# Region-Based Active Learning

**Corinna Cortes**
Google Research
New York, NY

**Giulia DeSalvo**
Google Research
New York, NY

**Claudio Gentile**
Google Research
New York, NY

**Mehryar Mohri**
Google Research & Courant
New York, NY

**Ningshan Zhang**
New York University
New York, NY

## Abstract

We study a scenario of active learning where the input space is partitioned into different regions and where a distinct hypothesis is learned for each region. We first introduce a new active learning algorithm (EIWAL), which is an enhanced version of the IWAL algorithm, based on a finer analysis that results in more favorable learning guarantees. Then, we present a new learning algorithm for region-based active learning, ORIWAL, in which either IWAL or EIWAL serve as a subroutine. ORIWAL optimally allocates points to the subroutine algorithm for each region. We give a detailed theoretical analysis of ORIWAL, including generalization error guarantees and bounds on the number of points labeled, in terms of both the hypothesis set used in each region and the probability mass of that region. We also report the results of several experiments for our algorithm which demonstrate substantial benefits over existing non-region-based active learning algorithms, such as IWAL, and over passive learning.

## 1 Introduction

Standard supervised learning algorithms often rely on large amounts of labeled samples to achieve a high performance. But labeling samples is often very costly since it typically requires human inspection and in some cases high human expertise. Can we learn with a limited labeling budget? This is the challenge of active learning, which remains an active area of research in machine learning, with substantial applications and benefits.

Active learning algorithms seek to request as few labels as possible to learn an accurate predictor. There are two standard settings of active learning: the so-called *pool setting* where the algorithm receives as input an i.i.d. pool of

unlabeled points and where it incrementally requests the label of a number of points; and the *on-line setting* where the algorithm receives one i.i.d. point at each round and must decide on whether to request its label. In both cases, after making a number of label requests within a budget, the algorithm returns a predictor chosen out of a hypothesis set, which is hoped to admit a small generalization error. Observe that an active learning algorithm for the on-line setting can also be applied to the pool setting.

In the last few decades, a number of active learning algorithms have been designed, some for specific tasks and requiring strong assumptions. When the problem is separable, Cohn et al. [1994] proposed an algorithm with logarithmic label complexity. A line of work [Dasgupta et al., 2005, Balcan et al., 2007, Balcan and Long, 2013, Awasthi et al., 2014, 2015, Zhang, 2018] studied learning linear separators by labeling samples close to the current estimate of decision boundary. This type of algorithms admits favorable label complexity on the uniform distribution over the unit sphere or on the log-concave distribution. In the pool setting, Dasgupta and Hsu [2008] proposed a hierarchical sampling approach which selectively queries labels from the pool of data and moves down the hierarchies until relatively pure clusters are uncovered. For this type of cluster-based active learning, Urner et al. [2013], Kpotufe et al. [2015] provided a label complexity analysis, but only under various conditions on the data distribution. In the on-line setting, general active learning algorithms [Balcan et al., 2006, Dasgupta et al., 2008, Beygelzimer et al., 2009, 2010, Huang et al., 2015, Zhang and Chaudhuri, 2014] with favorable guarantees both in terms of generalization and label complexity have been devised. These algorithms rely on efficient searching in the concept class, and request labels based on the "disagreement" among hypotheses in the current version space. Their label complexities are bounded in terms of an important quantity known as the disagreement coefficient [Hanneke, 2007]. Among these algorithms, some are computationally inefficient, however, for keeping track of the version space explicitly [Balcan et al., 2006], or for solving expensive optimization problems such as empirical risk minimization with 0-1 loss [Dasgupta et al., 2008, Zhang and Chaudhuri, 2014]. The issue of computational efficiency is one of the key research questions in this area.

This paper considers the on-line active learning in a novel scenario where the input space is partitioned into a finite number of regions. The problem then consists of requesting labels as in the standard case to learn one predictor for each region. This problem naturally arises in a number of applications, such as speech recognition where the regions are data sources (e.g. broadcast news, conversational speech, email, or dictation), and problems in recommendation systems, where the regions are general categories of an item (e.g., film genres). In all these cases, the regions of the input space are suggested by the application at hand. In other tasks, there may be a natural partitioning into regions based on the features used. Nevertheless, simple partitions of the input space, such as random partitions, are often convenient in the absence of prior knowledge about the nature of the input features, and still provide significant benefit in learning, as empirically shown by our experiments.

In all cases, a different hypothesis set can be used for each region and the hope is that often, but not always, the best-in-class predictor at each region will be very accurate, in fact achieving a loss of almost zero on its region. This is the main motivation for our study of *region-based active learning*. As we shall see, in many applications one can indeed achieve a substantially better performance via this formulation of the problem.

The idea of separating the input space in on-line active learning is novel, as all on-line active learning algorithms available in the literature focus on the standard single region input space. A related area in the pool active learning setting is hierarchical sampling (e.g., [Dasgupta and Hsu, 2008]), where the input space admits a hierarchical clustering structure. This scenario of disjoint input space is partially related to stratified sampling techniques in statistics [Neyman, 1934], where a statistical population is divided into disjoint and homogeneous subgroups. Each subgroup is sampled independently, and different criteria can be used to determine an optimal sample size for each group [Rossi et al., 1983]. One such criterion is the sample variance from an existing sample. While such a strategy will help minimize the overall variance, the technique does not address generalization and comes with no learning guarantees.

In this work, we first introduce a new active learning algorithm (EIWAL), which is an enhanced version of the IWAL algorithm from Beygelzimer et al. [2009], based on a finer analysis that results in more favorable learning guarantees. Then, we present a new learning algorithm for region-based active learning, ORIWAL, in which either IWAL or EIWAL serve as a subroutine. ORIWAL optimally allocates points to the subroutine algorithm for each region. We give a detailed theoretical analysis of ORIWAL, including generalization error guarantees and bounds on the number of points labeled, in terms of both the hypothesis set used in each region and the probability mass of that region. We also report the results of several experiments for our algorithm

which demonstrate substantial benefits over existing non-region-based active learning algorithms, such as IWAL, and over passive learning.

The rest of this paper is organized as follows. Section 2 introduces the definitions and notation needed for our analysis and specifies the learning scenario we consider. In Section 3, we introduce the EIWAL algorithm, and prove the associated theoretical guarantees. Section 4 presents our novel region-based active learning algorithm ORIWAL and its learning guarantees. In Section 5, we report the results of our experiments in several datasets. Section 6 concludes the paper with a discussion of future work.

