
Improved Algorithms for Agnostic Pool-based Active Classification

Julian Katz-Samuels¹ Jifan Zhang² Lalit Jain² Kevin Jamieson²

Abstract

We consider active learning for binary classification in the agnostic pool-based setting. The vast majority of works in active learning in the agnostic setting are inspired by the CAL algorithm where each query is uniformly sampled from the disagreement region of the current version space. The sample complexity of such algorithms is described by a quantity known as the disagreement coefficient which captures both the geometry of the hypothesis space as well as the underlying probability space. To date, the disagreement coefficient has been justified by minimax lower bounds only, leaving the door open for superior instance dependent sample complexities. In this work we propose an algorithm that, in contrast to uniform sampling over the disagreement region, solves an experimental design problem to determine a distribution over examples from which to request labels. We show that the new approach achieves sample complexity bounds that are never worse than the best disagreement coefficient-based bounds, but in specific cases can be dramatically smaller. From a practical perspective, the proposed algorithm requires no hyperparameters to tune (e.g., to control the aggressiveness of sampling), and is computationally efficient by means of assuming access to an empirical risk minimization oracle (without any constraints). Empirically, we demonstrate that our algorithm is superior to state of the art agnostic active learning algorithms on image classification datasets.

1. Introduction

Most applications of machine learning have an enormous amount of unlabeled data. Yet, many powerful machine

¹University of Wisconsin, Madison, WI ²Paul G. Allen School of Computer Science and Engineering, University of Washington, Seattle, WA. Correspondence to: Kevin Jamieson <jamieson@cs.washington.edu>.

learning methods require that this data be labeled and reliable labels are costly since they require human intervention. The cost of providing labels has become one of the main bottlenecks in applications of machine learning, generating much interest in the problem of *active classification* where the learner is given an unlabeled pool of examples and her goal is to identify an accurate hypothesis using the minimum number of labels possible (Settles, 2011).

One of the most popular algorithmic paradigms is disagreement-based active classification (Hanneke et al., 2014). Under this approach, after observing k labels a version space \mathcal{V}_k of the most promising classifiers is maintained, and the learner queries an example x if there are two hypotheses h_1 and h_2 belonging to \mathcal{V}_k that disagree on the label of x . This approach has received much attention because it applies to generic hypothesis classes, it can be made robust to label noise, and it can be efficient by using a constrained cost-sensitive classification oracle, a problem for which there are many reasonable heuristics (Agarwal et al., 2018; Beygelzimer et al., 2010).

However, disagreement-based active classification suffers from two significant shortcomings. First, it queries uniformly any example on which there is disagreement even though intuitively some of these examples may be much more informative than others. Second, disagreement-based active classification algorithms tend to take a naive union bound over all hypotheses, which ignores many of the dependencies among the hypotheses. Indeed, recent work in pure exploration combinatorial and linear bandits has shown that such naive union bounds can be highly suboptimal and have a significant impact on empirical performance (Cao & Krishnamurthy, 2019; Jain & Jamieson, 2019; Katz-Samuels et al., 2020). Given that these naive union bounds are very loose and appear in the confidence bounds used by the algorithms, in practice, many works instead replace these union bounds with a constant that can be tuned to control the aggressiveness of the algorithm (Beygelzimer et al., 2010; Huang et al., 2015). Unfortunately, this constant introduces a hyperparameter to the active learning algorithm that is difficult to set before seeing lots of data.

We design a new algorithm for pool-based active classification that addresses these shortcomings. It optimizes a novel experimental design objective that finds the best

subset of examples in the disagreement region to query in order to identify the best classifier. It avoids wasteful union bounds by adapting to the geometry of the hypothesis space and thus avoiding the need to choose hyperparameters. We introduce a new notion of sample complexity inspired by experimental design that improves on disagreement-based active classification by a factor up to \sqrt{n} where n is the size of the pool while being only a logarithmic factor worse than disagreement-based learning in the worst case.

1.1. Preliminaries

Let \mathcal{X} denote the input space, and let $\{x_1, \dots, x_n\} \subset \mathcal{X}$ denote a pool of examples. Let \mathcal{H} denote a class of hypotheses where each $h : \mathcal{X} \mapsto \{0, 1\}$ assigns a label to each example in the pool. Let $\mathcal{H}_x := \{(h(x_i))_{i \in [n]} : h \in \mathcal{H}\}$ denote the set of labelings over the pool induced by the hypothesis class \mathcal{H} . Let d denote the VC dimension of \mathcal{H} . When example $i \in [n]$ is queried, the agent receives label $Y_i \sim \text{Bern}(\eta_i)$ where $\eta = (\eta_i)_{i=1}^n \in [0, 1]^n$. We define the error of a hypothesis $h \in \mathcal{H}$ on the pool of examples as given by

$$\begin{aligned} \text{err}(h) &= \frac{1}{n} \sum_{i=1}^n \mathbb{P}(Y_i \neq h(x_i)) \\ &= \frac{1}{n} \sum_{i \in [n]} \eta_i(1 - h(x_i)) + (1 - \eta_i)h(x_i). \end{aligned} \quad (1)$$

Let $h_* := \arg \min_{h \in \mathcal{H}} \text{err}(h)$ be the hypothesis of minimum error, and let $\nu = \text{err}(h_*)$. The goal in active classification is to find an $h \in \mathcal{H}$ with error close to that of h_* using as few label queries as possible. In this paper, we quantify performance as follows:

Problem. Agnostic Pool Based PAC Active Classification: Given $\epsilon > 0, \delta \in (0, 1)$, identify an ϵ -good classifier, that is, an $h \in \mathcal{H}$ such that $\text{err}(h) - \text{err}(h_*) \leq \epsilon$ with probability at least $1 - \delta$ using as few labels as possible.

Remark 1. *The goal of finding an ϵ -good classifier over a pool of examples is closely related to the goal of using an active classification algorithm to find a classifier with good generalization. Suppose $VCdim(\mathcal{H}) = d$ and let \mathcal{D} be a distribution over $\mathcal{X} \times \{0, 1\}$. For $i = 1, \dots, n$ let $(x_i, y_i) \sim \mathcal{D}$. If \hat{h} satisfies $\text{err}(\hat{h}) \leq \min_{h \in \mathcal{H}} \text{err}(h) + \epsilon$, then with probability at least $1 - \delta$*

$$\begin{aligned} \mathbb{P}_{(x,y) \sim \mathcal{D}}(\hat{h}(x) \neq y) &\leq \\ \min_{h \in \mathcal{H}} \mathbb{P}_{(x,y) \sim \mathcal{D}}(h(x) \neq y) + \epsilon + O\left(\sqrt{\frac{d \ln(1/\delta)}{n}}\right). \end{aligned}$$

by standard passive generalization bounds (Boucheron et al., 2005).

1.2. Main contributions

We briefly summarize our contributions:

- We cast pool-based active binary classification as an *adaptive experimental design problem* that computes an optimal sampling distribution over the pool of unlabelled examples. We demonstrate that an ϵ -good classifier can be obtained with probability at least $1 - \delta$ by requesting just $\gamma^*(\epsilon) + \rho^*(\epsilon) \log(1/\delta)$ labels if examples to label are drawn from the optimal design, where $\gamma^*(\epsilon)$ and $\rho^*(\epsilon)$ are problem-dependent quantities defined in the next section.
- Since this optimal design uses problem dependent information like η , it is not a constructive strategy or algorithm for a learner. Treating the sample complexity achieved by this optimal design as a target, we design an algorithm that performs sequential stages of experimental design to match the sample complexity of the optimal design, $\gamma^*(\epsilon) + \rho^*(\epsilon) \log(1/\delta)$ up to a $\log(1/\epsilon)$ factor. The algorithm employs the use of a novel estimator that appeals to a chaining argument. Unfortunately, the method is not computationally efficient.
- We propose a second algorithm that is computationally efficient given access to an empirical risk minimization oracle. The price for computational tractability is a slightly worse sample complexity. Besides being computationally efficient, our approach avoids the need to tune hyperparameters and the use of a *constrained* empirical risk minimization oracle which are required by other active learning algorithms (Beygelzimer et al., 2010; Huang et al., 2015).
- We compare our sample complexity results to those of state-of-the-art disagreement-based learning algorithms that are given in terms of the so-called disagreement coefficient. We demonstrate that our results, up to log factors, are never worse than previous results, but can be substantially better in certain cases.
- Empirically, we compare our procedure to state-of-the-art algorithms for the agnostic setting including variants of the importance weighted active learning algorithm (IWAL) (Beygelzimer et al., 2010) and active cover (Huang et al., 2015). We demonstrate that our method is superior across four image classification tasks.¹

2. Experimental Design for Active Classification

We seek to identify an ϵ -good classifier by seeing as few labels as possible. To this end, we can take motivation from *experimental design* to consider the *optimal* sampling distribution over our pool of unlabeled examples $[n]$. For an arbitrary distribution $\lambda \in \Delta_n := \{p \in \mathbb{R}^n : p_i \geq 0, \forall i \in [n]; \sum_{i=1}^n p_i = 1\}$ suppose we sampled $I_1, \dots, I_t \sim \lambda$ and then observed y_s for each $s \in [t]$. Then an unbiased natural estimator for the error of a classifier $h \in \mathcal{H}$ defined by (1)

¹Code can be found at <https://github.com/jifanz/ACED>.

is given by

$$\widetilde{\text{err}}(h) = \frac{1}{t} \sum_{s=1}^t \frac{1/n}{\lambda_{I_s}} \mathbb{1}\{h(x_{I_s}) \neq y_s\}.$$

Indeed, by i.i.d. sampling from λ , we have for any $s \in [t]$

$$\begin{aligned} \mathbb{E}[\widetilde{\text{err}}(h)] &= \mathbb{E}\left[\frac{1/n}{\lambda_{I_s}} \mathbb{1}\{h(x_{I_s}) \neq y_s\}\right] \\ &= \sum_{i=1}^n \mathbb{P}(I_s = i) \frac{1/n}{\lambda_i} \mathbb{E}[\mathbb{1}\{h(x_i) \neq y_s\} | I_s = i] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{P}(Y_i \neq h(x_i)) = \text{err}(h) \end{aligned}$$

since by definition, $\mathbb{P}(I_s = i) = \lambda_i$. Likewise, an estimator for the excess risk is given by

$$\begin{aligned} \widetilde{\text{err}}(h) - \widetilde{\text{err}}(h_*) &= \quad (2) \\ \frac{1}{t} \sum_{s=1}^t \frac{1/n}{\lambda_i} (\mathbb{1}\{h(x_{I_s}) \neq y_s\} - \mathbb{1}\{h_*(x_{I_s}) \neq y_s\}). \end{aligned}$$

It is straightforward to show that the variance of $\widetilde{\text{err}}(h) - \widetilde{\text{err}}(h_*)$ is upper bounded by $\frac{1}{n} \sum_{i=1}^n \frac{1}{\lambda_i n^2} \mathbb{1}\{h_*(x_i) \neq h(x_i)\}$, using the upper bound $\mathbb{1}\{h(x_I) \neq y_s\} - \mathbb{1}\{h_*(x_I) \neq y_s\} \leq \mathbb{1}\{h_*(x_I) \neq h(x_I)\}$. Applying Bernstein's inequality (and ignoring the $1/t$ term) with probability at least $1 - \delta$

$$\begin{aligned} |\widetilde{\text{err}}(h) - \widetilde{\text{err}}(h_*) - (\text{err}(h) - \text{err}(h_*))| &\lesssim \\ \sqrt{\frac{\sum_{i=1}^n \frac{1}{\lambda_i n^2} \mathbb{1}\{h_*(x_i) \neq h(x_i)\} \log(|\mathcal{H}_x|/\delta)}{t}}. \quad (3) \end{aligned}$$

This then suggests that to estimate the excess error of this particular h with probability at least $1 - \delta$, it suffices to take t large enough to make the RHS of (3) less than ϵ . To upper bound the excess risk of every $h \in \mathcal{H}$ simultaneously, it suffices to take $t \geq \sup_{h \in \mathcal{H}} \frac{\sum_{i=1}^n \frac{1}{\lambda_i n^2} \mathbb{1}\{h_*(x_i) \neq h(x_i)\}}{\max\{\epsilon^2, (\text{err}(h) - \text{err}(h_*))^2\}} \log(|\mathcal{H}|/\delta)$. If we seek to *minimize* the total number of observations, we simply minimize over all $\lambda \in \Delta_n$, motivating the complexity measure:

$$\rho^*(\epsilon) := \inf_{\lambda \in \Delta_n} \sup_{h \in \mathcal{H} \setminus \{h_*\}} \frac{\sum_{i=1}^n \frac{1}{\lambda_i n^2} \mathbb{1}\{h_*(x_i) \neq h(x_i)\}}{\max(\text{err}(h) - \text{err}(h_*), \epsilon)^2}.$$

Thus, we'd expect that if $t \geq \rho^*(\epsilon) \log(|\mathcal{H}|/\delta)$ samples are drawn from the λ that minimizes $\rho^*(\epsilon)$, then $\widehat{h} = \arg \min_{h \in \mathcal{H}} \widetilde{\text{err}}(h)$ will be ϵ -good.

