathcal{P}(x : h_i^*(x) \neq h_{m+1}^*(x)) \leq m\epsilon/24$, then for $\hat{h} = \text{argmin}_{h \in \mathbb{C}} \sum_{i=1}^m \mathbb{1}[h(X_i) \neq Y_i]$, $\mathbb{P}(\hat{h}(X_{m+1}) \neq h_{m+1}^*(X_{m+1})) \leq \epsilon$. This result immediately inspires an algorithm \mathcal{A} which, at every time t , chooses a value $m_t \leq t-1$, and predicts $\hat{Y}_t = \hat{h}_t(X_t)$, for $\hat{h}_t = \text{argmin}_{h \in \mathbb{C}} \sum_{i=t-m_t}^{t-1} \mathbb{1}[h(X_i) \neq Y_i]$. We are then interested in choosing m_t to minimize the value of ϵ obtainable via the result of [15]. However, that method is based on the values $\mathcal{P}(x : h_i^*(x) \neq h_t^*(x))$, which would typically not be accessible to the algorithm. However, suppose instead we have access to a sequence Δ such that $\mathbf{h}^* \in S_\Delta$. In this case, we could approximate $\mathcal{P}(x : h_i^*(x) \neq h_t^*(x))$ by its *upper bound* $\sum_{j=i+1}^t \Delta_j$. In this case, we are interested choosing m_t to minimize the smallest value of ϵ such that $\sum_{i=t-m_t}^{t-1} \sum_{j=i+1}^t \Delta_j \leq m_t\epsilon/24$ and $m_t = \Omega\left(\frac{d}{\epsilon} \text{Log} \frac{1}{\epsilon}\right)$. One can easily verify that this minimum is obtained at a value

$$m_t = \Theta \left(\text{argmin}_{1 \leq m \leq t-1} \frac{1}{m} \sum_{i=t-m}^{t-1} \sum_{j=i+1}^t \Delta_j + \frac{d \text{Log}(m/d)}{m} \right),$$

and via the result of [15] (applied to X_{t-m_t}, \dots, X_t) the resulting algorithm has

$$\mathbb{P}(\hat{Y}_t \neq Y_t) \leq O \left(\min_{1 \leq m \leq t-1} \frac{1}{m} \sum_{i=t-m}^{t-1} \sum_{j=i+1}^t \Delta_j + \frac{d \text{Log}(m/d)}{m} \right). \quad (1)$$

As a special case, if every t has $\Delta_t = \Delta$ for a fixed value $\Delta \in [0, 1]$, this result recovers the bound $\sqrt{d\Delta \text{Log}(1/\Delta)}$, which is only slightly larger than the best bound of [16]. It also applies to far more general and more interesting sequences Δ , including some that allow periodic large jumps (i.e., $\Delta_t = 1$ for some indices t), others where the sequence Δ_t converges to 0, and so on. Note, however, that the algorithm obtaining this bound directly depends on the sequence Δ . One of the contributions of the present work is to remove this requirement, while maintaining essentially the same bound, though in a slightly different form.

Computational Efficiency: [15] also proposed a reduction-based approach, which sometimes yields computationally efficient methods, though the tolerable Δ value is smaller. Specifically, given any (randomized) polynomial-time algorithm \mathcal{A} that produces a classifier $h \in \mathbb{C}$ with $\sum_{t=1}^m \mathbb{1}[h(x_t) \neq y_t] = 0$ for any sequence $(x_1, y_1), \dots, (x_m, y_m)$ for which such a classifier h exists (called the *consistency problem*), they propose a polynomial-time algorithm that is (ϵ, S_Δ) -tracking for all values of $\Delta \leq \Delta'_\epsilon$, where $\Delta'_\epsilon = \Theta\left(\frac{\epsilon^2}{d^2 \text{Log}(1/\epsilon)}\right)$. This is slightly worse (by a factor of $d \text{Log}(1/\epsilon)$) than the drift rate tolerable by the (typically inefficient) algorithm mentioned above. However, it does sometimes yield computationally-efficient methods. For instance, there are known polynomial-time algorithms for the consistency problem for the classes of linear separators, conjunctions, and axis-aligned rectangles.

Lower Bounds: [15] additionally prove *lower bounds* for specific concept spaces: namely, linear separators and axis-aligned rectangles. They specifically argue that, for \mathbb{C} a concept space $\text{BASIC}_n = \{\cup_{i=1}^n [i/n, (i+a_i)/n) : \mathbf{a} \in [0, 1]^n\}$ on $[0, 1]$, under \mathcal{P} the uniform distribution on $[0, 1]$, for any $\epsilon \in [0, 1/e^2]$ and $\Delta_\epsilon \geq e^4 \epsilon^2/n$, for any algorithm \mathcal{A} , and any $T \in \mathbb{N}$, there exists a choice of $\mathbf{h}^* \in S_{\Delta_\epsilon}$ such that the prediction \hat{Y}_T produced by \mathcal{A} at time T satisfies $\mathbb{P}(\hat{Y}_T \neq Y_T) > \epsilon$. Based on this, they conclude that no $(\epsilon, S_{\Delta_\epsilon})$ -tracking algorithm exists. They further observe that BASIC_n is embeddable in many common concept spaces, including halfspaces and axis-aligned rectangles in \mathbb{R}^n , so that for \mathbb{C} equal to either of these, there also is no $(\epsilon, S_{\Delta_\epsilon})$ -tracking algorithm.

4 Adapting to Arbitrarily Varying Drift Rates

This section presents a general bound on the error rate at each time, expressed as a function of the rates of drift, which are allowed to be *arbitrary*. Most importantly, in contrast to the methods from the literature discussed above, the method achieving this general result is *adaptive* to the drift rates, so that it requires no information about the drift rates in advance. This is an appealing property, as it essentially allows the algorithm to learn under an *arbitrary* sequence \mathbf{h}^* of target concepts; the difficulty of the task is then simply reflected in the resulting bounds on the error rates: that is, faster-changing sequences of target functions result in larger bounds on the error rates, but do not require a change in the algorithm itself.

4.1 Adapting to a Changing Drift Rate

Recall that the method yielding (1) (based on the work of [15]) required access to the sequence Δ of changes to achieve the stated guarantee on the expected number of mistakes. That method is based on choosing a classifier to predict \hat{Y}_t by minimizing the number of mistakes among the previous m_t samples, where m_t is a value chosen based on the Δ sequence. Thus, the key to modifying this algorithm to make it adaptive to the Δ sequence is to determine a suitable choice of m_t without reference to the Δ sequence. The strategy we adopt here is to use the *data* to determine an appropriate value \hat{m}_t to use. Roughly (ignoring logarithmic factors for now), the insight that enables us to achieve this feat is that, for the m_t used in the above strategy, one can show that $\sum_{i=t-m_t}^{t-1} \mathbb{1}[h_t^*(X_i) \neq Y_i]$ is roughly $\tilde{O}(d)$, and that making the prediction \hat{Y}_t with *any* $h \in \mathbb{C}$ with roughly $\tilde{O}(d)$ mistakes on these samples will suffice to obtain the stated bound on the error rate (up to logarithmic factors). Thus, if we replace m_t with the largest value m for which $\min_{h \in \mathbb{C}} \sum_{i=t-m}^{t-1} \mathbb{1}[h(X_i) \neq Y_i]$ is roughly $\tilde{O}(d)$, then the above observation implies $m \geq m_t$. This then implies that, for $\hat{h} = \operatorname{argmin}_{h \in \mathbb{C}} \sum_{i=t-m}^{t-1} \mathbb{1}[h(X_i) \neq Y_i]$, we have that $\sum_{i=t-m_t}^{t-1} \mathbb{1}[\hat{h}(X_i) \neq Y_i]$ is also roughly $\tilde{O}(d)$, so that the stated bound on the error rate will be achieved (aside from logarithmic factors) by choosing \hat{h}_t as this classifier \hat{h} . There are a

few technical modifications to this argument needed to get the logarithmic factors to work out properly, and for this reason the actual algorithm below (and proof) is somewhat more involved. Specifically, consider the following algorithm (the value of the universal constant $K \geq 1$ will be specified below).