## 2 Preliminaries

In this section, we first introduce the definitions and notation relevant to our analysis and next describe the active learning scenario we consider.

**Definitions.** We denote by $\mathcal{X} \subseteq \mathbb{R}^d$ the input space and by $\mathcal{Y} = \{-1, +1\}$ the binary output space. We assume given a partitioning of $\mathcal{X}$ into $n$ disjoint regions: $\mathcal{X} = \bigcup_{k=1}^{n} \mathcal{X}_k$, with $\mathcal{X}_k \cap \mathcal{X}_{k'} = \emptyset$ for $k \neq k'$. This partitioning may have been generated at random or selected in some other way based on some prior knowledge about the task. In all cases, it is assumed to be fixed before receiving sample points.

As in standard supervised learning, we assume that training and test points are drawn i.i.d. according to some unknown distribution $\mathcal{D}$ over $\mathcal{X} \times \mathcal{Y}$. We will denote by $\mathsf{p}_k = \mathbb{P}(\mathcal{X}_k)$ the probability mass of region $\mathcal{X}_k$ with respect to the marginal distribution induced by $\mathcal{D}$ over $\mathcal{X}$. For each $k \in [n]$, we denote by $\mathcal{H}_k$ the hypothesis set used for region $\mathcal{X}_k$, which consists of functions mapping from $\mathcal{X}$ to some prediction space $\mathcal{Z} \subseteq \mathbb{R}$. In the simplest case, the same hypothesis set is chosen for all regions: $\mathcal{H}_1 = \cdots = \mathcal{H}_n$.

We denote by $\ell \colon \mathcal{Z} \times \mathcal{Y} \to [0, 1]$ the loss function. The loss function we adopt in the implementations run in our experiments is the standard logistic loss, defined for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$ and hypothesis $h$ by $\log(1 + e^{-yh(x)})$, which we then normalize to be in $[0, 1]$. We will denote by $R(h)$ the generalization error or expected loss of a hypothesis $h$: $R(h) = \mathbb{E}[\ell(h(x), y)]$. Similarly, for any $k \in [n]$, we denote by $R_k(h)$ the expected loss of $h$ on region $\mathcal{X}_k$: $R_k(h) = \mathbb{E}_{(x,y)\sim\mathcal{D}}[\ell(h(x), y) \mid x \in \mathcal{X}_k]$. Thus, for any hypothesis $h$, we have $R(h) = \sum_{k=1}^{n} \mathsf{p}_k R_k(h)$.

We will denote by $\mathcal{H}_{[n]}$ be the set of aggregate region-based hypotheses:

$$\mathcal{H}_{[n]} = \left\{ \sum_{k=1}^{n} \mathbb{1}_{x \in \mathcal{X}_k} h_k(x) \colon h_k \in \mathcal{H}_k \right\},$$

whose size $|\mathcal{H}_{[n]}|$ equals $\prod_{k=1}^{n} |\mathcal{H}_k|$. We denote by $h^*$ the best-in-class hypothesis in $\mathcal{H}_{[n]}$, that is, $h^* =$

$\operatorname{argmin}_{h \in \mathcal{H}_{[n]}} R(h)$, and similarly denote by $h_k^*$ the best-in-class hypothesis in region $\mathcal{X}_k$: $h_k^* = \operatorname{argmin}_{h \in \mathcal{H}_k} R_k(h)$. For simplicity, we denote by $R^* = R(h^*)$ and $R_k^* = R_k(h_k^*)$ the error of overall and region-specific best-in-class, respectively. The best-in-class hypothesis $h^* \in \mathcal{H}_{[n]}$ can be expressed as follows in terms of the $h_k^*$s:

$$h^*(x) = \operatorname*{argmin}_{h \in \mathcal{H}_{[n]}} \sum_{k=1}^n \mathbf{p}_k R_k(h) \tag{1}$$

$$= \sum_{k=1}^n \mathbf{1}_{x \in \mathcal{X}_k} \left[ \operatorname*{argmin}_{h \in \mathcal{H}_k} R_k(h) \right] = \sum_{k=1}^n \mathbf{1}_{x \in \mathcal{X}_k} h_k^*(x) \,.$$

Observe, however, that the risk minimization over each region individually is always more advantageous than the risk minimization over the entire space, for the minimal error within the aggregate region-based hypothesis set $\mathcal{H}_{[n]}$ is always less than or equal to the minimal error achieved by selecting each single hypothesis for all regions. Too see this, consider the simplest case where $\mathcal{H}_1 = \cdots = \mathcal{H}_n = \mathcal{H}$. Then, by the super-additivity of the $\min$ operator, the following holds:

$$R(h^*) = \sum_{k=1}^n \mathbf{p}_k \left[ \min_{h \in \mathcal{H}} R_k(h) \right]$$

$$\leq \min_{h \in \mathcal{H}} \left[ \sum_{k=1}^n \mathbf{p}_k R_k(h) \right] = \min_{h \in \mathcal{H}} R(h) \,.$$

In other words, the approximation error of $\mathcal{H}_{[n]}$ is always less than or equal to that of $\mathcal{H}$, implying that $\mathcal{H}_{[n]}$ is always significantly richer than any individual hypothesis set $\mathcal{H}$.

**Learning scenario.** We consider active learning in the *on-line setting*. Unlike the *pool-based setting* where the learner receives the full set of unlabeled points beforehand, in the on-line setting, at each round $t \in [T] = \{1, \ldots, T\}$, the learner receives a point $x_t$ drawn i.i.d. according to the marginal distribution induced by $\mathcal{D}$ on $\mathcal{X}$. She then either selects to request the label of $x_t$, in which case she receives its label $y_t$, or chooses not to solicit $x_t$'s label.

The quality of an active learning algorithm is measured by two quantities in this setting: the generalization error of the hypothesis $h \in \mathcal{H}_{[n]}$ it returns after the $T$ rounds, and the number of labels it requests after $T$ rounds.

## 3 Enhanced-IWAL Algorithm

In this section, we present an enhanced version of the IWAL (Importance Weighted Active Learning) algorithm of Beygelzimer et al. [2009], called EIWAL.

Algorithms such as IWAL use importance weights to address key the issue of sampling bias in active learning. Beygelzimer et al. [2009] gave theoretical guarantees both for the generalization error and the label complexity of IWAL.