2.1. Sidestepping the Naive Union Bound

A significant shortcoming of the standard approach of applying Bernstein's inequality with a naive union bound is

that the naive union bound incurs an additional factor of $\log(|\mathcal{H}_x|)$ in the sample complexity. For infinite classes, $\log(|\mathcal{H}_x|)$ can be replaced by the VC-dimension of \mathcal{H}_x , however this can still be very loose. In practice, active learning algorithms replace $\log(|\mathcal{H}_x|)$ by a tunable parameter C_0 (Beygelzimer et al., 2010; Huang et al., 2015). Ideally C_0 would be chosen via cross-validation but since our data is being chosen adaptively, under an active algorithm that depends on C_0 , it is unclear how to make the choice a priori.

To improve upon the naive union-bound we appeal to results from empirical process theory. Appealing to the Talagrand/Bousquet inequality (Boucheron et al., 2005), for all $h \in \mathcal{H}$, especially the empirical risk minimizer $\widehat{h} = \arg \min_{h \in \mathcal{H}} \widetilde{\text{err}}(h)$, we have

$$\begin{aligned} \widetilde{\text{err}}(h) - \widetilde{\text{err}}(h^*) - (\text{err}(h) - \text{err}(h^*)) &\leq 2\mathbb{E}[\sup_{h \in \mathcal{H}} |\widetilde{\text{err}}(h) - \widetilde{\text{err}}(h^*) - (\text{err}(h) - \text{err}(h^*))|] \\ &+ \sqrt{\frac{\sup_{h \in \mathcal{H}} \sum_{i=1}^n \frac{1}{\lambda_i n^2} \mathbb{1}\{h_*(x_i) \neq h(x_i)\} \log(1/\delta)}{t}} \\ &+ \frac{4 \sup_{i \in [n]} 1/\lambda_i \log(1/\delta)}{3t}. \end{aligned}$$

Traditionally, we compute the expectation of the suprema using symmetrization to obtain the Rademacher complexity of $\mathcal{H} \setminus \{h_*\}$. In general, the Rademacher complexity is within a $\log(n)$ factor of the Gaussian Width (Bartlett & Mendelson, 2002). In particular,

$$\begin{aligned} \mathbb{E}[\sup_{h \in \mathcal{H}} |\widetilde{\text{err}}(h) - \widetilde{\text{err}}(h^*) - (\text{err}(h) - \text{err}(h^*))|] &\leq \frac{1}{\sqrt{t}} \mathbb{E}_{\zeta \sim N(0, I)} \left[\sup_{h \in \mathcal{H}} \sum_{i \in [n]} \frac{\zeta_i}{n\lambda^{1/2}} (h_*(x_i) - h(x_i)) \right]. \end{aligned}$$

Using the same argument that motivated $\rho^*(\epsilon)$ but applying Bousquet's inequality instead of Bernstein's inequality, we introduce the following new complexity measure for active classification:

$$\gamma^*(\epsilon) := \inf_{\lambda \in \Delta_n} \mathbb{E}_{\zeta} \left[\sup_{h \in \mathcal{H}} \frac{\sum_{i \in [n]} \frac{\zeta_i (h_*(x_i) - h(x_i))}{n\lambda_i^{1/2}}}{\max(\text{err}(h) - \text{err}(h_*), \epsilon)} \right]^2.$$

Analogous to above, if we ignore the $1/t$ term, we'd expect that if $t \geq \gamma^*(\epsilon) + \rho^*(\epsilon) \log(1/\delta)$ samples are drawn from the λ that minimizes the maximum of $\gamma^*(\epsilon)$ and $\rho^*(\epsilon)$, then $\widehat{h} = \arg \min_{h \in \mathcal{H}} \widetilde{\text{err}}(h)$ will be ϵ -good.

We can relate $\gamma^*(\epsilon)$ to $\rho^*(\epsilon)$ in the following way.

Proposition 1 (Katz-Samuels et al. (2020)). $\gamma^*(\epsilon) \leq c \log(|\mathcal{H}_x|) \rho^*(\epsilon) \leq cd \log(\frac{n}{d}) \rho^*(\epsilon)$.

The first inequality parallels the application of Massart's finite class lemma to bound the Rademacher complexity in statistical learning theory and the second inequality follows

from the Sauer-Shelah Lemma. [Katz-Samuels et al. \(2020\)](#) also demonstrates a lower bound on $\gamma^*(\epsilon)$ that is dominated by $\rho^*(\epsilon)$. In the appendix, we show that $\gamma^*(\epsilon)$ matches the minimax rates for classification given for the hypothesis class of thresholds in ([Castro & Nowak, 2008](#)).

Main Takeaway: In Section 3 we will establish an algorithm that achieves a sample complexity of $(\gamma^*(\epsilon) + \rho^*(\epsilon) \log(1/\delta)) \log(1/\epsilon)$ to obtain an ϵ -good classifier with probability greater than $1 - \delta$. In the next section we compare this result to disagreement based methods. Note that we will write $\rho^* := \rho^*(0)$ and $\gamma^* := \gamma^*(0)$.

2.2. Comparison with the Disagreement Coefficient

To date, theoretically grounded active learning algorithms in the agnostic setting are disagreement region sampling methods. At the beginning of each round t these algorithms construct a version space $\mathcal{V} \subset \mathcal{H}$ which is defined to be the set of classifiers that have yet to be ruled out by the algorithm using the observed labels up to round $t - 1$. These algorithms then choose x_{I_t} to be uniformly sampled from $\text{DIS}(\mathcal{V})$, the disagreement region, which is the set of points on which any two hypotheses in \mathcal{V} disagree:

$$\text{DIS}(\mathcal{V}) = \{i : \exists h, h' \in \mathcal{V} \text{ s.t. } h(x_i) \neq h'(x_i)\}.$$

In the notation of the previous section, these algorithms are sampling from λ_t where λ_t is the uniform distribution supported on $\text{DIS}(\mathcal{V})$ ([Hanneke et al., 2014](#)).

The main complexity measure considered for disagreement based algorithms is the disagreement coefficient defined as

$$\theta(\xi) = \sup_{r \geq \xi} \frac{|\text{DIS}(B(h_*, r))|/n}{r}$$

where $B(h_*, r)$ is the ball of radius r centered at h_* :

$$B(h_*, r) = \{h \in \mathcal{H} : \frac{\sum_{i \in [n]} \mathbb{1}\{h_*(x_i) \neq h(x_i)\}}{n} \leq r\}.$$

We consider sample complexity results for finding an h with $\text{err}(h) \leq \nu + \epsilon$, where $\nu = \text{err}(h_*)$ under two common settings.

1. **The Agnostic Setting:** we make no assumptions on $\eta \in [0, 1]^n$. In this case the best known sample complexities scale like

$$\theta(\epsilon) \left(\frac{\nu^2}{\epsilon^2} + \log(1/\epsilon) \right) d$$

where d is the VC dimension of \mathcal{H} ([Hanneke et al., 2014](#)). Note that the noiseless setting of $\eta \in \{0, 1\}^n$ is a special case.

2. **The Tsybakov noise condition:** for some $a \in [1, \infty)$ and $\alpha \in (0, 1]$ every $h \in \mathcal{H} \setminus \{h_*\}$ satisfies

$$\frac{\sum_{i \in [n]} \mathbb{1}\{h_*(x_i) \neq h(x_i)\}}{n} \leq a(\text{err}(h) - \text{err}(h_*))^\alpha.$$

In this case the best known sample complexities scale like:

$$a^2 \frac{1}{\epsilon^{2-2\alpha}} \theta(a\epsilon^\alpha) d \log(1/\epsilon).$$

We now compare our claimed sample complexity of $\gamma^*(\epsilon) + \rho^*(\epsilon) \log(1/\delta)$ to these known sample complexity results. Define $\Delta_{\min} := \min_{h \in \mathcal{H} \setminus \{h_*\}} \text{err}(h) - \text{err}(h_*)$.

Proposition 2. • Suppose that $\eta \in \{0, 1\}^n$.

$$\rho^*(\epsilon) \leq c \log(n \Delta_{\min}^{-1} \vee \epsilon^{-1}) \theta(\epsilon) \left[1 + \frac{\nu^2}{\epsilon^2} \right].$$

- Suppose that the Tsybakov noise condition holds for some $a \in [1, \infty)$ and $\alpha \in (0, 1]$. Then,

$$\rho^*(\epsilon) \leq ca^2 \frac{1}{\epsilon^{2-2\alpha}} \theta(a\epsilon^\alpha) \log(n \Delta_{\min}^{-1} \vee \epsilon^{-1}).$$

Recall that Proposition 1 shows $\gamma^*(\epsilon) \leq cd\rho^*(\epsilon) \log(n/d)$. Hence from Proposition 2, we see that our sample complexity, $\gamma^* + \rho^* \log(1/\delta)$ is always as good as the state-of-the-art sample complexities of disagreement-based learning up to logarithmic factors in n and ϵ^{-1} in both settings.

However, the converse is not true. In general the disagreement based active classification sample complexities can be substantially larger than ρ^* and γ^* .

Proposition 3. There exists an instance where for sufficiently small ξ , $\theta(\xi) \geq \Omega(n^{1/2})$ while $\rho^* = O(1)$ and $\gamma^* = \log(n)$.

We emphasize that this is not just a feature of the analysis; any algorithm that selects queries uniformly at random in the region of disagreement will perform poorly on the instance in the proposition. This gap demonstrates a provable improvement over prior art.

3. Fixed Confidence Algorithm

Algorithm 1 is an elimination-style algorithm, in the style of A^2 ([Balcan et al., 2009](#); [Dasgupta et al., 2007](#); [Huang et al., 2015](#); [Jain & Jamieson, 2019](#)), but optimizes the querying distribution similarly to algorithms from the pure exploration linear bandits literature ([Fiez et al., 2019](#); [Katz-Samuels et al., 2020](#)). It chooses a distribution λ_k over the examples in (4) that minimizes the confidence bounds from Theorem 1 and queries enough random examples from λ_k to ensure that the estimates of the difference in error rates, $\text{err}(h) - \text{err}(h_*)$, improve at least by a factor of 2 for all remaining hypotheses $h \in \mathcal{H}_k$. Using these improved estimates of the gaps, it then eliminates all hypotheses that can be shown to be suboptimal using the confidence bound in Theorem 1.