0. For $T = 1, 2, \dots$
1. Let $\hat{m}_T = \max \left\{ m \in \{1, \dots, T-1\} : \min_{h \in \mathcal{C}} \max_{m' \leq m} \frac{\sum_{t=T-m'}^{T-1} \mathbb{1}[h(X_t) \neq Y_t]}{d \text{Log}(m'/d) + \text{Log}(1/\delta)} < K \right\}$
2. Let $\hat{h}_T = \operatorname{argmin}_{h \in \mathcal{C}} \max_{m' \leq \hat{m}_T} \frac{\sum_{t=T-m'}^{T-1} \mathbb{1}[h(X_t) \neq Y_t]}{d \text{Log}(m'/d) + \text{Log}(1/\delta)}$

Note that the classifiers \hat{h}_t chosen by this algorithm have no dependence on Δ , or anything other than the data $\{(X_i, Y_i) : i < t\}$, and the concept space \mathcal{C} . For space, the proof is deferred to the full version of this paper [13].

Theorem 1. Fix any $\delta \in (0, 1)$, and let \mathcal{A} be the above algorithm. For any sequence Δ in $[0, 1]$, for any \mathcal{P} and any choice of $\mathbf{h}^* \in S_\Delta$, for every $T \in \mathbb{N} \setminus \{1\}$, with probability at least $1 - \delta$,

$$\text{er}_T(\hat{h}_T) \leq O \left(\min_{1 \leq m \leq T-1} \frac{1}{m} \sum_{i=T-m}^{T-1} \sum_{j=i+1}^T \Delta_j + \frac{d \text{Log}(m/d) + \text{Log}(1/\delta)}{m} \right).$$

One immediate implication of Theorem 1 is that, if the sum of Δ_t values grows sublinearly, then there exists an algorithm achieving an expected number of mistakes growing sublinearly in the number of predictions. Formally, we have the following corollary. The proof is deferred to the full version [13].

Corollary 1. If $\sum_{t=1}^T \Delta_t = o(T)$, then there exists an algorithm \mathcal{A} such that, for every \mathcal{P} and every choice of $\mathbf{h}^* \in S_\Delta$, $\mathbb{E} \left[\sum_{t=1}^T \mathbb{1} \left[\hat{Y}_t \neq Y_t \right] \right] = o(T)$.

For many concept spaces of interest, the condition $\sum_{t=1}^T \Delta_t = o(T)$ in Corollary 1 is also a *necessary* condition for *any* algorithm to guarantee a sublinear number of mistakes. In the full version of this paper [13], we establish that for the class of *homogeneous linear separators* on \mathbb{R}^2 with \mathcal{P} the uniform distribution on the unit circle, there exists an algorithm with $\mathbb{E} \left[\sum_{t=1}^T \mathbb{1} \left[\hat{Y}_t \neq Y_t \right] \right] = o(T)$ for every choice of $\mathbf{h}^* \in S_\Delta$ if and only if $\sum_{t=1}^T \Delta_t = o(T)$.

5 Polynomial-Time Algorithms for Linear Separators

In this section, we suppose $\Delta_t = \Delta$ for every $t \in \mathbb{N}$, for a fixed constant $\Delta > 0$, and we consider the special case of learning homogeneous linear separators in \mathbb{R}^k under a uniform distribution on the origin-centered unit sphere. In this case, the analysis of [15] mentioned in Section 3 implies that it is possible to achieve a bound on the error rate that is $\tilde{O}(d\sqrt{\Delta})$, using an algorithm that runs

in time $\text{poly}(d, 1/\Delta, \log(1/\delta))$ (and independent of t) for each prediction. This also implies that it is possible to achieve expected number of mistakes among T predictions that is $\tilde{O}(d\sqrt{\Delta})T$. [8]² have since proven that a variant of the Perceptron algorithm achieves an expected number of mistakes $\tilde{O}((d\Delta)^{1/4})T$.

Below, we improve on this result by showing that there exists an efficient algorithm that achieves a bound on the error rate that is $\tilde{O}(\sqrt{d\Delta})$, as was possible with the inefficient algorithm of [15, 16]. This leads to a bound $\tilde{O}(\sqrt{d\Delta})T$ on the expected number of mistakes. Furthermore, our approach also allows us to present the method as an *active learning* algorithm, and to bound the expected number of queries, as a function of the number of samples T , by $\tilde{O}(\sqrt{d\Delta})T$. The technique is based on modifying the algorithm of [15], replacing an ERM step with (a modification of) the computationally-efficient algorithm of [1].

Formally, define the class of homogeneous linear separators as the set of classifiers $h_w : \mathbb{R}^d \rightarrow \{-1, +1\}$, for $w \in \mathbb{R}^d$ with $\|w\| = 1$, such that $h_w(x) = \text{sign}(w \cdot x)$ for every $x \in \mathbb{R}^d$. We have the following result.

Theorem 2. *When \mathbb{C} is the space of homogeneous linear separators (with $d \geq 4$) and \mathcal{P} is the uniform distribution on the surface of the origin-centered unit sphere in \mathbb{R}^d , for any fixed $\Delta > 0$, for any $\delta \in (0, 1/e)$, there is an algorithm that runs in time $\text{poly}(d, 1/\Delta, \log(1/\delta))$ for each time t , such that for any $\mathbf{h}^* \in S_\Delta$, for every sufficiently large $t \in \mathbb{N}$, with probability at least $1 - \delta$,*

$$\text{er}_t(\hat{h}_t) = O\left(\sqrt{\Delta d \log\left(\frac{1}{\delta}\right)}\right).$$

Also, choosing $\delta = \sqrt{\Delta d} \wedge 1/e$, the expected number of mistakes among the first T predictions is $O\left(\sqrt{\Delta d \log\left(\frac{1}{\Delta d}\right)}T\right)$. Furthermore, the algorithm can be run as an active learning algorithm, in which case, for this δ , the expected number of labels requested by the algorithm among the first T instances is $O\left(\sqrt{\Delta d} \log^{3/2}\left(\frac{1}{\Delta d}\right)T\right)$.