Our enhanced version of IWAL admits improved confidence intervals, and thus sharper performance guarantees than the original IWAL, especially in the case where the best-in-class error $R(h^*)$ is small. In that small error regime, EIWAL also improves upon a more recent and more refined importance-weighted active learning algorithm discussed in Beygelzimer et al. [2010] (Theorem 3 therein). This advantage is particularly significant in the scenario of region-based active learning that we are interested in where, often with a large number of regions, the region-based best-in-class errors $R_k(h_k^*)$ are small.

Given a finite hypothesis set $\mathcal{H}$, EIWAL operates on an i.i.d. sample $(x_1, y_1), (x_2, y_2), \ldots, (x_T, y_T)$ drawn according to $\mathcal{D}$. The algorithm maintains at any time $t$ a version space $\mathtt{H}_t$, with $\mathtt{H}_1 = \mathcal{H}$. At time $t$, the algorithm flips a coin $Q_t \in \{0, 1\}$ with bias $p_t = p_t(x_t)$ defined by

$$p_t = \max_{f, g \in \mathtt{H}_t} \max_{y \in \mathcal{Y}} \ell(f(x_t), y) - \ell(g(x_t), y) \,.$$

If $Q_t = 1$, then the label $y_t$ is requested and $\mathtt{H}_t$ is trimmed to $\mathtt{H}_{t+1}$ via an importance-weighted empirical risk minimization:

$$\mathtt{H}_{t+1} = \left\{ h \in \mathtt{H}_t : \frac{1}{t} \sum_{s=1}^t \frac{Q_s}{p_s} \ell(h(x_s), y_s) \leq L_t^* + \Delta_t \right\} \,,$$

where $L_t^*$ is given by

$$L_t^* = \min_{h \in \mathtt{H}_t} \frac{1}{t} \sum_{s=1}^t \frac{Q_s}{p_s} \ell(h(x_s), y_s) \,,$$

and where the slack term $\Delta_t$ is of the form[1]

$$\frac{1}{t} \left[ \sqrt{ \left[ \sum_{s=1}^t p_s \right] \log \left[ \frac{t|\mathcal{H}|}{\delta} \right] } + \log \left[ \frac{t|\mathcal{H}|}{\delta} \right] \right] \,.$$

The definition of the slack term $\Delta_t$ is the main significant difference between EIWAL and the original IWAL. In the latter, $(\sum_{s=1}^t p_s)$ is replaced by the crude upper bound $t$: $\Delta_t = \frac{1}{t} \sqrt{t \log(t|\mathcal{H}|/\delta)}$. The final hypothesis $h_T$ returned by EIWAL is defined as in IWAL:

$$h_T = \operatorname*{argmin}_{h \in \mathtt{H}_T} \frac{1}{T} \sum_{t=1}^T \frac{Q_t}{p_t} \ell(h(x_t), y_t) \,.$$

For our theoretical analysis of EIWAL, we will adopt the definitions and concepts in Beygelzimer et al. [2009]. Define the distance between two hypotheses $f, g \in \mathcal{H}$ as

$$\rho(f, g) = \mathbb{E}_{x \sim \mathcal{D}} \max_y |\ell(f(x), y) - \ell(g(x), y)| \,.$$

---

[1] See the exact expression in the proof of Theorem 1 in Appendix A.

Given $r > 0$, let $B(f, r)$ denote the ball of radius $r$ centered in $f \in \mathcal{H}$: $B(f, r) = \{g \in \mathcal{H} : \rho(f, g) \leq r\}$. The generalized disagreement coefficient $\theta(\mathcal{D}, \mathcal{H})$ can then be defined as follows:

$$\theta(\mathcal{D}, \mathcal{H}) = \inf_{\theta} \left\{ \forall r \geq 0, \right.$$

$$\left. \mathbb{E}_{x \sim \mathcal{D}} \sup_{h \in B(h^*, r)} \sup_y |\ell(h(x), y) - \ell(h^*(x), y)| \leq \theta r \right\},$$

where $h^* = \operatorname{argmin}_{h \in \mathcal{H}} R(h)$. The disagreement coefficient $\theta$ is a complexity measure widely used in disagreement-based active learning problems. In particular, Hanneke [2007] proved upper and lower bounds for the label complexity for the $A^2$ algorithm in terms of the disagreement coefficient $\theta$. Dasgupta et al. [2008] also gave an upper bound for the DHM algorithm using $\theta$. See [Hanneke, 2014] for a more extensive analysis of the disagreement coefficient and active learning.

Using the definitions and concepts just introduced, the following theoretical guarantees can be proven for EIWAL.[2]

**Theorem 1** (EIWAL). *Let $h_T$ denote the hypothesis returned by* EIWAL *after $T$ rounds and $\tau_T$ the total number of requested labels. Then, for all $\delta > 0$, with probability at least $1 - \delta$, for any $T > 0$ the following inequality holds:*

$$R(h_T) \leq R(h^*) + \frac{2}{T} \left[ \sqrt{\sum_{t=1}^{T} p_t} + 6 \sqrt{\log\left[\frac{2(3+T)T^2}{\delta}\right]} \right]$$

$$\times \sqrt{\log\left[\frac{16T^2|\mathcal{H}|^2\log(T)}{\delta}\right]}.$$

*Moreover, with probability at least $1 - \delta$, for any $T > 0$, the following inequality holds:*

$$\tau_T \leq 8\theta K_l \left( R(h^*)T + O(\sqrt{R(h^*)T \log(T|\mathcal{H}|/\delta)}) \right)$$

$$+ O(\log^3(T|\mathcal{H}|/\delta)),$$

*where $K_\ell$ is a constant that depends on the loss function $\ell$.*

For reference, the generalization bound given in [Beygelzimer et al., 2009] for IWAL admits the following form:

$$R(h_T) \leq R(h^*) + \sqrt{\frac{1}{T} \log\left[\frac{T^2|\mathcal{H}|^2}{\delta}\right]}, \qquad (2)$$

and the bound on the number of labels is given by

$$\tau_T = O\left(\theta K_l \left(R(h^*)T + \sqrt{T \log(T|\mathcal{H}|/\delta)}\right)\right). \qquad (3)$$

The comparison of Theorem 1 with (2) and (3), as well as with Beygelzimer et al. [2010] (Theorem 3 therein) shows the following: the bound on the generalization error $R(h_T)$

---

2Due to space limitations, the proofs of all our main results are given in the appendix.