Given an estimator $\hat{\eta}$ for η , denote the induced estimate for

Algorithm 1 ACED (Active Classification using Experimental Design).

Input: Confidence level $\delta \in (0, 1)$.
 $\mathcal{H}_1 \leftarrow \mathcal{H}$, $k \leftarrow 1$, $\delta_k \leftarrow \delta/2k^2$.
while $|\mathcal{H}_k| > 1$ **do**
 Let λ_k and τ_k be the solution and value of the following optimization problem

$$\inf_{\lambda \in \Delta_n} \mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}_k} \sum_{i \in [n]} h(x_i) \frac{\zeta_i}{n\lambda_i^{1/2}} \right]^2 \quad (4)$$

$$+ 2 \log\left(\frac{1}{\delta_k}\right) \max_{h, h' \in \mathcal{H}_k} \max_{h, h' \in \mathcal{G}} \sum_{i=1}^n \frac{1}{\lambda_i n^2} \mathbb{1}\{h(x_i) \neq h'(x_i)\}$$

Set $N_k \leftarrow c\tau_k 2^{2(k+1)}$ where c is a universal constant.
 Query $I_1, \dots, I_{N_k} \sim \lambda_k$ and receive rewards y_1, \dots, y_{N_k} .
 Let $\hat{\eta}_k := \hat{\eta}(\mathcal{H}_k, \delta_k)$ be the estimator defined in Theorem 1 for \mathcal{H}_k with failure probability δ_k using the samples $\{(x_{I_s}, y_s)\}_{s=1}^{N_k}$.
 $\mathcal{H}_{k+1} \leftarrow \mathcal{H}_k \setminus \{h \in \mathcal{H}_k : \exists h' \text{ such that } \widetilde{\text{err}}(h', \hat{\eta}_k) - \widetilde{\text{err}}(h, \hat{\eta}_k) + \frac{1}{2^{k+1}} \leq 0\}$.
 $k \leftarrow k + 1$
end while
Return: $\mathcal{H}_k = \{\hat{h}\}$.

the error as

$$\widetilde{\text{err}}(h, \hat{\eta}) = \frac{1}{n} \sum_{i \in [n]} \hat{\eta}_i (1 - h(x_i)) + (1 - \hat{\eta}_i) h(x_i).$$

Theorem 1. Let $\mathcal{G} \subset \mathcal{H}$. There exists an estimator $\hat{\eta}(\mathcal{G}, \delta)$ for η constructed from t samples drawn i.i.d. from λ such that with probability at least $1 - \delta$,

$$\sup_{h, h' \in \mathcal{G}} |[\widetilde{\text{err}}(h, \hat{\eta}) - \widetilde{\text{err}}(h', \hat{\eta})] - [\text{err}(h) - \text{err}(h')]|$$

$$\lesssim \sqrt{\frac{\log(2/\delta) \max_{h, h' \in \mathcal{G}} \sum_{i=1}^n \frac{1}{\lambda_i n^2} \mathbb{1}\{h(x_i) \neq h'(x_i)\}}{t}}$$

$$+ \sqrt{\frac{\mathbb{E}[\sup_{h \in \mathcal{G}} \sum_{i \in [n]} h(x_i) \frac{\zeta_i}{n\lambda_i^{1/2}}]^2}{t}}.$$

For now, we treat the estimator in Theorem 1 as a black-box and defer its discussion until Section 4.1. Note that unlike the Talagrand/Bousquet inequality presented before (3), *this confidence interval does not have a term depending on the inverse of the worst case importance weight.*

Algorithm 1 attains the following sample complexity.

Theorem 2. Let $\delta \in (0, 1)$ and $\epsilon > 0$. With probability at least $1 - \delta$ Algorithm 1 returns $\hat{h} \in \mathcal{H}$ after τ samples where $\text{err}(\hat{h}) \leq \text{err}(h_*) + \epsilon$ and

$$\tau \lesssim \log(1/\epsilon) [\log(1/\delta) \rho^*(\epsilon) + \gamma^*(\epsilon)].$$

4. Fixed Budget Algorithm

Algorithm 2 Fixed Budget ACED.

Input: Budget T , tolerance $\epsilon > 0$
 $N \leftarrow \lceil T / \log_2(\epsilon^{-1}) \rceil$, and $\hat{\eta}_0 = 0$
for $k = 1, 2, \dots, \lceil \log_2(\epsilon^{-1}) \rceil$ **do**
 $\tilde{h}_k \leftarrow \arg \min_{h \in \mathcal{H}} \widetilde{\text{err}}(h, \hat{\eta}_{k-1})$.
 Let λ_k be the solution of the following optimization problem

$$\inf_{\lambda \in \Delta_n} \mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}} \frac{\sum_{i \in [n]} (\tilde{h}_k(x_i) - h(x_i)) \frac{\zeta_i}{n\lambda_i^{1/2}}}{2^{-k+1} + \widetilde{\text{err}}(h, \hat{\eta}_{k-1}) - \widetilde{\text{err}}(\tilde{h}_k, \hat{\eta}_{k-1})} \right] \quad (5)$$

Sample $\{x_{I_1}, \dots, x_{I_N}\} \sim \lambda_k$.

Query x_{I_1}, \dots, x_{I_N} and observe y_1, \dots, y_N .

Compute an estimate $\hat{\eta}_k$.

end for

Return: $\arg \min_{h \in \mathcal{H}} \widetilde{\text{err}}(h, \hat{\eta}_k)$

In many applications, the agent is given a budget of T queries and a performance target $\epsilon > 0$, and the goal is to maximize the probability of outputting a classifier $\hat{h} \in \mathcal{H}$ such that $\text{err}(\hat{h}) \leq \text{err}(h_*) + \epsilon$. We design a new algorithm for this setting that can be made computationally efficient given access to a weighted classification oracle (defined shortly).

Algorithm 2 splits the budget into $\lceil \log(\epsilon^{-1}) \rceil$ phases. In each phase, the algorithm computes the design that optimizes (5), the objective of which approximates $\mathbb{E} \left[\max_{h \in \mathcal{H} \setminus \{h_*\}} \frac{\sum_{i \in [n]} \frac{\zeta_i}{n\lambda_i^{1/2}} (h_*(x_i) - h(x_i))}{\max(\text{err}(h) - \text{err}(h_*), 2^{-k+1})} \right]^2$. The algorithm can use any estimator $\hat{\eta}_k$ at each round k . The next theorem uses the estimator of Theorem 1.

Theorem 3. Let $T \in \mathbb{N}$ and $\epsilon > 0$. Let \hat{h} denote the $h \in \mathcal{H}$ returned by Algorithm 2. There exists an estimator $\hat{\eta}_k$ using the samples $\{(x_{I_s}, y_s)\}_{s=1}^N$ in round k of Algorithm 2 such that for an absolute constant $c > 0$

$$\mathbb{P}(\text{err}(\hat{h}) \geq \text{err}(h_*) + \epsilon)$$

$$\leq \log(n\epsilon^{-1})^2 \exp\left(-\frac{cT}{\log(\epsilon^{-1})[\gamma^*(\epsilon) + \rho^*(\epsilon)]}\right).$$

If $T \geq c \log(\log(\epsilon^{-1})) \log(1/\delta) \log(\epsilon^{-1}) [\gamma^*(\epsilon) + \rho^*(\epsilon)]$, then with probability at least $1 - \delta$, Algorithm 2 outputs $\hat{h} \in \mathcal{H}$ such that $\text{err}(\hat{h}) \leq \text{err}(h_*) + \epsilon$. The proof of Theorem 3 leverages the estimator defined in Theorem 1 for \mathcal{H} and failure probability $\delta_k = \exp(-\Theta(N/\gamma_k))$ with γ_k equal to the value of (5).

Remark 2. Given $\{(I_t, y_t)\}_{t=1}^T$ where $I_t \sim \lambda$ define

$$\hat{\eta}_\gamma^{(\text{Importance})} = \frac{1}{T} \sum_{t=1}^T \frac{y_t}{\lambda_{I_t} + \gamma} \mathbf{e}_{I_t}. \quad (6)$$

If importance-weighted estimator $\hat{\eta}_\gamma^{(\text{Importance})}$ with $\gamma = 0$ is used in Algorithm 2 (with a slightly modified objective func-

tion in (5), see the Supplementary Material), one can obtain a computationally efficient algorithm whose probability of error scales as

$$\mathbb{P}(\text{err}(\hat{h}) \geq \text{err}(h_*) + \epsilon) \leq \log(n\epsilon^{-1})^2 \exp\left(-\frac{T - \log(|\mathcal{H}_x|)\psi^*(\epsilon)}{\log(\epsilon^{-1})[\gamma^*(\epsilon) + \rho^*(\epsilon) + \psi^*(\epsilon)]}\right)$$

where

$$\psi^*(\epsilon) := \min_{\lambda \in \Delta_n} \max_{\substack{i \in [n]: \exists h \in \mathcal{H} \\ h_*(x_i) \neq h(x_i)}} \frac{1/n\lambda_i}{\max(\epsilon, \text{err}(h) - \text{err}(h_*))}.$$

There are instances where $\psi^*(\epsilon) \gg \gamma^*(\epsilon)$ and therefore the cost of computational efficiency is a worse sample complexity. See the appendix for more details.

4.1. Discussion of Theorem 1

Theorem 1 above demonstrates the existence of an estimator that avoids any dependence on $\log(|\mathcal{H}_x|)$. The construction of the estimator in Theorem 1 uses generic chaining, a technique that builds a highly optimized union bound to avoid extraneous logarithmic factors (Talagrand, 2014). Generic chaining is most easily applied when a given estimator $\hat{\eta}$ satisfies the property that $\widetilde{\text{err}}(h, \hat{\eta}) - \widetilde{\text{err}}(h', \hat{\eta})$ is sub-Gaussian for every “direction” $h - h'$ of interest (e.g., see (Katz-Samuels et al., 2020)). Though the $\hat{\eta}_{\gamma}^{(\text{Importance})}$ estimator has sub-Gamma tails in general, ruling out its use, the following result shows that for $h - h'$ in a ball under a certain norm, we can construct an estimator for $\widetilde{\text{err}}(h, \hat{\eta}) - \widetilde{\text{err}}(h', \hat{\eta})$ with a sub-Gaussian-like tail.

Proposition 4. Fix $\lambda \in \Delta_n$, $\delta \in (0, 1)$, and $h, h' \in \mathcal{H}$. If T samples are taken from λ and $\hat{\eta} := \hat{\eta}_{\gamma}^{(\text{Importance})}$ is computed with $\gamma = \sqrt{\frac{\log(2/\delta)}{3 \sum_{i=1}^n \frac{1}{\lambda_i n^2} \mathbb{1}\{h(x_i) \neq h'(x_i)\}}}$ then with probability at least $1 - \delta$

$$|\widetilde{\text{err}}(h, \hat{\eta}) - \widetilde{\text{err}}(h', \hat{\eta})| - |\text{err}(h) - \text{err}(h')| \leq \left(\sqrt{\frac{2}{3}} + 1\right) \sqrt{\frac{2 \sum_{i=1}^n \frac{1}{\lambda_i n^2} \mathbb{1}\{h(x_i) \neq h'(x_i)\} \log(\frac{2}{\delta})}{t}}.$$

The idea behind Theorem 1 is to apply generic chaining to all $h - h'$, but to use a different $\hat{\eta}$ (specifically, a different γ) based on the size of $h - h'$ prescribed by Proposition 4. Details of the technique can be found in the supplementary materials.