We first state the algorithm used to obtain this result. It is primarily based on a margin-based learning strategy of [1], combined with an initialization step based on a modified Perceptron rule from [8, 9]. For $\tau > 0$ and $x \in \mathbb{R}$, define $\ell_\tau(x) = \max\{0, 1 - \frac{x}{\tau}\}$. Consider the following algorithm and subroutine; parameters δ_k , m_k , τ_k , r_k , b_k , α , and κ will all be specified in the context of the proof (see Lemmas 2 and 6); we suppose $M = \sum_{k=0}^{\lceil \log_2(1/\alpha) \rceil} m_k$.

Algorithm: DriftingHalfspaces

0. Let \tilde{h}_0 be an arbitrary classifier in \mathbb{C}
1. For $i = 1, 2, \dots$
2. $\tilde{h}_i \leftarrow \text{ABL}(M(i-1), \tilde{h}_{i-1})$

² This work in fact studies a much broader model of drift, which allows the distribution \mathcal{P} to vary with time as well. However, this $\tilde{O}((d\Delta)^{1/4})T$ result can be obtained from their theorem by calculating the various parameters for this particular setting.

Subroutine: ModPerceptron(t, \tilde{h})

0. Let w_t be any element of \mathbb{R}^d with $\|w_t\| = 1$
1. For $m = t + 1, t + 2, \dots, t + m_0$
2. Choose $\hat{h}_m = \tilde{h}$ (i.e., predict $\hat{Y}_m = \tilde{h}(X_m)$ as the prediction for Y_m)
3. Request the label Y_m
4. If $h_{w_{m-1}}(X_m) \neq Y_m$
5. $w_m \leftarrow w_{m-1} - 2(w_{m-1} \cdot X_m)X_m$
6. Else $w_m \leftarrow w_{m-1}$
7. Return w_{t+m_0}

Subroutine: ABL(t, \tilde{h})

0. Let w_0 be the return value of ModPerceptron(t, \tilde{h})
1. For $k = 1, 2, \dots, \lceil \log_2(1/\alpha) \rceil$
2. $W_k \leftarrow \{\}$
3. For $s = t + \sum_{j=0}^{k-1} m_j + 1, \dots, t + \sum_{j=0}^k m_j$
4. Choose $\hat{h}_s = \tilde{h}$ (i.e., predict $\hat{Y}_s = \tilde{h}(X_s)$ as the prediction for Y_s)
5. If $|w_{k-1} \cdot X_s| \leq b_{k-1}$, Request label Y_s and let $W_k \leftarrow W_k \cup \{(X_s, Y_s)\}$
6. Find $v_k \in \mathbb{R}^d$ with $\|v_k - w_{k-1}\| \leq r_k$, $0 < \|v_k\| \leq 1$, and

$$\sum_{(x,y) \in W_k} \ell_{\tau_k}(y(v_k \cdot x)) \leq \inf_{v: \|v - w_{k-1}\| \leq r_k} \sum_{(x,y) \in W_k} \ell_{\tau_k}(y(v \cdot x)) + \kappa |W_k|$$
7. Let $w_k = \frac{1}{\|v_k\|} v_k$
8. Return $h_{w_{\lceil \log_2(1/\alpha) \rceil - 1}}$

The general idea here is to replace empirical risk minimization in the method of [15] discussed above with a computationally efficient method of [1]: namely, the subroutine ABL above. For technical reasons, we apply this method to batches of M samples at a time, and simply use the classifier learned from the previous batch to make the predictions. The method of [1] was originally proposed for the problem of *agnostic* learning, to error rate within a constant factor of the optimal. To use this for our purposes, we set up an analogy between the best achievable error rate in agnostic learning and a value $O(\Delta M)$ here (which bounds the best achievable *average* error rate in a given batch).

The analysis of [1] required this method to be initialized with a reasonably accurate classifier (constant bound on its error rate). For this, we find (in Lemma 1) that the modified Perceptron algorithm (of [8,9]) suffices. The ABL algorithm then iteratively refines a hypothesis w_k by taking a number of samples within a slab of width $b_{k-1} \propto 2^{-k}/\sqrt{d}$ around the previous hypothesis separator w_{k-1} , and optimizing a weighted hinge loss (subject to a constraint that the new hypothesis not be too far from the previous). The analysis (Lemma 6) then reveals that the hypothesis w_k approaches a classifier w^* with error rate $O(\Delta M)$ with respect to all of the target concepts in the batch.

We note that, even with the above-described analogy between $O(\Delta M)$ and the noise rate in agnostic learning, the analysis below does not follow immediately from that of [1]. This is because the sample size M that would be required by the analysis of [1] to achieve error rate within a constant factor of the noise

rate would be too large (by a factor of d) for our purposes. In particular, noting that ΔM is increasing in M , converting that original analysis to our present setting would result in a bound on $\text{er}_t(\hat{h}_t)$ larger than that stated in Theorem 2 by roughly a factor of \sqrt{d} . The analysis below refines several aspects of the analysis, using stronger concentration arguments for the weighted hinge loss, and being more careful in bounding the error rate in terms of the weighted hinge loss performance. We thereby reduce the bound to the result stated above.

We have a few lemmas that will be needed for the proof. With some effort, the following result can be derived from the analysis of ModPerceptron by [8]. The details are included in the full version of this article [13].

Lemma 1. *Suppose $\Delta \leq \frac{\pi^2}{400 \cdot 2^{27} (d + \ln(4/\delta))}$. For $m_0 = \max\{\lceil 128(1/c_1) \ln(32) \rceil, \lceil 512 \ln(\frac{4}{\delta}) \rceil\}$, with probability at least $1 - \delta/4$, ModPerceptron(t, \tilde{h}) returns a vector w with $\mathcal{P}(x : h_w(x) \neq h_{t+m_0+1}^*(x)) \leq 1/16$.*

Next, we consider the execution of ABL(t, \tilde{h}), and let the sets W_k be as in that execution. We will denote by w^* the weight vector with $\|w^*\| = 1$ such that $h_{t+m_0+1}^* = h_{w^*}$. Also denote by $M_1 = M - m_0$.

The proof relies on a few results proven in the work of [1], which we summarize in the following lemmas. Although the results were proven in a slightly different setting in that work (agnostic learning under a fixed joint distribution), one can easily verify that their proofs remain valid in our present context as well.