---

**Algorithm 1** ORIWAL$((\mathcal{H}_k)_{k\in[n]}, (\mathrm{p}_k)_{k\in[n]}, \delta, T)$

**for** $k \in [n]$ **do**
  $c_k \leftarrow \log\left[\frac{16T^2|\mathcal{H}_k|^2\log(T)n}{\delta}\right]$
  $\alpha_k \leftarrow \frac{(c_k/\mathrm{p}_k)^{\frac{1}{3}}}{\max_{k\in[n]}(c_k/\mathrm{p}_k)^{\frac{1}{3}}}$
**end for**
**for** $t \in [T]$ **do**
  RECEIVE$(x_t)$
  $k_t \leftarrow k$ such that $x_t \in \mathcal{X}_k$
  $B \sim$ Bernoulli$(\alpha_{k_t})$
  **if** $B = 1$ **then**
    $h_{k,t} \leftarrow$ EIWAL$_{k_t}(x_t)$
    Request $y_t$ according to EIWAL$_{k_t}$ on input $x_t$
    Update (if any) internal state of EIWAL$_{k_t}$
  **end if**
**end for**
**Return** $h_T \leftarrow \left[x \mapsto \sum_{k=1}^{n} 1_{x \in \mathcal{X}_k} h_{k,T}(x)\right]$

---

in Theorem 1 is at least as favorable as the $1/\sqrt{T}$ rate of these previous results, since $\sum_{t=1}^{T} p_t \leq T$. Furthermore, the bound on the number of labels $\tau_T$ is better than both (3) and the results in Beygelzimer et al. [2010] when $R(h^*)$ is small, since we have an extra $R(h^*)$ inside the square root. In fact, in the separable case where $R(h^*) = 0$, our label complexity bound is $\log^3(T)$, which is only poly-logarithmic in $T$, as opposed to the $\sqrt{T}$ guarantee of both (3) and Beygelzimer et al. [2010]. Similarly, when $R(h^*) = 0$, one can see that the generalization error bound of EIWAL has the form $\log^2(T)/T$, rather than $1/\sqrt{T}$ of (2) and Beygelzimer et al. [2010]. This is because $\sum_{t=1}^{T} p_t$ concentrates fast around $\tau_T$ which, as we just said, is only $O(\log^3(T))$ when $R(h^*) = 0$.

## 4  Region-Based Active Learning

In this section, we describe an active learning algorithm, ORIWAL (Optimal Region-based IWAL), under the region-based setting. The algorithm works by running a separate subroutine EIWAL on each of the $n$ regions, while carefully allocating labeling resources across the regions.

### 4.1  The ORIWAL Algorithm

At each time $t$, ORIWAL receives an unlabeled point $x_t$ that belongs to region $\mathcal{X}_{k_t}$, for some $k_t \in [n]$. Then, with some probability $\alpha_k$, ORIWAL decides whether to send $x_t$ to subroutine EIWAL$_{k_t}$, the EIWAL algorithm running on region $\mathcal{X}_{k_t}$. If $x_t$ is sent to EIWAL$_{k_t}$, then it is EIWAL$_{k_t}$ that determines whether to request the associated label $y_t$. Thus, $y_t$ is requested only if $x_t$ is passed to EIWAL$_{k_t}$ (probability $\alpha_k$) *and* EIWAL$_{k_t}$ happens to ask for this label (probability $p_t$ depending on the current state of EIWAL$_{k_t}$). The pseudocode of ORIWAL is given in Algorithm 1.

In what follows, when ORIWAL passes $x_t$ to $\text{EIWAL}_{k_t}$, we say that ORIWAL *queries* $\text{EIWAL}_{k_t}$. Notice that querying $\text{EIWAL}_k$ to determine whether to ask for a label is computationally much more expensive than determining whether or not to pass a point to $\text{EIWAL}_k$. Thus, we will discuss the learning guarantees and label complexity bounds of ORIWAL in terms of the number of queries to the EIWAL subroutines.

The crux of the ORIWAL algorithm rests on finding the probabilities $\alpha_k$, which determine how many points in expectation are passed to the $k$-th region, so as to optimize learning guarantees. Ideally, the algorithm should not pass points to a region where the subroutine has already found a good hypothesis. Regions in need for labels are those where the corresponding subroutines have received few points or where a larger number of points is needed to identify an accurate hypothesis.

To determine the probabilities $\alpha_k$, we first use the theoretical guarantees derived for EIWAL to determine $T_k$, the number of queries made to $\text{EIWAL}_k$ operating in region $\mathcal{X}_k$. At a high level, the optimal setting of $T_k$s, which translates into an optimal setting of $\alpha_k$s, is one that admits the best trade-off between generalization guarantee and label complexity bound. By Theorem 1, the generalization bound of $\text{EIWAL}_k$ is proportional to a complexity term $c_k$ of the form [3] $c_k = \log\left[\frac{16T^2|\mathcal{H}_k|^2\log(T)n}{\delta}\right]$, where we upper bound $\log T_k$ by $\log T$, and further upper bound all label requesting probabilities $p_t$ by 1. [4] Hence, to determine the optimal setting of $T_k$s, we need to find $T_1, T_2, \ldots, T_n$ satisfying:

$$\min_{T_1,\cdots,T_n} \sum_{k=1}^{n} p_k\sqrt{\frac{c_k}{T_k}}, \quad \text{s.t.} \sum_{k=1}^{n} T_k \leq T,$$

where $p_k = \mathbb{P}(\mathcal{X}_k)$. It is straightforward to show that the optimal solution $T_k^*$ admits the following form:

$$T_k^* = \left[\frac{p_k^{\frac{2}{3}}c_k^{\frac{1}{3}}}{\sum_{k'=1}^{n} p_{k'}^{\frac{2}{3}}c_{k'}^{\frac{1}{3}}}\right]T.$$

We then choose the probabilities $\alpha_k$s such that, given the total number $T$ of possible queries, the expected number of queries to $\text{EIWAL}_k$ matches $T_k^*$. That is, $\alpha_k$ should satisfy

$$\frac{p_k\alpha_k}{\sum_{k'=1}^{n} p_{k'}\alpha_{k'}} = \frac{T_k^*}{T} = \frac{p_k^{\frac{2}{3}}c_k^{\frac{1}{3}}}{\sum_{k'=1}^{n} p_{k'}^{\frac{2}{3}}c_{k'}^{\frac{1}{3}}}, \qquad (4)$$

where the left-most side is the conditional probability of querying $\text{EIWAL}_k$, conditioning on a total number $T$ of queries, and the right-most side is the optimal allocation proportion determined by $T_k^*$. It is straightforward to show that for any $\lambda > 0$, $\alpha_k = \lambda(c_k/p_k)^{\frac{1}{3}}$ would satisfy (4).