4.2. Computationally Efficient Experimental Design

In this section, we discuss how to solve (5) efficiently given access to a weighted empirical risk minimization oracle, which we will introduce shortly. First, note that minimizing (5) is equivalent to minimizing

$\mathbb{E}_{\zeta \sim \mathcal{N}(0, I)}[\max_{h \in \mathcal{H}} f(\lambda; h; \zeta)]$ with respect to λ where

$$\begin{aligned} f(\lambda; h; \zeta) &:= \frac{\sum_{i \in [n]} (\tilde{h}_k(x_i) - h(x_i)) \frac{\zeta_i}{n\lambda_i^{1/2}}}{2^{-k+1} + \widetilde{\text{err}}(h, \hat{\eta}_{k-1}) - \widetilde{\text{err}}(\tilde{h}_k, \hat{\eta}_{k-1})} \\ &:= \frac{\sum_{i \in [n]} (\tilde{h}_k(x_i) - h(x_i)) \frac{\zeta_i}{n\lambda_i^{1/2}}}{2^{-k+1} + \sum_{i \in [n]} (1 - 2\hat{\eta}_{k-1, i})(\tilde{h}_k(x_i) - h(x_i))}. \end{aligned}$$

It is known that $\mathbb{E}_{\zeta \sim \mathcal{N}(0, I)}[\max_{h \in \mathcal{H}} f(\lambda; h; \zeta)]$ is convex in λ (Katz-Samuels et al., 2020), hence we perform the minimization over λ via stochastic mirror descent with stochastic gradient $g(\lambda, \zeta) = \nabla f(\lambda, \tilde{h}; \zeta)$ where $\zeta \sim \mathcal{N}(0, I)$ and $\tilde{h} \in \arg \max_{h \in \mathcal{H}} f(\lambda, h; \zeta)$. To obtain \tilde{h} for a fixed λ and ζ , first note that the value $\max_{h \in \mathcal{H}} f(\lambda; h; \zeta)$ is equal to

$$\min_{r \in \mathbb{R}^+} r \text{ subject to } ar + b + \max_{h \in \mathcal{H}} \sum_{i \in [n]} (c_i r + d_i) h(x_i) \leq 0$$

where $a = -2^{-k+1} - \sum_{i \in [n]} (1 - 2\hat{\eta}_{k-1, i}) \tilde{h}_k(x_i)$, $b = \sum_{i \in [n]} \frac{\zeta_i}{n\lambda_i^{1/2}} \tilde{h}_k(x_i)$, $c_i = 1 - 2\hat{\eta}_{k-1, i}$ and $d_i = -\frac{\zeta_i}{n\lambda_i^{1/2}}$.

For any fixed positive value of r it suffices to check the constraint. We can then use a line search procedure to find the minimizing value of r (details in Appendix K).

Thus we have reduced to checking the constraint for a fixed $r \in \mathbb{R}^+$. Specifically, the difficulty is to solve for $\max_{h \in \mathcal{H}} \sum_{i \in [n]} w_i \cdot h(x_i)$ where w_i are arbitrary weights. This can be reduced to weighted 0/1-loss minimization problem that is solvable by a weighted classification oracle:

$$\text{oracle}(\{\tilde{w}_i, \tilde{x}_i, \tilde{y}_i\}_{i=1}^n) := \arg \min_{h \in \mathcal{H}} \sum_{i \in [n]} \tilde{w}_i \cdot \mathbf{1}\{h(\tilde{x}_i) \neq \tilde{y}_i\}$$

for inputs $\{\tilde{w}_i, \tilde{x}_i, \tilde{y}_i\}_{i=1}^n$. Then,

$$\max_{h \in \mathcal{H}} \sum_{i \in [n]} w_i \cdot h(x_i) = \text{oracle}(\{|w_i|, x_i, \mathbf{1}\{w_i \geq 0\}\}_{i=1}^n).$$

5. Implementation and Experiments

In the previous section we reduced the experimental design objective of (5) to a weighted 0/1 loss classification problem using weights that are functions of the estimated vector $\hat{\eta}$. In practice we replace this 0/1 loss with a surrogate convex loss, namely the logistic loss. However, to implement Algorithm 2 we still have to specify the choice of estimator $\hat{\eta}$. Though the estimator specified in Theorem 1 is theoretically grounded, it is difficult to implement in practice since it involves a costly constrained linear optimization problem over the set of hypothesis in \mathcal{H}_k . As described in Remark 2, it is still possible to have a theoretical guarantee for other estimators such as the IPS estimator. As described precisely in Appendix K, in our implementation we take the estimate for $\hat{\eta}_k$ to be

$$[\hat{\eta}_k^{(\text{Naive})}]_i = \text{average}(\{y_s^{(j)} : I_s^{(j)} = i, s \in [N_j], j \in [k]\}),$$

i.e. a simple average of the labels we see. Here $I_s^{(j)}$ indexes the s -th query we made in round j . In our experiments we only considered the persistent noise setting (i.e., querying the same image more than once would always return the same label as before, or formally, $\eta_i \in \{0, 1\}$). Thus, if we sample a point $x_{I_s}^{(j)}$ (i.e., $I_s^{(j)}$) more than once, we set $y_s^{(j)}$ to be the previously observed label and we did not count this observation in our count of total labels taken. To take advantage of all of the labels observed so far, we also employ a water-filling technique for sampling in practice (details in Appendix K).

Baselines. To validate Algorithm 2 we conducted a set of experiments against the following baselines that are considered to be state-of-the-art theoretically-justified methods in disagreement based active learning. Our set of methods are chosen based on the ones considered in (Huang et al., 2015), the most recent work of relevance. Details on the precise implementations of these methods are available in the supplementary materials in Appendix K.

- **Passive:** We considered a passive baseline where we uniformly at random choose samples from our pool, retrain our model on our current samples and report the accuracy.
- **Importance Weighted Active Learning (IWAL)** : IWAL was originally introduced in Beygelzimer et al. (2009) and is an active learning algorithm in the streaming setting. Our implementation is based on the algorithm presented in Beygelzimer et al. (2010) which we refer to as IWAL0. We also consider variants, IWAL1, and oracular versions ORA-IWAL0, ORA-IWAL1 detailed in Huang et al. (2015).
- **Online Active Cover (OAC):** OAC is described in Huang et al. (2015). We used the implementation of OAC that is available in Vowpal Wabbit (Agarwal et al., 2014).

Datasets. We evaluate on the following four real datasets.

- **MNIST 0-4 vs 5-9** (LeCun et al., 1998). We considered the standard MNIST dataset but in a binary setting where digits 0-4 are labelled as 0, and 5-9 are labelled as 1. Our pool has 50000 images in total, and we classified based on the flattened images (784 dimensions).
- **SVHN 2 vs 7** (Netzer et al., 2011). We considered the binary classification problem of determining whether a digit was a 2 or a 7 (ignoring all other images). To prevent the logistic classifier from overfitting to arbitrary labels and to restrict the hypothesis class \mathcal{H} , we downsample the images to 512 dimensional feature vectors through PCA. There are 16180 images in total.
- **CIFAR Bird vs Plane** (Netzer et al., 2011). We considered the binary classification problem of determining whether a digit was a bird or a plane (ignoring all other images). To prevent the logistic classifier from overfitting to arbitrary labels and to restrict the hypothesis class \mathcal{H} ,

we downsample the images to 576 dimensional feature vectors through PCA. There are 10000 images in total.

- **FashionMNIST T-shirt vs Pants** (Xiao et al., 2017). We considered the binary classification problem of T-shirt vs Pants. Our pool has 12000 images in total, and we classified based on the flattened images (784 dimensions).

Implementations. We use two implementations to measure the performances of the algorithms.

- Implementation from Vowpal Wabbit (Agarwal et al., 2014) that is used by Huang et al. (2015). The implementation employs an online learner that only updates based on the latest queried label, therefore has time complexity that scales linearly in the number of images n .
- For our implementation in a batched setting, we retrain the entire classifier to convergence every time new labels become available. We find that the online learner of above can perform significantly better than our batched learner during the first few batches of training. However, our implementation has more stable accuracies during the course of training and performs slightly superior ($< 1\%$) in final accuracy. This comes at a cost of an $O(n^2)$ time complexity, which is too expensive in some of the settings.

In particular, we only use the Vowpal Wabbit implementation for the OAC experiments and the oracular variants of IWAL algorithms for our MNIST experiment, due to the high computation cost for running these algorithms with exhaustive hyperparameter search. However, we think this is still a fair comparison when evaluating some baselines using the the two implementations since it is the best one can achieve for those baselines within a computation budget (single machine with state-of-art commercialized CPU that runs for a month).

Hypothesis Class. In our implementation, we took the hypothesis space to be the set of linear separators in the underlying feature space. We used the logistic regression implementation in Scikit-learn (Pedregosa et al., 2011) for our underlying classification oracle.