Lemma 2. *[1] Fix any $k \in \{1, \dots, \lceil \log_2(1/\alpha) \rceil\}$. For a universal constant $c_7 > 0$, suppose $b_{k-1} = c_7 2^{1-k} / \sqrt{d}$, and let $z_k = \sqrt{r_k^2 / (d-1) + b_{k-1}^2}$. For a universal constant $c_1 > 0$, if $\|w^* - w_{k-1}\| \leq r_k$,*

$$\left| \mathbb{E} \left[\sum_{(x,y) \in W_k} \ell_{\tau_k}(|w^* \cdot x|) \middle| w_{k-1}, |W_k| \right] - \mathbb{E} \left[\sum_{(x,y) \in W_k} \ell_{\tau_k}(y(w^* \cdot x)) \middle| w_{k-1}, |W_k| \right] \right| \leq c_1 |W_k| \sqrt{2^k \Delta M_1} \frac{z_k}{\tau_k}.$$

Lemma 3. *[4] $\forall c > 0$, there exists $c' > 0$ depending only on c (i.e., not depending on d) such that, for any $u, v \in \mathbb{R}^d$ with $\|u\| = \|v\| = 1$, letting $\sigma = \mathcal{P}(x : h_u(x) \neq h_v(x))$, if $\sigma < 1/2$, then $\mathcal{P}(x : h_u(x) \neq h_v(x) \text{ and } |v \cdot x| \geq c' \frac{\sigma}{\sqrt{d}}) \leq c\sigma$.*

The following is a well-known lemma concerning concentration around the equator for the uniform distribution (see e.g., [1, 3, 9]).

Lemma 4. *For any $C > 0$, there are constants $c_2, c_3 > 0$ depending only on C (i.e., independent of d) such that, for any $w \in \mathbb{R}^d$ with $\|w\| = 1$, $\forall \gamma \in [0, C/\sqrt{d}]$,*

$$c_2 \gamma \sqrt{d} \leq \mathcal{P}(x : |w \cdot x| \leq \gamma) \leq c_3 \gamma \sqrt{d}.$$

Based on this lemma, [1] prove the following.

Lemma 5. *[1] For $X \sim \mathcal{P}$, $\forall w \in \mathbb{R}^d$ with $\|w\| = 1$, $\forall C > 0$ and $\tau, b \in [0, C/\sqrt{d}]$, for c_2, c_3 as in Lemma 4, $\mathbb{E} \left[\ell_{\tau}(|w^* \cdot X|) \middle| |w \cdot X| \leq b \right] \leq \frac{c_3 \tau}{c_2 b}$.*

The following is a stronger version of a result of [1]; specifically, the size of m_k , and the bound on $|W_k|$, are smaller by a factor of d compared to the original.

Lemma 6. *Fix any $\delta \in (0, 1/e)$. For universal constants $c_4, c_5, c_6, c_7, c_8, c_9, c_{10} \in (0, \infty)$, for an appropriate choice of $\kappa \in (0, 1)$ (a universal constant), if $\alpha = c_9 \sqrt{\Delta d \log(\frac{1}{\kappa \delta})}$, for every $k \in \{1, \dots, \lceil \log_2(1/\alpha) \rceil\}$, if $b_{k-1} = c_7 2^{1-k} / \sqrt{d}$, $\tau_k = c_8 2^{-k} / \sqrt{d}$, $r_k = c_{10} 2^{-k}$, $\delta_k = \delta / (\lceil \log_2(4/\alpha) \rceil - k)^2$, and $m_k = \left\lceil c_5 \frac{2^k}{\kappa^2} d \log\left(\frac{1}{\kappa \delta_k}\right) \right\rceil$, and if $\mathcal{P}(x : h_{w_{k-1}}(x) \neq h_{w^*}(x)) \leq 2^{-k-3}$, then with probability at least $1 - (4/3)\delta_k$, $|W_k| \leq c_6 \frac{1}{\kappa^2} d \log\left(\frac{1}{\kappa \delta_k}\right)$ and $\mathcal{P}(x : h_{w_k}(x) \neq h_{w^*}(x)) \leq 2^{-k-4}$.*

Proof. By Lemma 4, and a Chernoff and union bound, for an appropriately large choice of c_5 and any $c_7 > 0$, letting c_2, c_3 be as in Lemma 4 (with $C = c_7 \vee (c_8/2)$), with probability at least $1 - \delta_k/3$,

$$c_2 c_7 2^{-k} m_k \leq |W_k| \leq 4 c_3 c_7 2^{-k} m_k. \tag{2}$$

The claimed upper bound on $|W_k|$ follows from this second inequality.

Next note that, if $\mathcal{P}(x : h_{w_{k-1}}(x) \neq h_{w^*}(x)) \leq 2^{-k-3}$, then

$$\max\{\ell_{\tau_k}(y(w^* \cdot x)) : x \in \mathbb{R}^d, |w_{k-1} \cdot x| \leq b_{k-1}, y \in \{-1, +1\}\} \leq c_{11} \sqrt{d}$$

for some universal constant $c_{11} > 0$. Furthermore, since $\mathcal{P}(x : h_{w_{k-1}}(x) \neq h_{w^*}(x)) \leq 2^{-k-3}$, we know that the angle between w_{k-1} and w^* is at most $2^{-k-3}\pi$, so that $\|w_{k-1} - w^*\| = \sqrt{2 - 2w_{k-1} \cdot w^*} \leq \sqrt{2 - 2 \cos(2^{-k-3}\pi)} \leq \sqrt{2 - 2 \cos^2(2^{-k-3}\pi)} = \sqrt{2} \sin(2^{-k-3}\pi) \leq 2^{-k-3}\pi\sqrt{2}$. For $c_{10} = \pi\sqrt{2}2^{-3}$, this is r_k . By Hoeffding’s inequality (under the conditional distribution given $|W_k|$), the law of total probability, Lemma 2, and linearity of conditional expectations, with probability at least $1 - \delta_k/3$, for $X \sim \mathcal{P}$,

$$\begin{aligned} \sum_{(x,y) \in W_k} \ell_{\tau_k}(y(w^* \cdot x)) &\leq |W_k| \mathbb{E}[\ell_{\tau_k}(|w^* \cdot X|) | w_{k-1}, |w_{k-1} \cdot X| \leq b_{k-1}] \\ &\quad + c_1 |W_k| \sqrt{2^k \Delta M_1 \frac{z_k}{\tau_k}} + \sqrt{|W_k| (1/2) c_{11}^2 d \ln(3/\delta_k)}. \end{aligned} \tag{3}$$

We bound each term on the right separately. By Lemma 5, the first term is at most $|W_k| \frac{c_3 \tau_k}{c_2 b_{k-1}} = |W_k| \frac{c_3 c_8}{2 c_2 c_7}$. Next, $\frac{z_k}{\tau_k} = \frac{\sqrt{c_{10}^2 2^{-2k} / (d-1) + 4c_7^2 2^{-2k} / d}}{c_8 2^{-k} / \sqrt{d}} \leq \frac{\sqrt{2c_{10}^2 + 4c_7^2}}{c_8}$, while $2^k \leq \frac{2}{\alpha}$, so the second term is at most $\sqrt{2} c_1 \frac{\sqrt{2c_{10}^2 + 4c_7^2}}{c_8} |W_k| \sqrt{\frac{\Delta m}{\alpha}}$. Noting

$$M_1 = \sum_{k'=1}^{\lceil \log_2(1/\alpha) \rceil} m_{k'} \leq \frac{32c_5}{\kappa^2} \frac{1}{\alpha} d \log\left(\frac{1}{\kappa \delta}\right), \tag{4}$$