---

[3] The extra factor $n$ is due to a union bound over the $n$ regions, so as to make Theorem 1 hold for all regions simultaneously.

[4] In Section 4.4, we will present the version of ORIWAL derived from using the original requesting probabilities $p_t$.

Finally, to determine the optimal setting of $\alpha_k$s, we need to determine the last parameter $\lambda$. Observe that, for a given $\lambda$ and its corresponding $\alpha_k$s, a total of $\sum_{k=1}^{n} p_k(1 - \alpha_k)$ unlabeled points will be discarded due to the "**if** $B = 1$ **then** ..." step of ORIWAL (Algorithm 1). Thus, we choose $\lambda$ that minimizes the number of discarded unlabeled points:

$$\min_{\lambda \geq 0} \sum_{k=1}^{n} p_k\big(1-\lambda(c_k/p_k)^{\frac{1}{3}}\big), \text{ s.t. } \lambda\big(c_k/p_k\big)^{\frac{1}{3}} \leq 1, \forall k \in [n].$$

The constraint on $\lambda$ ensures that $\alpha_k$s are valid probabilities: $\alpha_k \leq 1$, $\forall k \in [n]$. Solving the above problem yields the optimal setting of $\alpha_k$s:

$$\lambda = \frac{1}{\max_{k\in[n]}(c_k/p_k)^{\frac{1}{3}}}, \quad \alpha_k = \frac{(c_k/p_k)^{\frac{1}{3}}}{\max_{k\in[n]}(c_k/p_k)^{\frac{1}{3}}}. \tag{5}$$

Observe that in the expression of $\alpha_k$s, we assumed access to the probability mass $p_k$ of each region. This is a reasonable assumption in many applications of active learning, since accurately estimating $p_k$ only requires unlabeled data. Hence, we can conceive a preprocessing stage where the probabilities $p_k$ are accurately estimated from large amounts of unlabeled data. Alternatively, these probabilities can be estimated incrementally, and our analysis can be extended to cover that way of proceeding as well.

## 4.2 Theoretical Analysis

For $\alpha_k$s defined as in (5), the following theoretical guarantees hold for the returned hypothesis and label complexity. The guarantees of ORIWAL depend on region-based disagreement coefficient $\theta_k = \theta(\mathcal{D}_k, \mathcal{H}_k)$, where $\mathcal{D}_k = \mathcal{D}|\mathcal{X}_k$ is defined as the conditional distribution of $x$ on region $k$.

**Theorem 2.** *For any $\delta > 0$, with probability at least $1 - \delta$, for any $T > 0$, the following inequality holds for the hypothesis returned by ORIWAL (Algorithm 1) at time $T$:*

$$R(h_T) \leq R(h^*)$$
$$+ \sum_{k=1}^{n} 2p_k\sqrt{\frac{4\theta_k K_\ell R_k^*}{T_k} \log\left[\frac{16T_k^2|\mathcal{H}_k|^2\log(T_k)n}{\delta}\right]}$$
$$+ \left(\sum_{k=1}^{n} \frac{p_k}{T_k}\right)O\left(\log^2\big(\max_{k\in[n]} T_k|\mathcal{H}_k|n/\delta\big)\right),$$

*where $T_k$ is number of queries made to $\text{IWAL}_k$. Moreover, with probability at least $1 - \delta$, for any $T > 0$, the following inequality holds for the number of requested labels $\tau_T$:*

$$\tau_T \leq \sum_{k=1}^{n}\Big(8\theta_k K_l\Big[R_k^*T_k + O\big(\sqrt{R_k^*T_k\log(T_k|\mathcal{H}_k|n/\delta)}\big)\Big]$$
$$+ O(\log^3(T_k|\mathcal{H}_k|n/\delta))\Big).$$

The generalization bound is the sum of the generalization error of the best in class $h^* \in \mathcal{H}_{[n]}$ and the sum of the

complexity terms of the hypothesis sets $\mathcal{H}_k$. In particular, if the probability mass $\mathsf{p}_k$ of region $\mathfrak{X}_k$ is small, then the corresponding complexity term of set $\mathcal{H}_k$ is given less weight. Moreover, as one could expect, the overall bound becomes tighter as the number of queries $T_k$ made to EI-WAL$_k$ increases. For the label complexity bound of $\tau_T$, the term inside the bracket is of the same form as the term in the label complexity bound of EIWAL for a single region. In this case, the region-specific disagreement coefficients $\theta_k$, best-in-class error $R_k^*$, and complexity terms $\log|\mathcal{H}_k|$, scale the contribution of each region accordingly.

We can also derive guarantees that do not depend on the empirical quantities $T_k$, but only on $T$. When the expected number of passed samples per region is at least $O(\log n)$, we have the following result. For sake of brevity, we denote by $\mathsf{q}_k$ the optimal allocation proportion in Equation (4):

$$\mathsf{q}_k = \frac{\mathsf{p}_k^{2/3} c_k^{1/3}}{\sum_{k'=1}^n \mathsf{p}_{k'}^{2/3} c_{k'}^{1/3}}, \qquad k \in [n] \ .$$

**Corollary 3.** *For all $\delta > 0$, with probability at least $1 - \delta$, for any $T \geq \frac{4\log(2n/\delta)}{\min_{k\in[n]} \mathsf{q}_k}$ the following inequality holds:*

$$R(h_T)$$

$$\leq R(h^*) + 2\sum_{k=1}^n \mathsf{p}_k \sqrt{\frac{4\theta_k K_l R_k^*}{T\mathsf{q}_k} \log\left[\frac{32T^2|\mathcal{H}_k|^2\log(T)n}{\delta}\right]}$$

$$+ \Big(\sum_{k=1}^n \frac{\mathsf{p}_k}{T\mathsf{q}_k}\Big) O\Big(\log^2\big(\max_{k\in[n]} T|\mathcal{H}_k|n/\delta\big)\Big) .$$

*Moreover, with probability at least $1 - 2\delta$, for all $T > 0$, the following inequality holds:*

$$\tau_T \leq 8K_\ell\left[\sum_{k=1}^n \theta_k R_k^* T\mathsf{q}_k\right]$$

$$+ \sum_{k=1}^n O\left(\sqrt{R_k^* T\mathsf{q}_k \log\left[\frac{T|\mathcal{H}_k|n}{\delta}\right]}\right)$$

$$+ O\Big(n\log^3\big(T\max_{k\in[n]}|\mathcal{H}_k|n/\delta\big)\Big) .$$

We have been discussing the learning guarantees in terms of the number of queries to the EIWAL subroutines, and we do not take into account the number of rounds in which the ORIWAL decides not to query EIWAL. This is because, as we have mentioned earlier, querying the EIWAL subroutine consumes a significantly larger amount of computational resources than determining whether to make a query. This view of the learning problem naturally arises in applications where the unlabeled samples are inexpensive and are processed beforehand, so it takes no time to determine their regions and to sample a Bernoulli random variable to decide whether to query. In other words, given a limited amount of resources, we are more interested in the performance of the

algorithm in terms of the number of expensive operations, i.e., queries to the subroutines, than in terms of the number of rounds where no expensive operations are made.