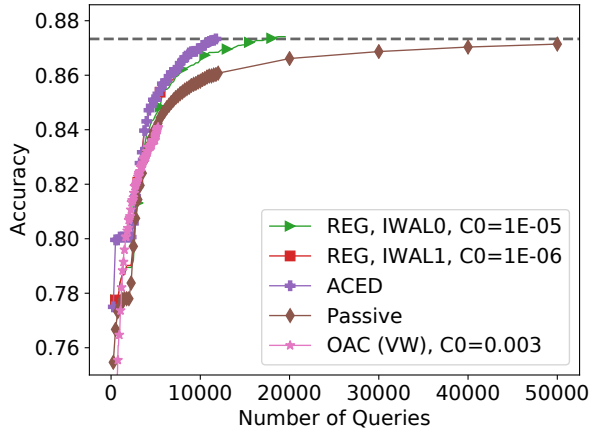


Figure 1. MNIST Performance

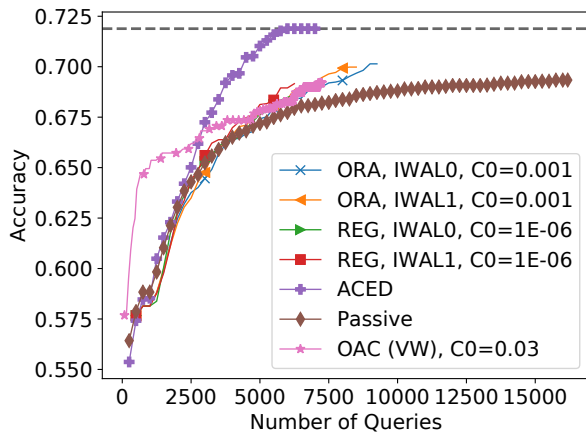


Figure 2. SVHN Performance

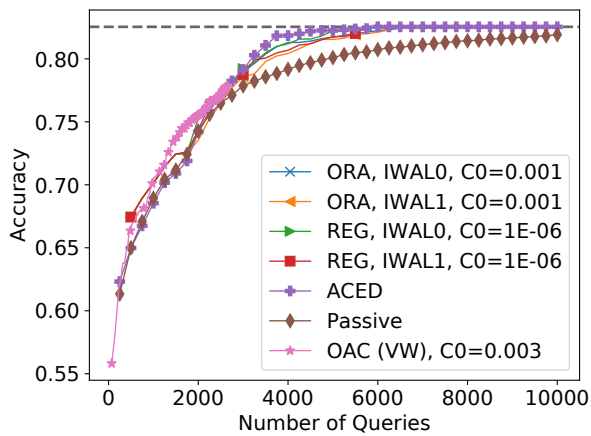


Figure 3. CIFAR Performance

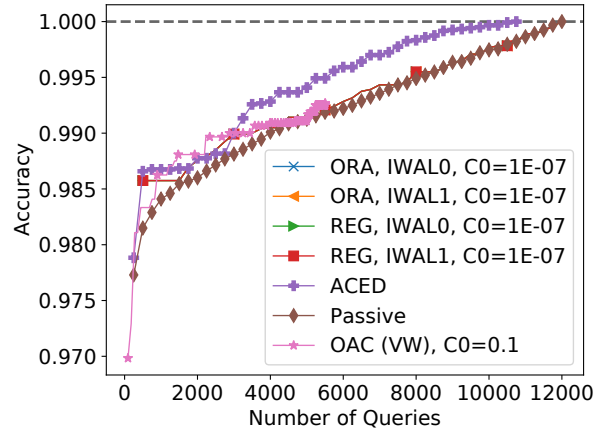


Figure 4. FashionMNIST

Discussion. For each of the binary classification datasets, we plot the running maximum accuracy on the unlabelled pool against the number of queries taken as in Figure 1,2,4,3 (full scale images included in Appendix L). The passive curves are evaluated based on the averages of 10 runs. In the CIFAR experiment, ACED is an average over 5 runs. We find the curve in this setting to be very consistent, and that the standard deviations are minimal for visualization. All of the other curves are evaluated based on a single run. For baselines algorithms proposed in the streaming setting (variants of IWAL and OAC), in each round we uniformly sample an example from the pool, and feed a fixed number of passes. We select the best C_0 based on which hyperparameter setting takes the least amount of queries to reach the same level of accuracy. Detailed hyperparameters considered for the baselines are included in Appendix M. Furthermore, to demonstrate active gains in generalization, we include plots on holdout test sets in Appendix N.

On all four datasets, our algorithm outperforms other baselines by taking much fewer queries to reach the passive accuracy on the entire dataset. Sometimes the active learning algorithms even beat the passive accuracy on the whole dataset, which is a known phenomenon of active learning studied by [Mussmann & Liang \(2018\)](#). For the MNIST dataset, we do not include performance curves for the oracular variants of IWAL, since the Vowpal Wabbit implementation turns out to be performing at random chance. We also notice that OAC stops taking queries very early on (no longer making queries when given more passes over the pool). However, when increasing C_0 , the aggressiveness to make a query, OAC starts performing worse than passive pretty easily. We include Figure 9 in the appendix to demonstrate how sensitive the OAC curves are to the hyperparameter C_0 , which one cannot tune in real applications.

As a special case, on the FashionMNIST dataset, our binary classification task is linearly separable and the baseline

methods fail miserably. For all of the IWAL algorithms on this dataset, we searched in an extended range of hyperparameters than the ones used in the other three tasks. When fixing the random order of the stream, however, all of the baselines become equivalent, and perform almost identical to passive. Since in practice, only one set of hyperparameters can be deployed, this again demonstrates the shortcoming of these baseline algorithms, whereas our method does not rely on any aggressiveness hyperparameter.

6. Related Work and Discussion

Active Classification: Active classification has received much attention with a large number of theoretical and empirical works (see (Hanneke et al., 2014) and (Settles, 2011) for excellent surveys). Cohn et al. (1994) initiated research into the study of disagreement based active classification algorithms, proposing CAL for the realizable setting. Balcan et al. (2009) extended disagreement-based active classification to the agnostic case, introducing the method, A^2 . Hanneke (2007) provided a general analysis of A^2 in terms of the disagreement coefficient, with follow-up works improving on the sample complexity of this approach (Dasgupta et al., 2007; Hanneke, 2009; Hanneke et al., 2011; Koltchinskii, 2010; Hanneke et al., 2014). The results in Section 2.2 show that our sample complexities are never worse than the ones obtained by this line of work.

An extension of this line of work has aimed to attain similar sample complexities, while leveraging an empirical risk minimization oracle to design more practical algorithms (Dasgupta et al., 2007; Hsu, 2010; Beygelzimer et al., 2010; Huang et al., 2015). With the exception of Huang et al. (2015), these methods tend to have a conservative query policy that samples uniformly in the disagreement region, leading to an onerous label requirement. While Huang et al. (2015) has a more aggressive query policy that does not sample uniformly in the disagreement region, their sample complexity result could also be obtained by sampling uniformly in the disagreement region and, therefore, their theoretical result does not reflect gains from a careful selection of points in the disagreement region. In particular, the dominant term is still the disagreement coefficient and, hence, it can be much worse than our sample complexity on instances such as the one in Proposition 3.

Recently, Jain & Jamieson (2019) showed that active classification in the pool-based setting is an instance of combinatorial bandits, an observation that is central to our analysis. They provided the first analysis that shows the contribution of each example to the sample complexity providing a more fine-grained result than the disagreement coefficient. We improve on this work by optimizing the sampling distribution in the region of disagreement and using improved estimators such as the one in Theorem 1. Proposition 4 of

Katz-Samuels et al. (2020) implies that our sample complexity is always better than the sample complexity in Jain & Jamieson (2019).

Finally, we also note that Zhang & Chaudhuri (2014) also give an algorithm that improves on disagreement-based active learning, but the sample complexity of their algorithm is difficult to interpret and their algorithm is not computationally efficient.

Linear and Combinatorial Bandits. ρ^* has been shown to be the dominant term in a lower bound for pure exploration linear bandits and combinatorial bandits (Soare et al., 2014; Chen et al., 2017; Fiez et al., 2019). Recently Katz-Samuels et al. (2020) introduced the notion of γ^* for linear and combinatorial bandits, showing that it is a lower bound for any non-interactive oracle MLE algorithm. One of our contributions is making the connection between the active classification and linear/combinatorial bandit literature, and showing that we can leverage the results from this work to obtain improved sample complexities for agnostic active classification.

Acknowledgements

The authors would like to thank Tzu-Kuo Huang, Alekh Agarwal and John Langford for their help with Vowpal Wabbit baseline experiments. Computational resources from Amazon Web Services were generously gifted as part of an Amazon Research Award. The work of KJ is supported in part by grants NSF RI 1907907 and NSF CCF 2007036.

References

- Agarwal, A., Chapelle, O., Dudík, M., and Langford, J. A reliable effective terascale linear learning system. *The Journal of Machine Learning Research*, 15(1):1111–1133, 2014.
- Agarwal, A., Beygelzimer, A., Dudík, M., Langford, J., and Wallach, H. A reductions approach to fair classification. In *International Conference on Machine Learning*, pp. 60–69. PMLR, 2018.
- Balcan, M.-F., Beygelzimer, A., and Langford, J. Agnostic active learning. *Journal of Computer and System Sciences*, 75(1):78–89, 2009.
- Bartlett, P. L. and Mendelson, S. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- Beygelzimer, A., Dasgupta, S., and Langford, J. Importance weighted active learning. In *Proceedings of the 26th annual international conference on machine learning*, pp. 49–56, 2009.

- Beygelzimer, A., Hsu, D. J., Langford, J., and Zhang, T. Agnostic active learning without constraints. In *Advances in neural information processing systems*, pp. 199–207, 2010.
- Boucheron, S., Bousquet, O., and Lugosi, G. Theory of classification: A survey of some recent advances. *ESAIM: probability and statistics*, 9:323–375, 2005.
- Cao, T. and Krishnamurthy, A. Disagreement-based combinatorial pure exploration: Sample complexity bounds and an efficient algorithm. In *Conference on Learning Theory*, pp. 558–588, 2019.
- Castro, R. M. and Nowak, R. D. Minimax bounds for active learning. *IEEE Transactions on Information Theory*, 54(5):2339–2353, 2008.
- Chen, L., Gupta, A., Li, J., Qiao, M., and Wang, R. Nearly optimal sampling algorithms for combinatorial pure exploration. In *Conference on Learning Theory*, pp. 482–534, 2017.
- Cohn, D., Atlas, L., and Ladner, R. Improving generalization with active learning. *Machine learning*, 15(2): 201–221, 1994.
- Dasgupta, S., Hsu, D. J., and Monteleoni, C. A general agnostic active learning algorithm. *Advances in neural information processing systems*, 20:353–360, 2007.
- Fiez, T., Jain, L., Jamieson, K. G., and Ratliff, L. Sequential experimental design for transductive linear bandits. In *Advances in Neural Information Processing Systems*, pp. 10666–10676, 2019.
- Hanneke, S. A bound on the label complexity of agnostic active learning. In *Proceedings of the 24th international conference on Machine learning*, pp. 353–360, 2007.
- Hanneke, S. Adaptive rates of convergence in active learning. In *COLT*. Citeseer, 2009.
- Hanneke, S. et al. Rates of convergence in active learning. *The Annals of Statistics*, 39(1):333–361, 2011.
- Hanneke, S. et al. Theory of disagreement-based active learning. *Foundations and Trends® in Machine Learning*, 7(2-3):131–309, 2014.
- Hsu, D. J. *Algorithms for active learning*. PhD thesis, UC San Diego, 2010.
- Huang, T.-K., Agarwal, A., Hsu, D. J., Langford, J., and Schapire, R. E. Efficient and parsimonious agnostic active learning. In *Advances in Neural Information Processing Systems*, pp. 2755–2763, 2015.
- Jain, L. and Jamieson, K. G. A new perspective on pool-based active classification and false-discovery control. In *Advances in Neural Information Processing Systems*, pp. 13992–14003, 2019.
- Karampatziakis, N. and Langford, J. Online importance weight aware updates. *arXiv preprint arXiv:1011.1576*, 2010.
- Katz-Samuels, J., Jain, L., Karnin, Z., and Jamieson, K. An empirical process approach to the union bound: Practical algorithms for combinatorial and linear bandits. *arXiv preprint arXiv:2006.11685*, 2020.
- Koltchinskii, V. Rademacher complexities and bounding the excess risk in active learning. *The Journal of Machine Learning Research*, 11:2457–2485, 2010.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Ledoux, M. *The Concentration of Measure Phenomenon*. Mathematical surveys and monographs. American Mathematical Society, 2001. ISBN 9780821837924. URL https://books.google.com/books?id=mCX_cWL6rqwC.
- Musmann, S. and Liang, P. Uncertainty sampling is pre-conditioned stochastic gradient descent on zero-one loss. *arXiv preprint arXiv:1812.01815*, 2018.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images with unsupervised feature learning. 2011.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Settles, B. From theories to queries: Active learning in practice. In *Active Learning and Experimental Design workshop In conjunction with AISTATS 2010*, pp. 1–18, 2011.
- Soare, M., Lazaric, A., and Munos, R. Best-arm identification in linear bandits. In *Advances in Neural Information Processing Systems*, pp. 828–836, 2014.
- Talagrand, M. *Upper and lower bounds for stochastic processes: modern methods and classical problems*, volume 60. Springer Science & Business Media, 2014.
- Vershynin, R. *High-Dimensional Probability*. 2019.

Xiao, H., Rasul, K., and Vollgraf, R. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

Zhang, C. and Chaudhuri, K. Beyond disagreement-based agnostic active learning. *Advances in Neural Information Processing Systems*, 27:442–450, 2014.