the second term on the right of (3) is at most $\sqrt{\frac{c_5}{c_9} \frac{8c_1}{\kappa} \frac{\sqrt{2c_{10}^2 + 4c_7^2}}{c_8}} |W_k| \sqrt{\frac{\Delta d \log(\frac{1}{\kappa \delta})}{\alpha^2}}$
 $= \frac{8c_1 \sqrt{c_5}}{\kappa} \frac{\sqrt{2c_{10}^2 + 4c_7^2}}{c_8 c_9} |W_k|$. Since $d \ln(3/\delta_k) \leq 2d \ln(1/\delta_k) \leq \frac{2\kappa^2}{c_5} 2^{-k} m_k$,

and (2) implies $2^{-k}m_k \leq \frac{1}{c_2 c_7} |W_k|$, the third term on the right of (3) is at most $|W_k| \frac{c_{11} \kappa}{\sqrt{c_2 c_5 c_7}}$. Altogether, $\sum_{(x,y) \in W_k} \ell_{\tau_k}(y(w^* \cdot x)) \leq |W_k| \left(\frac{c_3 c_8}{2c_2 c_7} + \frac{8c_1 \sqrt{c_5}}{\kappa} \frac{\sqrt{2c_{10}^2 + 4c_7^2}}{c_8 c_9} + \frac{c_{11} \kappa}{\sqrt{c_2 c_5 c_7}} \right)$. For $c_9 = 1/\kappa^3$, $c_8 = \kappa$, this is at most $\kappa |W_k| \left(\frac{c_3}{2c_2 c_7} + 8c_1 \sqrt{c_5} \sqrt{2c_{10}^2 + 4c_7^2} + \frac{c_{11}}{\sqrt{c_2 c_5 c_7}} \right)$.

Next, note that because $h_{w_k}(x) \neq y \Rightarrow \ell_{\tau_k}(y(v_k \cdot x)) \geq 1$, and because (as proven above) $\|w^* - w_{k-1}\| \leq r_k$, $|W_k| \text{er}_{W_k}(h_{w_k}) \leq \sum_{(x,y) \in W_k} \ell_{\tau_k}(y(v_k \cdot x)) \leq \sum_{(x,y) \in W_k} \ell_{\tau_k}(y(w^* \cdot x)) + \kappa |W_k|$. Combined with the above, we have

$$|W_k| \text{er}_{W_k}(h_{w_k}) \leq \kappa |W_k| \left(1 + \frac{c_3}{2c_2 c_7} + 8c_1 \sqrt{c_5} \sqrt{2c_{10}^2 + 4c_7^2} + \frac{c_{11}}{\sqrt{c_2 c_5 c_7}} \right).$$

Let $c_{12} = 1 + \frac{c_3}{2c_2 c_7} + 8c_1 \sqrt{c_5} \sqrt{2c_{10}^2 + 4c_7^2} + \frac{c_{11}}{\sqrt{c_2 c_5 c_7}}$. Furthermore, $|W_k| \text{er}_{W_k}(h_{w_k}) = \sum_{(x,y) \in W_k} \mathbb{1}[h_{w_k}(x) \neq y] \geq \sum_{(x,y) \in W_k} \mathbb{1}[h_{w_k}(x) \neq h_{w^*}(x)] - \sum_{(x,y) \in W_k} \mathbb{1}[h_{w^*}(x) \neq y]$. For an appropriately large value of c_5 , by a Chernoff bound, with probability at least $1 - \delta_k/3$, $\sum_{s=t+\sum_{j=0}^{k-1} m_{j+1}}^{t+\sum_{j=0}^k m_j} \mathbb{1}[h_{w^*}(X_s) \neq Y_s] \leq 2e\Delta M_1 m_k + \log_2(3/\delta_k)$. In particular, this implies $\sum_{(x,y) \in W_k} \mathbb{1}[h_{w^*}(x) \neq y] \leq 2e\Delta M_1 m_k + \log_2(3/\delta_k)$, so that $\sum_{(x,y) \in W_k} \mathbb{1}[h_{w_k}(x) \neq h_{w^*}(x)] \leq |W_k| \text{er}_{W_k}(h_{w_k}) + 2e\Delta M_1 m_k + \log_2(3/\delta_k)$. Noting that (4) and (2) imply

$$\begin{aligned} \Delta M_1 m_k &\leq \Delta \frac{32c_5}{\kappa^2} \frac{d \log\left(\frac{1}{\kappa \delta}\right)}{c_9 \sqrt{\Delta d \log\left(\frac{1}{\kappa \delta}\right)}} \frac{2^k}{c_2 c_7} |W_k| \leq \frac{32c_5}{c_2 c_7 c_9 \kappa^2} \sqrt{\Delta d \log\left(\frac{1}{\kappa \delta}\right)} 2^k |W_k| \\ &= \frac{32c_5}{c_2 c_7 c_9 \kappa^2} \alpha 2^k |W_k| = \frac{32c_5 \kappa^4}{c_2 c_7} \alpha 2^k |W_k| \leq \frac{32c_5 \kappa^4}{c_2 c_7} |W_k|, \end{aligned}$$

and (2) implies $\log_2(3/\delta_k) \leq \frac{2\kappa^2}{c_2 c_5 c_7} |W_k|$, altogether we have

$$\begin{aligned} \sum_{(x,y) \in W_k} \mathbb{1}[h_{w_k}(x) \neq h_{w^*}(x)] &\leq |W_k| \text{er}_{W_k}(h_{w_k}) + \frac{64ec_5 \kappa^4}{c_2 c_7} |W_k| + \frac{2\kappa^2}{c_2 c_5 c_7} |W_k| \\ &\leq \kappa |W_k| \left(c_{12} + \frac{64ec_5 \kappa^3}{c_2 c_7} + \frac{2\kappa}{c_2 c_5 c_7} \right). \end{aligned}$$

Letting $c_{13} = c_{12} + \frac{64ec_5}{c_2 c_7} + \frac{2}{c_2 c_5 c_7}$, and noting $\kappa \leq 1$, we have $\sum_{(x,y) \in W_k} \mathbb{1}[h_{w_k}(x) \neq h_{w^*}(x)] \leq c_{13} \kappa |W_k|$.

Applying a classic ratio-type VC bound (see [17], Section 4.9.2) under the conditional distribution given $|W_k|$, combined with the law of total probability, we have that with probability at least $1 - \delta_k/3$,

$$\begin{aligned} |W_k| \mathcal{P}(x : h_{w_k}(x) \neq h_{w^*}(x) \mid |w_{k-1} \cdot x| \leq b_{k-1}) \\ \leq \sum_{(x,y) \in W_k} \mathbb{1}[h_{w_k}(x) \neq h_{w^*}(x)] + c_{14} \sqrt{|W_k| (d \log(|W_k|/d) + \log(1/\delta_k))}, \end{aligned}$$

for a universal constant $c_{14} > 0$. Combined with the above, and the fact that (2) implies $\log(1/\delta_k) \leq \frac{\kappa^2}{c_2 c_5 c_7} |W_k|$ and $d \log(|W_k|/d) \leq d \log\left(\frac{8c_3 c_5 c_7 \log\left(\frac{1}{\kappa \delta_k}\right)}{\kappa^2}\right) \leq$