### 4.3 Discussion

The advantage of ORIWAL over non-region-based algorithms is twofold: it seeks region-specific best-in-class hypotheses, and it controls the number of queries on each region in an optimal way. If ORIWAL does not optimize for query allocations but instead sets $\alpha_k = 1$ for all $k \in [n]$, ORIWAL reduces to a special region-based algorithm we call RIWAL (Region-based IWAL). RIWAL still enjoys the advantage of region-based hypotheses, but it simply passes on all the points to the subroutines. The only algorithmic difference between ORIWAL and RIWAL is that the former generates a Bernoulli random variable for each incoming sample point, which only consumes a negligible amount of time compared to querying subroutine EIWAL. Given the same number of queries to EIWAL, the two algorithms therefore have comparable computational cost.

Yet, the learning guarantee of ORIWAL is potentially more favorable than that of RIWAL, since ORIWAL explicitly optimizes for the allocations $T_k$ among a fixed budget of $T$ queries to EIWAL. Given a total of $T$ queries, Corollary 3 provides the generalization error of the hypothesis returned after $T$ rounds, in terms of $\mathsf{q}_k = \mathsf{p}_k\alpha_k/(\sum_{k'}^n \mathsf{p}_{k'}\alpha_{k'})$, the probability of querying EIWAL$_k$, conditioned on a query being made. Upper bounding the constants $4\theta_k K_\ell R_k^*$ by 1 gives the following:

$$R(h_T^{\text{RIWAL}}) \leq R(h^*) + 2\sum_{k=1}^n \mathsf{p}_k \sqrt{\frac{c_k}{T\mathsf{q}_k^{\text{RIWAL}}}} + O\Big(\frac{\mathsf{p}_k}{T\mathsf{q}_k^{\text{RIWAL}}}\Big),$$

$$R(h_T^{\text{ORIWAL}}) \leq R(h^*) + 2\sum_{k=1}^n \mathsf{p}_k \sqrt{\frac{c_k}{T\mathsf{q}_k^{\text{ORIWAL}}}} + O\Big(\frac{\mathsf{p}_k}{T\mathsf{q}_k^{\text{ORIWAL}}}\Big).$$

RIWAL sets $\alpha_k = 1$ and thus $\mathsf{q}_k^{\text{RIWAL}} = \mathsf{p}_k$. Meanwhile, by definition, $\mathsf{q}_k^{\text{ORIWAL}} = \mathsf{p}_k^{2/3} c_k^{1/3}/(\sum_{k'=1}^n \mathsf{p}_{k'}^{2/3} c_{k'}^{1/3})$. Disregarding lower order terms, i.e., the third term in the two upper bounds above, the application of Jensen's inequality to the convex function $x \mapsto x^{\frac{3}{2}}$ yields

$$\sum_{k=1}^n \mathsf{p}_k \sqrt{\frac{c_k}{T\mathsf{q}_k^{\text{ORIWAL}}}} = \sqrt{\frac{1}{T}\left[\sum_{k=1}^n \mathsf{p}_k\left[\frac{c_k}{\mathsf{p}_k}\right]^{\frac{1}{3}}\right]^{\frac{3}{2}}} \qquad (6)$$

$$\leq \sqrt{\frac{1}{T}\left[\sum_{k=1}^n \mathsf{p}_k\left[\frac{c_k}{\mathsf{p}_k}\right]^{\frac{1}{2}}\right]} = \sum_{k=1}^n \mathsf{p}_k\sqrt{\frac{c_k}{T\mathsf{q}_k^{\text{RIWAL}}}}. \qquad (7)$$

Thus ORIWAL yields a potentially more favorable learning guarantee than RIWAL given the same number of $T$ queries to subroutines. Note that the potential improvement of ORI-WAL over RIWAL, that is the difference between (6) and (7), depends on how the ratios $c_k/\mathsf{p}_k$ vary across regions. Unbalanced ratio values across regions make (6) significantly

smaller than (7), while in the case where $c_k/\mathtt{p}_k$ coincide for all $k \in [n]$, there is no improvement.

### 4.4 ORIWAL with time-varying $\alpha_k$

When deriving the optimal value of $\alpha_k$s, we upper bounded $\sum_{t=1}^{T} p_t$ by $T$ in order to simplify the discussion, but there is a finer analysis based on a tighter bound on the complexity term, which results in finding an optimal time-varying $\alpha_k(t)$. Without upper bounding this sum of probabilities, the complexity term of Theorem 1 is $C_k(T_k) = c_k\beta_k(T_k)$, where $\beta_k(T_k) = \left(\sum_t p_t 1_{x_t \in \mathcal{X}_k}\right)/T_k$ is the label requesting probability on region $k$, averaged over the $T_k$ queries. Note that $C_k(T_k) \leq c_k$. Now, since $C_k(T_k)$ depends on $T_k$ which is an unknown quantity at any given round $t \in [T]$, we cannot directly use it to solve the optimization problem:

$$\min_{T_1,\cdots,T_n} \sum_{k=1}^{n} \mathtt{p}_k \sqrt{\frac{C_k(T_k)}{T_k}}, \text{ s.t. } \sum_{k=1}^{n} T_k \leq T.$$

However, by definition, the label requesting probabilities of EIWAL are non-increasing, which implies that $\beta_k(T_k)$ as well as $C_k(T_k)$ are also non-increasing. Thus, at a given current round $t \in [T]$, we can upper bound the above optimization problem by

$$\min_{T_k \geq t_k, k \in [n]} \sum_{k=1}^{n} \mathtt{p}_k \sqrt{\frac{C_k(t_k)}{T_k}}, \text{ s.t. } \sum_{k=1}^{n} T_k \leq T,$$

where $t_k$ denotes the number of queries made for region $k$ at a time $t$. Via a similar reasoning as before, the solution of this optimization problem leads to setting $\alpha_k(t_k) = \frac{(C_k(t_k)/\mathtt{p}_k)^{\frac{1}{3}}}{\max_{k \in [n]}(C_k(t_k)/\mathtt{p}_k)^{\frac{1}{3}}}$. ORIWAL therefore uses these time-varying quantities $\alpha_k(t_k)$ at each time $t$ instead of $\alpha_k$ in Algorithm 1 to determine whether to query EIWAL$_k$.