A. Generalization

Proof of Remark 1. Define the event

$$\mathcal{E} := \{\forall h \in \mathcal{H} : |\text{err}(h) - \mathbb{P}(h(x) \neq f(x))| \leq c[\sqrt{\frac{\ln(1/\delta)}{n}} + \sqrt{\frac{d}{n}}]\}.$$

where the randomness is over the draw of the pool $\{x_1, \dots, x_n\} \sim \mathcal{D}_{\mathcal{X}}$ and c is a universal positive constant. Using the bounded differences inequality and 3.4 in (Boucheron et al., 2005), we have that if c is a sufficiently large universal positive constant, then by a standard argument $\mathbb{P}(\mathcal{E}) \geq 1 - \delta$. Suppose \mathcal{E} holds for the remainder of the proof. Let $\bar{h} = \arg \min_{h \in \mathcal{H}} \mathbb{P}_{(x,y) \sim \mathcal{D}}(h(x) \neq f(x))$. Then,

$$\begin{aligned} \mathbb{P}(\hat{h}(x) \neq f(x)) &\leq \text{err}(\hat{h}) + c[\sqrt{\frac{\ln(1/\delta)}{n}} + \sqrt{\frac{d}{n}}] \\ &\leq \min_{h \in \mathcal{H}} \text{err}(h) + \epsilon + c[\sqrt{\frac{\ln(1/\delta)}{n}} + \sqrt{\frac{d}{n}}] \\ &\leq \text{err}(\bar{h}) + \epsilon + c[\sqrt{\frac{\ln(1/\delta)}{n}} + \sqrt{\frac{d}{n}}] \\ &\leq \mathbb{P}(\bar{h}(x) \neq f(x)) + \epsilon + 2c[\sqrt{\frac{\ln(1/\delta)}{n}} + \sqrt{\frac{d}{n}}]. \end{aligned}$$

□

B. Reduction to Combinatorial Bandits

We state and prove our results for active classification in the language of combinatorial bandits, a strictly more general problem, which we now introduce.

Combinatorial Bandits: There are n distributions ν_1, \dots, ν_n supported on $[-1, 1]$ with mean $\mu_i = \mathbb{E}_{\zeta \sim \nu_i} \zeta$. \mathcal{H} is a collection of subsets of $[n]$. At each round t , the agent queries a distribution (or arm) I_t and observes $y_t \sim \nu_{I_t}$. Given $\epsilon > 0, \delta \in (0, 1)$, the goal is to identify $h \in \mathcal{H}$ that satisfies

$$\sum_{i \in h} \mu_i \geq \sum_{i \in h_*} \mu_i - \epsilon$$

with probability at least $1 - \delta$ using as few samples as possible.

We also write $\mu := (\mu_1, \dots, \mu_n)^\top$. We interchangeably treat each $h \in \mathcal{H}$ as a set in $[n]$ or as a vector in $\{0, 1\}^n$ with $h_i = 1$ if $i \in h$ and $h_i = 0$ otherwise. Using this vector notation, we often write $h^\top \mu = \sum_{i \in h} \mu_i$. We use the notation

$$\Delta_h := h_*^\top \mu - h^\top \mu$$

where $h_* \in \arg \max_{h \in \mathcal{H}} h^\top \mu$.

Reduction to combinatorial bandits: We use the reduction of active classification to combinatorial bandits from (Jain & Jamieson, 2019). Note that

$$\text{err}(h) = \frac{1}{n} \left[\sum_{i \in [n]: h(x_i)=0} \eta_i + \sum_{i \in [n]: h(x_i)=1} (1 - \eta_i) \right] = \frac{1}{n} \left[\sum_{i \in [n]} \eta_i - \sum_{i \in [n]: h(x_i)=1} \mu_i \right]$$

where $\mu_i := 2\eta_i - 1$. Thus, treating each $h \in \mathcal{H}$ as a set where $i \in h$ iff $h(x_i) = 1$, we observe $\arg \min_{h \in \mathcal{H}} \text{err}(h) = \arg \max_{h \in \mathcal{H}} \sum_{i \in h} \mu_i$ and that finding a hypothesis h such that

$$\text{err}(h) - \min_{h' \in \mathcal{H}} \text{err}(h') \leq \epsilon$$

is equivalent to finding h such that $\sum_{i \in h} \mu_i \geq \max_{h' \in \mathcal{H}} \sum_{i \in h'} \mu_i - n\epsilon$. Thus, active classification can be viewed as an instance of combinatorial bandits where each ν_i is a random variable supported on $\{-1, 1\}$ with mean $\mu_i = 2\eta_i - 1$. Therefore, any algorithm for ϵ -good arm identification for combinatorial bandits yields an algorithm for active binary classification in the pool-based setting. Finally, we note that

$$\Delta_h = n[\text{err}(h) - \text{err}(h_*)]$$

C. Disagreement Coefficient

We now introduce equivalent definitions of $\rho^*(\epsilon)$, $\gamma^*(\epsilon)$ and the disagreement coefficient in the combinatorial bandit setting.

$$\rho^*(\epsilon) := \inf_{\lambda \in \Delta_n} \sup_{h \in \mathcal{H} \setminus \{h_*\}} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{\max(\mu^\top(h_* - h), \epsilon)^2}$$

$$\gamma^*(\epsilon) := \inf_{\lambda \in \Delta_n} \mathbb{E}_{\zeta \sim N(0, I)} \left[\sup_{h \in \mathcal{H} \setminus \{h_*\}} \frac{(h_* - h)^\top A(\lambda)^{-1/2} \zeta}{\max(\mu^\top(h_* - h), \epsilon)} \right]^2.$$

Define the ball of radius r centered at h_*

$$B(h_*, r) = \{h \in \mathcal{H} : \frac{|h \Delta h_*|}{n} \leq r\}$$

and

$$\text{DIS}(B(h_*, r)) = \{i : \exists h, h' \in B(h_*, r) \text{ s.t. } i \in h \Delta h'\}.$$

The disagreement coefficient is defined as

$$\theta(\epsilon) = \sup_{r \geq \epsilon} \frac{|\text{DIS}(B(h_*, r))|/n}{r}$$

$$= \sup_{r \geq \epsilon} \frac{|\{i : i \in h \Delta h' \text{ for some } h, h' \in \mathcal{H} \text{ s.t. } \max(|h_* \Delta h|, |h_* \Delta h'|) \leq nr\}|}{nr}.$$

The proof of Proposition 2 follows by a peeling argument and the application of the following lemma.

Lemma 1. *Let $\epsilon \in [\frac{\Delta_{\min}}{n}, 1]$.*

- *Suppose the noiseless case holds, i.e., $\eta \in \{0, 1\}^n$. If $\epsilon \in [\frac{\Delta_{\min}}{n}, \nu)$, then*

$$\sup_{\xi \geq \epsilon} \min_{\lambda} \max_{h: \Delta_h \leq n\xi} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{(n\xi)^2} \leq 9\theta(\epsilon) \frac{\nu^2}{\epsilon^2}$$

and if $\epsilon \in [\nu, 1]$, then

$$\sup_{\xi \geq \epsilon} \min_{\lambda} \max_{h: \Delta_h \leq n\xi} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{(n\xi)^2} \leq 9\theta(\epsilon).$$

- *Suppose that the Tsabokov noise condition holds for some $a \in [1, \infty)$ and $\alpha \in (0, 1]$. Then,*

$$\min_{\lambda} \max_{h: \Delta_h \leq n\epsilon} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{(n\epsilon)^2} \leq ca^2 \frac{1}{\epsilon^{2-2\alpha}} \theta(a\epsilon^\alpha) \log(n\Delta_{\min}^{-1} \vee \epsilon^{-1})$$

Proof. Case 1: $\eta \in \{0, 1\}^n$ (Noiseless). We begin by noting that when $\eta \in \{0, 1\}^n$, we have the following equality that we will use repeatedly:

$$\text{err}(h) = \frac{1}{n} \left[\sum_{i \in [n]: h(x_i)=0} \eta_i + \sum_{i \in [n]: h(x_i)=1} (1 - \eta_i) \right] = \frac{1}{n} |\eta \Delta h|.$$

Then,

$$\text{err}(h) = \frac{1}{n} |\eta \Delta h| = \frac{1}{n} \left[\sum_{i \in \eta \setminus h} 1 + \sum_{i \in h \setminus \eta} 1 \right] = \frac{1}{n} \left[\sum_{i \in \eta} 1 + \sum_{i \in h \setminus \eta} 1 - \sum_{i \in \eta \cap h} 1 \right] = \frac{1}{n} |\eta| - \frac{1}{n} \sum_{i \in h} \mu_i \quad (7)$$

where we used $\mu_i = 2\eta_i - 1$.

Recall $\nu = \text{err}(h_*)$. Fix $\xi \geq \epsilon$. Suppose $\Delta_h \leq n\xi$. We have by (7) $\frac{1}{n} \sum_{i \in h} \mu_i = \frac{1}{n} |\eta| - \frac{1}{n} |\eta \Delta h|$ and thus

$$n\xi \geq \Delta_h = |\eta \Delta h| - |\eta \Delta h_*| = |\eta \Delta h| - n\nu$$

and thus $|\eta \Delta h| \leq n(\xi + \nu)$. Thus, if $\Delta_h \leq n\xi$, we have that

$$|h_* \Delta h| \leq |h_* \Delta \eta| + |h \Delta \eta| \leq n(2\nu + \xi) \quad (8)$$

Furthermore, by (8)

$$\begin{aligned} \min_{\lambda} \max_{h: \Delta_h \leq n\xi} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{(n\xi)^2} &= \min_{\lambda} \max_{h: \Delta_h \leq n\xi} \frac{\sum_{i \in h_* \Delta h} \frac{1}{\lambda_i}}{(n\xi)^2} \\ &\leq \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } \Delta_h \leq n\xi\}| \cdot \max_{h: \Delta_h \leq n\xi} |h_* \Delta h|}{(n\xi)^2} \\ &\leq \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } |h_* \Delta h| \leq n(2\nu + \xi)\}| \cdot n(2\nu + \xi)}{(n\xi)^2}. \end{aligned} \quad (9)$$

where the first inequality takes λ to be the uniform distribution over $|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } \Delta_h \leq n\xi\}|$.