$$\begin{aligned}
d \log\left(\frac{8c_3c_5c_7}{\kappa^3\delta_k}\right) &\leq 3 \log(8 \max\{c_3, 1\}c_5)c_5d \log\left(\frac{1}{\kappa\delta_k}\right) \leq 3 \log(8 \max\{c_3, 1\})\kappa^22^{-k}m_k \\
&\leq \frac{3 \log(8 \max\{c_3, 1\})}{c_2c_7}\kappa^2|W_k|, \text{ we have } |W_k|\mathcal{P}(x : h_{w_k}(x) \neq h_{w^*}(x))|w_{k-1} \cdot x| \leq b_{k-1}) \\
&\leq c_{13}\kappa|W_k| + c_{14}\sqrt{|W_k|\left(\frac{3 \log(8 \max\{c_3, 1\})}{c_2c_7}\kappa^2|W_k| + \frac{\kappa^2}{c_2c_5c_7}|W_k|\right)}. \text{ Letting } c_{15} = \\
&\left(c_{13} + c_{14}\sqrt{\frac{3 \log(8 \max\{c_3, 1\})}{c_2c_7} + \frac{1}{c_2c_5c_7}}\right), \text{ this is } c_{15}\kappa|W_k|, \text{ which implies}
\end{aligned}$$

$$\mathcal{P}(x : h_{w_k}(x) \neq h_{w^*}(x))|w_{k-1} \cdot x| \leq b_{k-1}) \leq c_{15}\kappa. \quad (5)$$

Next, note that $\|v_k - w_{k-1}\|^2 = \|v_k\|^2 + 1 - 2\|v_k\| \cos(\pi\mathcal{P}(x : h_{w_k}(x) \neq h_{w_{k-1}}(x)))$. Thus, one implication of the fact that $\|v_k - w_{k-1}\| \leq r_k$ is that $\frac{\|v_k\|}{2} + \frac{1-r_k^2}{2\|v_k\|} \leq \cos(\pi\mathcal{P}(x : h_{w_k}(x) \neq h_{w_{k-1}}(x)))$; since the left hand side is positive, we have $\mathcal{P}(x : h_{w_k}(x) \neq h_{w_{k-1}}(x)) < 1/2$. Additionally, by differentiating, one can easily verify that for $\phi \in [0, \pi]$, $x \mapsto \sqrt{x^2 + 1 - 2x \cos(\phi)}$ is minimized at $x = \cos(\phi)$, in which case $\sqrt{x^2 + 1 - 2x \cos(\phi)} = \sin(\phi)$. Thus, $\|v_k - w_{k-1}\| \geq \sin(\pi\mathcal{P}(x : h_{w_k}(x) \neq h_{w_{k-1}}(x)))$. Since $\|v_k - w_{k-1}\| \leq r_k$, we have $\sin(\pi\mathcal{P}(x : h_{w_k}(x) \neq h_{w_{k-1}}(x))) \leq r_k$. Since $\sin(\pi x) \geq x$ for all $x \in [0, 1/2]$, combining this with the fact (proven above) that $\mathcal{P}(x : h_{w_k}(x) \neq h_{w_{k-1}}(x)) < 1/2$ implies $\mathcal{P}(x : h_{w_k}(x) \neq h_{w_{k-1}}(x)) \leq r_k$.

In particular, we have that both $\mathcal{P}(x : h_{w_k}(x) \neq h_{w_{k-1}}(x)) \leq r_k$ and $\mathcal{P}(x : h_{w^*}(x) \neq h_{w_{k-1}}(x)) \leq 2^{-k-3} \leq r_k$. Now Lemma 3 implies that, for any universal constant $c > 0$, there exists a corresponding universal constant $c' > 0$ such that $\mathcal{P}(x : h_{w_k}(x) \neq h_{w_{k-1}}(x) \text{ and } |w_{k-1} \cdot x| \geq c' \frac{r_k}{\sqrt{d}}) \leq cr_k$ and $\mathcal{P}(x : h_{w^*}(x) \neq h_{w_{k-1}}(x) \text{ and } |w_{k-1} \cdot x| \geq c' \frac{r_k}{\sqrt{d}}) \leq cr_k$, so that $\mathcal{P}(x : h_{w_k}(x) \neq h_{w^*}(x) \text{ and } |w_{k-1} \cdot x| \geq c' \frac{r_k}{\sqrt{d}}) \leq \mathcal{P}(x : h_{w_k}(x) \neq h_{w_{k-1}}(x) \text{ and } |w_{k-1} \cdot x| \geq c' \frac{r_k}{\sqrt{d}}) + \mathcal{P}(x : h_{w^*}(x) \neq h_{w_{k-1}}(x) \text{ and } |w_{k-1} \cdot x| \geq c' \frac{r_k}{\sqrt{d}}) \leq 2cr_k$. In particular, letting $c_7 = c'c_{10}/2$, we have $c' \frac{r_k}{\sqrt{d}} = b_{k-1}$. Combining this with (5), Lemma 4, and a union bound, we have that $\mathcal{P}(x : h_{w_k}(x) \neq h_{w^*}(x)) \leq \mathcal{P}(x : h_{w_k}(x) \neq h_{w^*}(x) \text{ and } |w_{k-1} \cdot x| \geq b_{k-1}) + \mathcal{P}(x : h_{w_k}(x) \neq h_{w^*}(x) \text{ and } |w_{k-1} \cdot x| \leq b_{k-1}) \leq 2cr_k + \mathcal{P}(x : h_{w_k}(x) \neq h_{w^*}(x))|w_{k-1} \cdot x| \leq b_{k-1}) \leq 2cr_k + c_{15}\kappa c_3 b_{k-1} \sqrt{d} = (2^5 c c_{10} + c_{15} \kappa c_3 c_7 2^5) 2^{-k-4}$. Taking $c = \frac{1}{2^6 c_{10}}$ and $\kappa = \frac{1}{2^6 c_3 c_7 c_{15}}$, we have $\mathcal{P}(x : h_{w_k}(x) \neq h_{w^*}(x)) \leq 2^{-k-4}$, as required.

By a union bound, this occurs with probability at least $1 - (4/3)\delta_k$. \square

Proof (Proof of Theorem 2). We begin with the bound on the error rate. If $\Delta > \frac{\pi^2}{400 \cdot 2^{27}(d + \ln(4/\delta))}$, the result trivially holds, since then $1 \leq \frac{400 \cdot 2^{27}}{\pi^2} \sqrt{\Delta(d + \ln(4/\delta))}$.