By using the time-varying and algorithm-dependent quantities $C_k(t_k)$, ORIWAL gains more information about the current state of each region, and uses it to more efficiently allocate labeling resources. More concretely, according to Lemma 6 in Appendix A, $\beta_k(t_k) = 4\theta K_l R_k^* + O(\sqrt{R_k^*/t_k})$. Thus, when $C_k(t_k)/\mathtt{p}_k$ is relatively large for region $k$ (which implies that $\alpha_k(t_k)$ is relatively large), then either $t_k$ is small and $O(\sqrt{R_k^*/t_k})$ is large, so that EIWAL$_k$ is still learning, or $t_k$ is large but the best-in-class error scaled by the probability of that region, $R_k^*/\mathtt{p}_k$, is large. In both cases, ORIWAL allocates more weight to this region, which needs more labeling resources to learn. In the experiments, we ran ORIWAL with the time-varying $\alpha_k(t_k)$s.

Finally, in Appendix C, we present another extension of IWAL to the region-based setting, called NAIVE-IWAL, which simply runs IWAL with the composite hypothesis set $\mathcal{H}_{[n]}$. We show that NAIVE-IWAL is less favorable in terms of theoretical guarantees than RIWAL, thus is less favorable than ORIWAL as well.

Table 1: Binary classification dataset summary: number of observations ($N$), number of features ($d$), proportion of minority class ($r$). Datasets are ordered by number of observations.

| Dataset | $N$ | $d$ | $r$ |
|---|---|---|---|
| magic04 | 19,020 | 10 | 0.352 |
| nomao | 34,465 | 118 | 0.286 |
| shuttle | 43,500 | 9 | 0.216 |
| a9a | 48,842 | 123 | 0.239 |
| ijcnn1 | 49,990 | 22 | 0.097 |
| codrna | 59,535 | 8 | 0.333 |
| skin | 245,057 | 3 | 0.208 |
| covtype | 581,012 | 54 | 0.488 |

## 5 Experiments

In this section, we report the results of experiments comparing the ORIWAL, RIWAL, and IWAL algorithms. We also compared these active learning algorithms with two passive learning algorithms: PASSIVE, which simply requests the label for all points and finds the hypothesis with the smallest empirical logistic loss, and RPASSIVE, which runs PASSIVE on each region separately.

We experimented with the algorithms just mentioned in 8 binary classification datasets from the UCI repository: magic04, nomao, shuttle, a9a, ijcnn1, codrna, skin, covtype. Table 1 gives summary statistics for these 8 datasets. Note that, for each dataset, we kept the first 10 principal components of the original features. For each dataset, we randomly shuffled the data and ran the algorithms on the first 50% of the data, and tested the learned classifier on the remaining 50%. This was repeated 50 times on each dataset, and the results were averaged.

We randomly drew 3,000 hyperplanes with bounded norms as our base hypothesis set, which we call $\mathcal{H}$, and used these 3,000 hyperplanes as $\mathcal{H}_k$ for all regions $\mathcal{X}_k$, thus, we chose $\mathcal{H}_k = \mathcal{H}$ for all $k \in [n]$. To generate disjoint regions, for each dataset we constructed random binary trees, i.e., binary trees with random splitting criteria, and used the resulting terminal nodes as the disjoint regions. Note that these regions are generated without using any labels.

Below, we present these results for four datasets with 10 disjoint regions. The results for the remaining datasets, as well as for the case where we instead have 20 disjoint regions are provided in Appendix D. In Appendix D, we also contrast the performance of ORIWAL with 10 regions vs. ORIWAL with 20 regions.

We first compared the two region-based active learning algorithms, RIWAL and ORIWAL, and the region-based passive learning algorithm RPASSIVE. Both RIWAL and RPASSIVE were run with the same regions and hypothesis sets as ORIWAL, thus all three algorithms have the same model complexity. Figure 1 plots the misclassification loss on held-out
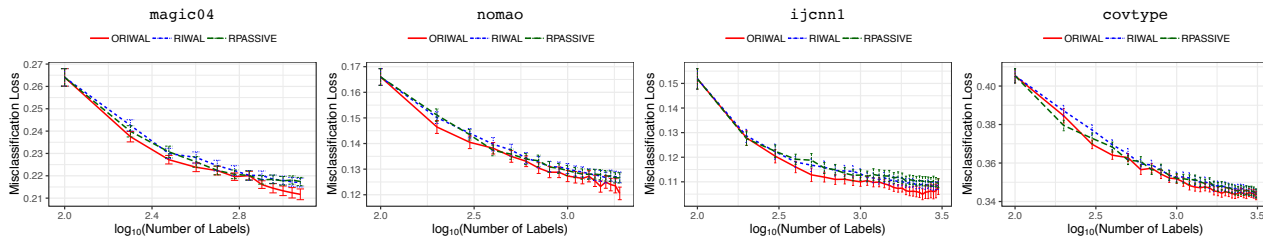
Figure 1: Misclassification loss of ORIWAL, RIWAL, and RPASSIVE on hold out test data versus number of labels requested ($\log_{10}$ scale). The input space was divided into 10 regions. The figures show that ORIWAL typically has a lower misclassification loss than RIWAL and RPASSIVE.
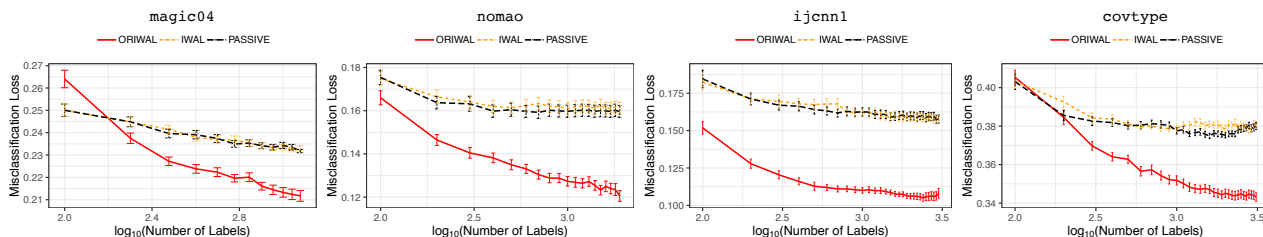


Figure 2: Misclassification loss of ORIWAL (our algorithm), non-region-based IWAL, and non-region-based passive learning PASSIVE on held-out test data, plotted as a function of the number of labels requested ($\log_{10}$ scale). The input space was divided into 10 regions. The curves for ORIWAL are repetitions from Figure 1. The figures show that, given a fixed number of labels, ORIWAL achieves a substantially smaller misclassification loss than IWAL and PASSIVE.

test data against the number of labels requested (on $\log_{10}$ scale), averaged over 50 runs. The error bars indicate $\pm$ standard error. ORIWAL shows consistent advantage over RIWAL and RPASSIVE on most datasets, such as `magic04`, `nomao`, and `ijcnn1`, and matches the performance of RIWAL or RPASSIVE on a few others. Since ORIWAL is significantly outperforming the other two region-based algorithms RIWAL and RPASSIVE, for the rest of our experiments we focused on ORIWAL.