Case 1.1. Suppose $\epsilon \geq \nu$. Then, $\xi \geq \epsilon \geq \nu$, and we have

$$\begin{aligned} &\frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } |h_* \Delta h| \leq n(2\nu + \xi)\}| n(2\nu + \xi)}{(n\xi)^2} \\ &\leq 3 \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } |h_* \Delta h| \leq 3n\xi\}|}{n\xi} \\ &= 9 \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } |h_* \Delta h| \leq 3n\xi\}|}{3n\xi} \end{aligned}$$

and, using (9) in addition,

$$\begin{aligned} \sup_{\xi \geq \epsilon} \min_{\lambda} \max_{h: \Delta_h \leq n\xi} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{(n\xi)^2} &\leq \sup_{\xi \geq \epsilon} 9 \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } |h_* \Delta h| \leq 3n\xi\}|}{3n\xi} \\ &\leq 9\theta(3\epsilon) \\ &\leq 9\theta(\epsilon). \end{aligned}$$

Case 1.2. Now, suppose $\epsilon \in [\frac{\Delta_{\min}}{n}, \nu]$. We may suppose wlog that $\xi \geq \epsilon$ satisfies $\xi \in [\frac{\Delta_{\min}}{n}, \nu]$ since otherwise it reduces to case 1.1. Then,

$$\begin{aligned} &\frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } |h_* \Delta h| \leq n(2\nu + \xi)\}| n(2\nu + \xi)}{(n\xi)^2} \\ &= \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } |h_* \Delta h| \leq n(2\nu + \xi)\}| n(2\nu + \xi)}{(n\nu)^2} \frac{\nu^2}{\xi^2} \\ &\leq 3 \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } |h_* \Delta h| \leq 3n\nu\}|}{n\nu} \frac{\nu^2}{\xi^2} \\ &\leq 9 \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } |h_* \Delta h| \leq 3n\nu\}|}{3n\nu} \frac{\nu^2}{\xi^2} \\ &\leq 9\theta(3\nu) \frac{\nu^2}{\xi^2} \\ &\leq 9\theta(\xi) \frac{\nu^2}{\xi^2}. \end{aligned}$$

Combining this with (9), this implies that

$$\begin{aligned} \sup_{\xi \geq \epsilon} \min_{\lambda} \max_{h: \Delta_h \leq n\xi} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{(n\xi)^2} &\leq \sup_{\xi \geq \epsilon} \min_{\lambda} \max_{h: \Delta_h \leq n\xi} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{(n\xi)^2} \\ &\leq 9\theta(\epsilon) \frac{\nu^2}{\epsilon^2}. \end{aligned}$$

Case 2: Tsabakov Noise. Suppose Tsabakov's noise condition is satisfied with $a \in [1, \infty)$ and $\alpha \in (0, 1]$. Fix $\xi \geq \epsilon$. If $\Delta_h \leq n\xi$, then Tsabakov's noise condition implies that

$$\frac{|h_* \Delta h|}{n} \leq a \left(\frac{\Delta_h}{n}\right)^\alpha \leq a\xi^\alpha$$

Then,

$$\begin{aligned} \min_{\lambda} \max_{h: \Delta_h \leq n\xi} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{(n\xi)^2} &= \min_{\lambda} \max_{h: \Delta_h \leq n\xi} \frac{\sum_{i \in h_* \Delta h} \frac{1}{\lambda_i}}{(n\xi)^2} \\ &\leq \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } \Delta_h \leq n\xi\}| \cdot \max_{h: \Delta_h \leq n\xi} |h_* \Delta h|}{(n\xi)^2} \\ &\leq \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } \frac{|h_* \Delta h|}{n} \leq a\xi^\alpha\}| \cdot \max_{h: \Delta_h \leq n\xi} |h_* \Delta h|}{(n\xi)^2} \\ &\leq \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } \frac{|h_* \Delta h|}{n} \leq a\xi^\alpha\}| a n \xi^\alpha}{(n\xi)^2} \\ &= a^2 \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } \frac{|h_* \Delta h|}{n} \leq a\xi^\alpha\}|}{a \xi^\alpha n} \frac{1}{\xi^{2-2\alpha}} \\ &\leq a^2 \frac{1}{\xi^{2-2\alpha}} \theta(a\xi^\alpha) \end{aligned}$$

Taking the sup over $\xi \geq \epsilon$ of both sides implies the result. □

Proof of Proposition 2. The argument follows by a peeling argument. Define

$$\begin{aligned} \lambda^{(k)} &:= \arg \min_{\lambda} \max_{h: \Delta_h \in (n2^{-k-1}, n2^{-k}]} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{\Delta_h^2} \\ \bar{\lambda} &:= \frac{1}{\lceil \log_2(n\Delta_{\min}^{-1} \vee \epsilon^{-1}) \rceil} \sum_{k=0}^{\lceil \log_2(n\Delta_{\min}^{-1} \vee \epsilon^{-1}) \rceil} \lambda^{(k)}. \end{aligned}$$

Notice that $\frac{1}{\lceil \log_2(n\Delta_{\min}^{-1} \vee \epsilon^{-1}) \rceil} A(\lambda^{(k)}) \preceq A(\bar{\lambda})$, which implies that

$$A(\bar{\lambda})^{-1} \preceq \lceil \log_2(n\Delta_{\min}^{-1} \vee \epsilon^{-1}) \rceil A(\lambda^{(k)})^{-1}. \quad (10)$$

Thus,

$$\begin{aligned}
 \rho^*(n\epsilon) &= \min_{\lambda} \max_{h \in \mathcal{H} \setminus \{h_*\}} \frac{\|h_* - h\|_{A(\lambda)}^2}{\max(n\epsilon, \Delta_h)^2} \\
 &\leq \max_{h \in \mathcal{H} \setminus \{h_*\}} \frac{\|h_* - h\|_{A(\bar{\lambda})}^2}{\max(n\epsilon, \Delta_h)^2} \\
 &= \max_{k=0,1,\dots, \lceil \log_2(n\Delta_{\min}^{-1} \vee \epsilon^{-1}) \rceil} \max_{h: \Delta_h \in (n2^{-k-1}, n2^{-k}]} \frac{\|h_* - h\|_{A(\bar{\lambda})}^2}{\max(n\epsilon, \Delta_h)^2} \\
 &\leq \lceil \log_2(n\Delta_{\min}^{-1} \vee \epsilon^{-1}) \rceil \max_{k=0,1,\dots, \lceil \log_2(n\Delta_{\min}^{-1} \vee \epsilon^{-1}) \rceil} \max_{h: \Delta_h \in (n2^{-k-1}, n2^{-k}]} \frac{\|h_* - h\|_{A(\lambda^{(k)})}^2}{\max(n\epsilon, \Delta_h)^2} \tag{11}
 \end{aligned}$$

$$\begin{aligned}
 &= \lceil \log_2(n\Delta_{\min}^{-1} \vee \epsilon^{-1}) \rceil \max_{k=0,1,\dots, \lceil \log_2(n\Delta_{\min}^{-1} \vee \epsilon^{-1}) \rceil} \min_{\lambda} \max_{h: \Delta_h \in (n2^{-k-1}, n2^{-k}]} \frac{\|h_* - h\|_{A(\lambda^{(k)})}^2}{\max(n\epsilon, \Delta_h)^2} \tag{12} \\
 &\leq 2 \lceil \log_2(n\Delta_{\min}^{-1} \vee \epsilon^{-1}) \rceil \sup_{\xi \geq \epsilon \vee \frac{\Delta_{\min}}{n}} \min_{\lambda} \max_{h: \Delta_h \leq n\xi} \frac{\|h_* - h\|_{A(\lambda)}^2}{(n\xi)^2}
 \end{aligned}$$

where inequality (11) follows by (10) and (12) follows by the definition of $\bar{\lambda}^{(k)}$.

The result now follows by applying Lemma 1 in each case (noiseless $\eta \in \{0, 1\}$ with $\epsilon > \nu$ and $\epsilon < \nu$, and Tsbokov noise condition). \square

Proof of Proposition 3. Step 1: Define the instance. Let $m \in \mathbb{N}$. Define $h_i = [m] \cup \{m+i\}$ for $i = 1, \dots, m^2$ and let $n = m + m^2$. Define $h_0 = \emptyset$. Let $\mu_i = -1$ for all $i \in [n]$. Note that h_0 is the best set and that $\mu^\top(h_0 - h_i) = m + 1$ for all $i \neq 0$.

Step 2: Compute problem-dependent quantities. We have that

$$\rho^* = \inf_{\lambda} \max_{i=1,\dots,m^2} \frac{\sum_{j \in [m] \cup \{m+i\}} \frac{1}{\lambda_j}}{(m+1)^2} \leq \frac{2m^2 + 2m^2}{(m+1)^2} \leq O(1)$$

where we used

$$\lambda_i = \begin{cases} \frac{1}{2m} & i \in [m] \\ \frac{1}{2m^2} & i \in \{m+1, \dots, m+m^2\} \end{cases}$$

Let ξ such that $n\xi \leq m + 1$. Letting $r = m + 1$, we have that

$$\begin{aligned}
 \theta(\xi) &\geq \frac{|\{i : i \in h_* \Delta h \text{ for some } h \in \mathcal{H} \text{ s.t. } |h_* \Delta h| \leq r\}|}{r} \\
 &= \frac{m + m^2}{m + 1} \\
 &\geq \frac{1}{2}m \\
 &\geq \frac{1}{2\sqrt{2}}\sqrt{n}.
 \end{aligned}$$

\square

D. Ridge IPS Estimator Tail Bound

Now, we introduce the ridge IPS estimator, which we leverage in constructing our generic chaining estimator.

Proposition 5. Fix $\mathcal{X} = \{\mathbf{e}_i : i \in [n]\}$ and $\mathcal{H} \subset \{0, 1\}^n$ as well as some $\mu \in [-1, 1]^n$. Let $v \in \mathbb{R}^n$. Fix some $\lambda \in \Delta_n$ and draw $\{x_s\}_{i=1}^t \sim \lambda$ and then observe y_s with mean $x_s^\top \mu$ and $|y_i| \leq 1$ with probability 1 for $s = 1, \dots, t$.

For any $\alpha \in \mathbb{R}_+^n$ define

$$A(\alpha) := \sum_{i=1}^n \alpha_i \mathbf{e}_i \mathbf{e}_i^\top$$

For some $s > 0$ define

$$\hat{\mu} = (A(t\lambda) + sI)^{-1} X^\top y.$$

where X, y, ϵ are the $\{(x_i, y_i, \epsilon_i)\}_i$ stacked. For any $s > 0$ let $A(\alpha + s) := A(\alpha) + sI_n$. If we take $s = \sqrt{\frac{\log(2/\delta)}{3\|v\|_{(nA(\lambda))^{-1}}^2}}$ then

$$|\langle v, \hat{\mu} - \mu \rangle| \leq (\sqrt{2/3} + 1) \sqrt{\frac{2\|v\|_{A(\lambda)^{-1}}^2 \log(2/\delta)}{t}}$$

Proof of Proposition 5. Fix $\mathcal{X} = \{\mathbf{e}_i : i \in [n]\}$ and $\mathcal{H} \subset \{0, 1\}^n$ as well as some $\mu \in [-1, 1]^n$. Note that for any $v \in \{-1, 0, 1\}^n$ we have

$$\begin{aligned} |\mathbb{E}[\langle v, \hat{\mu} - \mu \rangle]| &= |\langle v, (A(t\lambda) + sI)^{-1} A(t\lambda) \mu - \mu \rangle| \\ &= |\langle v, (A(t\lambda) + sI)^{-1} (A(t\lambda) + sI - sI) \mu - \mu \rangle| \\ &= |s \langle v, (A(t\lambda) + sI)^{-1} \mu \rangle| \\ &\leq s \|v\|_{(tA(\lambda) + sI)^{-1}}^2 \end{aligned}$$

using the fact that $\mu \in [-1, 1]^n$. Define

$$\begin{aligned} S &= \langle v, \hat{\mu} \rangle \\ &= \sum_{i=1}^n v^\top (A(t\lambda) + sI)^{-1} x_i y_i \\ &=: \sum_{i=1}^n X_i \end{aligned}$$

Note that

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}[X_i^2] &= \sum_{i=1}^n \mathbb{E}[\langle v, (A(t\lambda) + sI)^{-1} x_i y_i \rangle^2] \\ &\leq v^\top (A(t\lambda) + sI)^{-1} A(t\lambda) (A(t\lambda) + sI)^{-1} v \\ &\leq \|v\|_{(nA(\lambda) + sI)^{-1}}^2 \end{aligned}$$

and

$$|X_i| \leq |\langle v, (A(t\lambda) + sI)^{-1} x_i y_i \rangle| \leq 1/s$$

for all i . We have by Bernstein's inequality that with probability at least $1 - \delta$

$$\sum_{i=1}^n X_i - \mathbb{E}[X_i] \leq \sqrt{2\|v\|_{(nA(\lambda) + sI)^{-1}}^2 \log(1/\delta)} + \frac{\log(1/\delta)}{3s}.$$