Otherwise, suppose $\Delta \leq \frac{\pi^2}{400 \cdot 2^{27}(d + \ln(4/\delta))}$. Fix any $i \in \mathbb{N}$. Lemma 1 implies that, with probability at least $1 - \delta/4$, the w_0 returned in Step 0 of $\text{ABL}(M(i-1), \tilde{h}_{i-1})$ satisfies $\mathcal{P}(x : h_{w_0}(x) \neq h_{M(i-1)+m_0+1}^*(x)) \leq 1/16$. Taking this as a base case, Lemma 6 then inductively implies that, with probability at least $1 - \frac{\delta}{4} - \sum_{k=1}^{\lceil \log_2(1/\alpha) \rceil} (4/3) \frac{\delta}{2^{(\lceil \log_2(4/\alpha) \rceil - k)^2}} \geq 1 - \delta$, $\forall k \in \{0, 1, \dots, \lceil \log_2(1/\alpha) \rceil\}$,

$$\mathcal{P}(x : h_{w_k}(x) \neq h_{M(i-1)+m_0+1}^*(x)) \leq 2^{-k-4}, \quad (6)$$

and furthermore the number of labels requested during $\text{ABL}(M(i-1), \tilde{h}_{i-1})$ total to at most (for appropriate universal constants \hat{c}_1, \hat{c}_2) $m_0 + \sum_{k=1}^{\lceil \log_2(1/\alpha) \rceil} |W_k| \leq \hat{c}_1 \left(d + \ln\left(\frac{1}{\delta}\right) + \sum_{k=1}^{\lceil \log_2(1/\alpha) \rceil} d \log\left(\frac{(\lceil \log_2(4/\alpha) \rceil - k)^2}{\delta}\right) \right) \leq \hat{c}_2 d \log\left(\frac{1}{\Delta d}\right) \log\left(\frac{1}{\delta}\right)$. In particular, by a union bound, (6) implies that $\forall k \in \{1, \dots, \lceil \log_2(1/\alpha) \rceil\}$, $\forall m \in \left\{ M(i-1) + \sum_{j=0}^{k-1} m_j + 1, \dots, M(i-1) + \sum_{j=0}^k m_j \right\}$, $\mathcal{P}(x: h_{w_{k-1}}(x) \neq h_m^*(x)) \leq \mathcal{P}(x: h_{w_{k-1}}(x) \neq h_{M(i-1)+m_0+1}^*(x)) + \mathcal{P}(x: h_{M(i-1)+m_0+1}^*(x) \neq h_m^*(x)) \leq 2^{-k-3} + \Delta M$. Since $M = \sum_{k=0}^{\lceil \log_2(1/\alpha) \rceil} m_k = \Theta\left(d + \log\left(\frac{1}{\delta}\right) + \sum_{k=1}^{\lceil \log_2(1/\alpha) \rceil} 2^k d \log\left(\frac{\lceil \log_2(1/\alpha) \rceil - k}{\delta}\right)\right) = \Theta\left(\frac{1}{\alpha} d \log\left(\frac{1}{\delta}\right)\right) = \Theta\left(\sqrt{(d/\Delta) \log(1/\delta)}\right)$, with probability at least $1 - \delta$, $\mathcal{P}(x: h_{w_{\lceil \log_2(1/\alpha) \rceil - 1}}(x) \neq h_{M_i}^*(x)) \leq O(\alpha + \Delta M) = O(\sqrt{\Delta d \log(1/\delta)})$. This implies that, with probability at least $1 - \delta$, $\forall t \in \{Mi + 1, \dots, M(i+1) - 1\}$, $\text{er}_t(\hat{h}_t) \leq \mathcal{P}(x: h_{w_{\lceil \log_2(1/\alpha) \rceil - 1}}(x) \neq h_{M_i}^*(x)) + \mathcal{P}(x: h_{M_i}^*(x) \neq h_t^*(x)) \leq O(\sqrt{\Delta d \log(1/\delta)}) + \Delta M = O\left(\sqrt{\Delta d \log\left(\frac{1}{\delta}\right)}\right)$, which completes the proof of the error rate bound.

Setting $\delta = \sqrt{\Delta d}$, and noting that $\mathbb{1}[\hat{Y}_t \neq Y_t] \leq 1$, we have that for any $t > M$, $\mathbb{P}\left(\hat{Y}_t \neq Y_t\right) = \mathbb{E}\left[\text{er}_t(\hat{h}_t)\right] \leq O\left(\sqrt{\Delta d \log\left(\frac{1}{\delta}\right)}\right) + \delta = O\left(\sqrt{\Delta d \log\left(\frac{1}{\Delta d}\right)}\right)$. The bound on $\mathbb{E}\left[\sum_{t=1}^T \mathbb{1}[\hat{Y}_t \neq Y_t]\right]$ follows by linearity of the expectation. Furthermore, as mentioned, with probability at least $1 - \delta$, an execution of $\text{ABL}(M(i-1), \tilde{h}_{i-1})$ requests at most $O\left(d \log\left(\frac{1}{\Delta d}\right) \log\left(\frac{1}{\delta}\right)\right)$ labels. Thus, since the number of queries during the execution of $\text{ABL}(M(i-1), \tilde{h}_{i-1})$ cannot exceed M , letting $\delta = \sqrt{\Delta d}$, the expected number of queries during an execution is at most $O\left(d \log^2\left(\frac{1}{\Delta d}\right) + \sqrt{\Delta d} M\right) \leq O\left(d \log^2\left(\frac{1}{\Delta d}\right)\right)$. The bound on the expected number of queries among T samples follows by linearity of the expectation. \square

Remark: The original work of [8] additionally allowed some number K of *jumps*: times t with $\Delta_t = 1$. In the above algorithm, since the influence of each sample is localized to the predictors trained within that batch of M instances, the effect of allowing such jumps would only change the bound on the number of mistakes to $\tilde{O}(\sqrt{d\Delta T} + \sqrt{d/\Delta}K)$. This compares favorably to the result of [8], which is roughly $O((d\Delta)^{1/4}T + \frac{d^{1/4}}{\Delta^{3/4}}K)$. However, that result was proven for a more general setting, allowing certain nonuniform distributions \mathcal{P} (though they do require a relation between the angle between separators and the probability mass they disagree on, similar to that holding for the uniform distribution). It is not clear whether Theorem 2 generalizes to this larger family of distributions.

6 General Results for Active Learning

As mentioned, the above results on linear separators also provide results for the number of queries in *active learning*. One can also state quite general results on the expected number of queries and mistakes achievable by an active learning algorithm. This section provides such results, for an algorithm based on the the well-known strategy of *disagreement-based* active learning. Throughout this section, we suppose $\mathbf{h}^* \in S_\Delta$, for a given $\Delta \in (0, 1]$.

First, a few definitions. For any set $\mathcal{H} \subseteq \mathbb{C}$, define the *region of disagreement*:

$$\text{DIS}(\mathcal{H}) = \{x \in \mathcal{X} : \exists h, g \in \mathcal{H} \text{ s.t. } h(x) \neq g(x)\}.$$

This section focuses on the following algorithm. The Active subroutine is from the work of [12] (slightly modified here), and is a variant of the A^2 (Agnostic Active) algorithm of [2]; the values of M and $\hat{T}_k(\cdot)$ are discussed below.