We then compared our proposed algorithm ORIWAL with two baselines: the non-region-based IWAL, and the non-region-based passive learning algorithm, PASSIVE. Both IWAL and PASSIVE were run using the hypothesis set $\mathcal{H}$, which is the hypothesis set used in each region of ORIWAL. Figure 2 plots the misclassification error rate achieved by the three algorithms. The optimal region-based algorithm ORIWAL achieves from the beginning a significantly superior prediction accuracy than the two non region-based algorithms, IWAL and PASSIVE. Given the limited space for improvement when working with the single hypothesis set $\mathcal{H}$, IWAL shows no significant improvement over PASSIVE, and stops improving early on. On the other hand, while the learning curve of non region-based algorithms has plateaued, ORIWAL continues to improve in accuracy by leveraging more labels, and manages to significantly outperform PASSIVE and IWAL.

# 6 Conclusion

We presented a detailed analysis of the scenario of region-based active learning for which we gave a new algorithm, ORIWAL. This algorithm is based on an optimal allocation of points to the underlying region-dependent active learning algorithms. We showed that ORIWAL admits favorable theoretical guarantees, and further demonstrated empirically its substantial improvement over non-region-based algorithms such as IWAL or passive learning in a series of experiments.

Along the way, we also introduced a new active learning algorithm, EIWAL, that benefits from more favorable guarantees than the original IWAL algorithm, and that can be used as a subroutine in our region-based ORIWAL. More generally, other subroutine active learning algorithms can be used with our algorithm, which could lead to further performance improvements in some cases.

We hope to have shown the benefits of region-based active learning and prompted interest in research questions related to this problem. Several crucial questions arise, including the following: How should the regions be chosen? Which hypothesis set should be selected for each? Can we adaptively modify the original partitioning by merging or splitting regions? We have already initiated the study of all of these questions with some preliminary theoretical results. A more complete answer to these and other related questions could lead to significant improvements in active learning.

# References

P. Awasthi, M. F. Balcan, and P. M. Long. The power of localization for efficiently learning linear separators with noise. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 449–458. ACM, 2014.

P. Awasthi, M.-F. Balcan, N. Haghtalab, and R. Urner. Efficient learning of linear separators under bounded noise. In *Conference on Learning Theory*, pages 167–190, 2015.

M.-F. Balcan and P. Long. Active and passive learning of linear separators under log-concave distributions. In *Conference on Learning Theory*, pages 288–316, 2013.

M.-F. Balcan, A. Beygelzimer, and J. Langford. Agnostic active learning. In *ICML*, 2006.

M.-F. Balcan, A. Broder, and T. Zhang. Margin based active learning. In *International Conference on Computational Learning Theory*, pages 35–50. Springer, 2007.

A. Beygelzimer, S. Dasgupta, and J. Langford. Importance weighted active learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 49–56. ACM, 2009.

A. Beygelzimer, D. J. Hsu, J. Langford, and T. Zhang. Agnostic active learning without constraints. In *Advances in Neural Information Processing Systems*, pages 199–207, 2010.

L. Breiman, J. Friedman, R. Olshen, and C. Stone. Classification and regression trees. *CRC Press*, 1984.

N. Cesa-Bianchi and C. Gentile. Improved risk tail bounds for on-line algorithms. *IEEE Transactions on Information Theory*, 54(1):386–390, 2008.

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

S. Dasgupta and D. Hsu. Hierarchical sampling for active learning. In *Proceedings of the 25th international conference on Machine learning*, pages 208–215. ACM, 2008.

S. Dasgupta, A. T. Kalai, and C. Monteleoni. Analysis of perceptron-based active learning. In *International Conference on Computational Learning Theory*, pages 249–263. Springer, 2005.

S. Dasgupta, D. J. Hsu, and C. Monteleoni. A general agnostic active learning algorithm. In *Advances in neural information processing systems*, pages 353–360, 2008.

D. A. Freedman. On tail probabilities for martingales. *the Annals of Probability*, pages 100–118, 1975.

S. Hanneke. A bound on the label complexity of agnostic active learning. In *Proceedings of the 24th international conference on Machine learning*, pages 353–360. ACM, 2007.

S. Hanneke. Theory of disagreement-based active learning. *Foundations and Trends in Machine Learning*, 7(2-3): 131–309, 2014.

T.-K. Huang, A. Agarwal, D. Hsu, J. Langford, and R. E. Schapire. Efficient and parsimonious agnostic active learning. In *Advances in Neural Information Processing Systems 28*, 2015.

S. M. Kakade and A. Tewari. On the generalization ability of online strongly convex programming algorithms. In *Advances in Neural Information Processing Systems*, pages 801–808, 2009.

S. Kpotufe, R. Urner, and S. Ben-David. Hierarchical label queries with data-dependent partitions. In *Conference on Learning Theory*, pages 1176–1189, 2015.

J. Neyman. On the two different aspects of the representative method: the method of stratified sampling and the method of purposive selection. *Journal of the Royal Statistical Society*, 97(4):558–625, 1934.

P. H. Rossi, J. D. Wright, and A. B. Anderson. *Handbook of Survey Research*. Academic Press, 1983.

R. Urner, S. Wulff, and S. Ben-David. Plal: Cluster-based active learning. In *Conference on Learning Theory*, pages 376–397, 2013.

C. Zhang. Efficient active learning of sparse halfspaces. In *COLT*, 2018.

C. Zhang and K. Chaudhuri. Beyond disagreement-based agnostic active learning. In *Advances in Neural Information Processing Systems*, pages 442–450, 2014.