Algorithm: DriftingActive

0. For $i = 1, 2, \dots$
1. Active($M(i - 1)$)

Subroutine: Active(t)

0. Let \hat{h}_0 be an arbitrary element of \mathbb{C} , and let $V_0 \leftarrow \mathbb{C}$
1. Predict $\hat{Y}_{t+1} = \hat{h}_0(X_{t+1})$ as the prediction for the value of Y_{t+1}
2. For $k = 0, 1, \dots, \log_2(M/2)$
3. $Q_k \leftarrow \{\}$
4. For $s = 2^k + 1, \dots, 2^{k+1}$
5. Predict $\hat{Y}_s = \hat{h}_k(X_s)$ as the prediction for the value of Y_s
6. If $X_s \in \text{DIS}(V_k)$
7. Request the label Y_s and let $Q_k \leftarrow Q_k \cup \{(X_s, Y_s)\}$
8. Let $\hat{h}_{k+1} = \operatorname{argmin}_{h \in V_k} \sum_{(x,y) \in Q_k} \mathbb{1}[h(x) \neq y]$
9. Let $V_{k+1} \leftarrow \{h \in V_k : \sum_{(x,y) \in Q_k} \mathbb{1}[h(x) \neq y] - \mathbb{1}[\hat{h}_{k+1}(x) \neq y] \leq \hat{T}_k\}$

As in the DriftingHalbspaces algorithm above, this DriftingActive algorithm proceeds in batches, and in each batch runs an active learning algorithm designed to be robust to classification noise. This robustness to classification noise translates into our setting as tolerance for the fact that there is no classifier in \mathbb{C} that perfectly classifies all of the data. The specific algorithm employed here maintains a set $V_k \subseteq \mathbb{C}$ of candidate classifiers, and requests the labels of samples X_s for which there is some disagreement on the classification among classifiers in V_k . We maintain the invariant that there is a low-error classifier contained in V_k at all times, and thus the points we query provide some information to help us determine which among these remaining candidates has low error rate. Based on these queries, we periodically (in Step 9) remove from V_k those classifiers making a relatively excessive number of mistakes on the queried samples, relative to the minimum among classifiers in V_k . Predictions are made with an element of V_k .

We establish an abstract bound on the number of labels requested by this algorithm, expressed in terms of the *disagreement coefficient* [11]. Specifically, for any $r \geq 0$ and any classifier h , define $B(h, r) = \{g \in \mathbb{C} : \mathcal{P}(x : g(x) \neq h(x)) \leq r\}$. Then for $r_0 \geq 0$ and any classifier h , define the disagreement coefficient of h with respect to \mathbb{C} under \mathcal{P} : $\theta_h(r_0) = \sup_{r > r_0} \frac{\mathcal{P}(\text{DIS}(B(h, r)))}{r}$. Usually, the disagreement coefficient would be used with h equal the target concept; however, since the target concept is not fixed in our setting, we will use the worst-case value: $\theta_{\mathbb{C}}(r_0) = \sup_{h \in \mathbb{C}} \theta_h(r_0)$. This quantity has been bounded for a variety of \mathbb{C} and \mathcal{P} (see e.g., [4, 10, 11]). It is useful in bounding how quickly the region

$\text{DIS}(V_k)$ collapses in the algorithm. Thus, since the probability the algorithm requests the label of the next instance is $\mathcal{P}(\text{DIS}(V_k))$, the value $\theta_{\mathbb{C}}(r_0)$ naturally arises in bounding the number of labels the algorithm requests. Specifically, we have the following result. For space, the proof is deferred to the full version [13].

Theorem 3. *For an appropriate universal constant $c_1 \in [1, \infty)$, if $\mathbf{h}^* \in S_{\Delta}$ for a $\Delta \in (0, 1]$, then with³ $M = \left\lceil c_1 \sqrt{\frac{d}{\Delta}} \right\rceil_2$ and $\hat{T}_k = \log_2(1/\sqrt{d\Delta}) + 2^{2k+2}e\Delta$, defining $\epsilon_{\Delta} = \sqrt{d\Delta}\text{Log}(1/(d\Delta))$, among the first T instances, the expected number of mistakes by `DriftingActive` is $O(\epsilon_{\Delta}\text{Log}(d/\Delta)T) = \tilde{O}\left(\sqrt{d\Delta}\right)T$, and the expected number of label requests is $O(\theta_{\mathbb{C}}(\epsilon_{\Delta})\epsilon_{\Delta}\text{Log}(d/\Delta)T) = \tilde{O}\left(\theta_{\mathbb{C}}(\sqrt{d\Delta})\sqrt{d\Delta}\right)T$.*

References

1. Awasthi, P., Balcan, M.F., Long, P.M.: The power of localization for efficiently learning linear separators with noise. [arXiv:1307.8371v2](#) (2013)
2. Balcan, M.F., Beygelzimer, A., Langford, J.: Agnostic active learning. In: Proceedings of the 23rd International Conference on Machine Learning (2006)
3. Balcan, M.-F., Broder, A., Zhang, T.: Margin based active learning. In: Bshouty, N.H., Gentile, C. (eds.) COLT. LNCS (LNAI), vol. 4539, pp. 35–50. Springer, Heidelberg (2007)
4. Balcan, M.F., Long, P.M.: Active and passive learning of linear separators under log-concave distributions. In: Proceedings of the 26th Conference on Learning Theory (2013)
5. Bartlett, P.L., Ben-David, S., Kulkarni, S.R.: Learning changing concepts by exploiting the structure of change. *Machine Learning* **41**, 153–174 (2000)
6. Bartlett, P.L., Helmbold, D.P.: Learning changing problems (1996) (unpublished)
7. Barve, R.D., Long, P.M.: On the complexity of learning from drifting distributions. *Information and Computation* **138**(2), 170–193 (1997)
8. Crammer, K., Mansour, Y., Even-Dar, E., Vaughan, J.W.: Regret minimization with concept drift. In: Proceedings of the 23rd Conference on Learning Theory, pp. 168–180 (2010)
9. Dasgupta, S., Kalai, A., Monteleoni, C.: Analysis of perceptron-based active learning. *Journal of Machine Learning Research* **10**, 281–299 (2009)
10. El-Yaniv, R., Wiener, Y.: Active learning via perfect selective classification. *Journal of Machine Learning Research* **13**, 255–279 (2012)
11. Hanneke, S.: A bound on the label complexity of agnostic active learning. In: Proceedings of the 24th International Conference on Machine Learning (2007)
12. Hanneke, S.: Activized learning: Transforming passive to active with improved label complexity. *Journal of Machine Learning Research* **13**(5), 1469–1587 (2012)
13. Hanneke, S., Kanade, V., Yang, L.: Learning with a drifting target concept. [arXiv:1505.05215](#) (2015)
14. Haussler, D., Littlestone, N., Warmuth, M.: Predicting $\{0, 1\}$ -functions on randomly drawn points. *Information and Computation* **115**, 248–292 (1994)

³ Here, we define $\lceil x \rceil_2 = 2^{\lceil \log_2(x) \rceil}$, for $x \geq 1$.

15. Helmbold, D.P., Long, P.M.: Tracking drifting concepts by minimizing disagreements. *Machine Learning* **14**(1), 27–45 (1994)
16. Long, P.M.: The complexity of learning according to two models of a drifting environment. *Machine Learning* **37**(3), 337–354 (1999)
17. Vapnik, V.: *Statistical Learning Theory*. John Wiley & Sons Inc., New York (1998)
18. Vapnik, V., Chervonenkis, A.: On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications* **16**, 264–280 (1971)