Thus,

$$\begin{aligned}
 \langle v, \hat{\mu} \rangle &\leq \mathbb{E}[\langle v, \hat{\mu} \rangle] + \sqrt{2\|v\|_{(nA(\lambda)+sI)^{-1}}^2 \log(1/\delta)} + \frac{\log(1/\delta)}{3s} \\
 &= \langle v, \mu \rangle + \mathbb{E}[\langle v, \hat{\mu} - \mu \rangle] + \sqrt{2\|v\|_{(nA(\lambda)+sI)^{-1}}^2 \log(1/\delta)} + \frac{\log(1/\delta)}{3s} \\
 &\leq \langle v, \mu \rangle + s\|v\|_{(nA(\lambda)+sI)^{-1}}^2 + \sqrt{2\|v\|_{(nA(\lambda)+sI)^{-1}}^2 \log(1/\delta)} + \frac{\log(1/\delta)}{3s}
 \end{aligned}$$

from which we conclude, that with probability at least $1 - 2\delta$

$$\begin{aligned}
 |\langle v, \hat{\mu} - \mu \rangle| &\leq s\|v\|_{(nA(\lambda)+sI)^{-1}}^2 + \frac{\log(2/\delta)}{3s} + \sqrt{2\|v\|_{(nA(\lambda)+sI)^{-1}}^2 \log(2/\delta)} \\
 &\leq s\|v\|_{(nA(\lambda))^{-1}}^2 + \frac{\log(2/\delta)}{3s} + \sqrt{2\|v\|_{(nA(\lambda))^{-1}}^2 \log(2/\delta)}.
 \end{aligned}$$

If we take $s = \sqrt{\frac{\log(2/\delta)}{3\|v\|_{(nA(\lambda))^{-1}}^2}}$ then

$$|\langle v, \hat{\mu} - \mu \rangle| \leq (\sqrt{2/3} + 1) \sqrt{\frac{2\|v\|_{A(\lambda)^{-1}}^2 \log(2/\delta)}{n}}$$

□

E. Looseness of Bernstein's Bound for Importance Sampling Estimator

A natural approach to combinatorial bandits is to use the importance sampling estimator and apply Bernstein's bound in an algorithm like RAGE from (Fiez et al., 2019). Applying the standard analysis would yield a term in the sample complexity that scales as:

$$\inf_{\lambda} \max_{h \neq h_*} \max_{j \in h_* \Delta h} \frac{1}{\lambda_j} \frac{1}{\sum_{k \in h_*} \mu_k - \sum_{k \in h} \mu_k}.$$

The following proposition shows that there exists instances where such a sample complexity is suboptimal by a polynomial factor in the dimension.

Proposition 6. *There exists a combinatorial bandit problem where $\rho^* = O(1)$ and*

$$\inf_{\lambda} \max_{h \neq h_*} \max_{j \in h_* \Delta h} \frac{1}{\lambda_j} \frac{1}{\sum_{k \in h_*} \mu_k - \sum_{k \in h} \mu_k} \geq \Omega(\sqrt{n})$$

Proof. Step 1: Define the instance. Let $m \in \mathbb{N}$. Define $h_i = [m] \cup \{m+i\}$ for $i = 1, \dots, m^2$ and let $n = m + m^2$. Define $h_0 = \emptyset$. Let $\mu_i = -1$ for all $i \in [n]$. Note that $\mu^\top(h_0 - h_i) = m + 1$ for all $i \neq 0$.

Step 2: Compute problem-dependent quantities. Then,

$$\rho^* = \inf_{\lambda} \max_{i=1, \dots, m^2} \frac{\sum_{j \in [m] \cup \{m+i\}} \frac{1}{\lambda_j}}{(m+1)^2} \leq \frac{2m^2 + 2m^2}{(m+1)^2} \leq O(1)$$

where we used

$$\lambda_i = \begin{cases} \frac{1}{2m} & i \in [m] \\ \frac{1}{2m^2} & i \in \{m+1, \dots, m+m^2\} \end{cases}.$$

On the other hand,

$$\begin{aligned}
 \inf_{\lambda} \max_{i=1, \dots, m^2} \max_{j \in [m] \cup \{m+i\}} \frac{1}{\lambda_j} &\geq \inf_{\lambda} \max_{i=1, \dots, m^2} \max_{j=m+i} \frac{1}{\lambda_j} \\
 &= \frac{m^2}{m+1} \\
 &\geq \frac{1}{2}m \\
 &\geq \frac{1}{2\sqrt{2}}\sqrt{n}
 \end{aligned}$$

□

F. Proof of Theorem 1

Before proving Theorem 1, we introduce some machinery from the theory of generic chaining (see e.g. (Vershynin, 2019) for more details). Fix a set $T \subset \mathbb{R}^n$. Consider a sequence of subset $(T_k)_{k=0}^{\infty}$ such that $T_k \subset T$

$$|T_0| = 1, \quad |T_k| \leq 2^{2^k}.$$

A sequence $(T_k)_{k=0}^{\infty}$ satisfying the above properties is called *admissible*.

Definition 1. Let d be a metric on \mathbb{R}^n . The γ_2 -functional of T is defined as

$$\gamma_2(T, d) = \inf_{(T_k)} \sup_{t \in T} \sum_{k=0}^{\infty} 2^{k/2} d(t, T_k)$$

where the infimum is taken over all admissible sequences.

Here, we state and prove Theorem 4, the combinatorial bandit counterpart of Theorem 1. The proof uses the technique of generic chaining to avoid a naive union bound, which would introduce a dependence on $\log(|\mathcal{G}|)$. Unfortunately, the IPS estimator has an excessively large sub-Gaussian norm, and the concentration inequality for ridge IPS estimator from Proposition 5 decays at a suitably fast sub-Gaussian rate for only a subset of pairs $h, h' \in \mathcal{G}$ —not all pairs—implying that neither of these estimators can be used directly. To sidestep this issue, we apply the estimator from Proposition 5 to all pairs of $h, h' \in \mathcal{G}$ to construct a feasibility program that yields the estimator in Theorem 4. The proof has similarities to proof techniques in the theory of generic chaining (Vershynin, 2019).

Theorem 4. Let $\mathcal{G} \subset \mathcal{H}$. Fix some $\lambda \in \Delta_n$ and draw $\{x_s\}_{i=1}^t \sim \lambda$ and then observe y_s with mean $x_s^\top \mu$ and $|y_s| \leq 1$ with probability 1 for $s = 1, \dots, t$. There exists an estimator $\hat{\mu} \in [-1, 1]^n$ such that with probability at least $1 - \delta$,

$$\sup_{h, h' \in \mathcal{G}} |(h - h')^\top (\hat{\mu} - \mu)| \leq c[\log(1/\delta)] \max_{h, h' \in \mathcal{G}} \|h - h'\|_{A(t\lambda)^{-1}} + \mathbb{E}_{\zeta \sim N(0, I)} [\sup_{h \in \mathcal{G}} h^\top A(t\lambda)^{-1/2} \zeta]$$

Proof of Theorem 4. Step 1: Pick the admissible sequence. Fix $h_0 \in \mathcal{G}$. Note that for any $\tilde{\mu} \in \mathbb{R}^n$

$$\sup_{h, h' \in \mathcal{G}} |(h - h')^\top (\tilde{\mu} - \mu)| \leq 2 \sup_{h \in \mathcal{G}} |(h - h_0)^\top (\tilde{\mu} - \mu)| \quad (13)$$

and thus we focus on upper bounding the RHS.

Let $(\mathcal{R}_k)_{k=0}^K$ be an admissible sequence of \mathcal{G} where $\mathcal{R}_0 = \{h_0\}$ and $\mathcal{R}_K = \mathcal{G}$ and $\mathcal{R}_k \subset \mathcal{G}$ such that

$$\sup_{h \in \mathcal{G}} \sum_{k=1}^K 2^{k/2} \inf_{h' \in \mathcal{R}_k} \|h - h'\|_{A(t\lambda)^{-1}} \leq 2\gamma(\mathcal{G}, \|\cdot\|_{A(t\lambda)^{-1}})$$

Let $\mathcal{T}_k = \cup_{l=1}^k \mathcal{R}_l$. Note that $\mathcal{T}_1 \subset \mathcal{T}_2 \subset \dots \subset \mathcal{T}_K$ and $K = \log(\log(|\mathcal{G}|)) \leq \log(n)$.

Step 2: Defining the event. Let X , and y be the $\{(x_s, y_s)\}_{s=1}^t$ stacked, respectively. For each $h, h' \in \mathcal{G}$ and $k \in [K]$ such that $h, h' \in \mathcal{T}_k$, define the estimator

$$\hat{\mu}_{h,h',k} = (A(t\lambda) + s_{h,h',k}I)^{-1}X^\top y$$

where $s_{h,h',k} = \sqrt{1/3} \frac{u+2^{k/2}}{\|h-h'\|_{A(t\lambda)^{-1}}}$ where $s_{h,h',k}$ is chosen based on Proposition 5. Define the events

$$\mathcal{E}_{h,h',k} = \{ |(h-h')^\top (\hat{\mu}_{h,h',k} - \mu)| \leq (\sqrt{2/3} + 1)\sqrt{4}(u + 2^{k/2}) \|h-h'\|_{A(t\lambda)^{-1}} \} \quad (14)$$

$$\mathcal{E} = \bigcap_{k=1}^K \bigcap_{h,h' \in \mathcal{T}_k} \mathcal{E}_{h,h',k} \quad (15)$$

Applying the bound for the Ridge IPS from Proposition 5 and rearranging, for every $k \in [K]$ and $h, h' \in \mathcal{T}_k$,

$$\Pr(|(h-h')^\top (\hat{\mu}_{h,h',k} - \mu)| > (\sqrt{2/3} + 1)\sqrt{4}(u + 2^{k/2}) \|h-h'\|_{A(t\lambda)^{-1}}) \leq 2 \exp(-2(u + 2^{k/2})^2). \quad (16)$$

Then, by the union bound, we have that

$$\begin{aligned} \Pr(\mathcal{E}^c) &\leq \sum_{k \geq 1} |\mathcal{T}_k|^2 \Pr(\mathcal{E}_{h,h',k}^c) \\ &\leq c \sum_{k \geq 1} |\mathcal{T}_k|^2 \exp(-2(2^k + u^2)) \end{aligned} \quad (17)$$

$$\begin{aligned} &\leq c \sum_{k \geq 1} 2^{2k+1} \exp(-2(2^k + u^2)) \\ &\leq c' \exp(-2u^2). \end{aligned} \quad (18)$$

where (17) follows by (16) and where line (18) follows since by construction

$$|\mathcal{T}_k| \leq \sum_{l=1}^k |\mathcal{R}_k| \leq \sum_{l=1}^k 2^{2l} \leq 2^{2k+1}.$$

Suppose \mathcal{E} holds for the remainder of the proof with u taking the value of $\bar{u} = \sqrt{\frac{\log(c'/\delta)}{2}}$.

Step 3: Define the estimator. Define the polyhedron

$$\hat{\mu} \in P = \{z \in [-1, 1]^n : \forall k \in [K], \forall h, h' \in \mathcal{T}_k : |(h-h')^\top \hat{\mu}_{h,h',k} - z| \leq c[\bar{u} + 2^{k/2}] \|h-h'\|_{A(t\lambda)^{-1}}\}.$$

We define the estimator $\hat{\mu}$ to be any point in P if it is nonempty and, otherwise we let $\hat{\mu}$ be any point in \mathbb{R}^n .

Note that on the event \mathcal{E} , $\mu \in P$ and hence P is nonempty. Furthermore, by the triangle inequality, $\forall k \in [K], \forall h, h' \in \mathcal{T}_k$

$$\begin{aligned} |(h-h')^\top (\hat{\mu} - \mu)| &\leq |(h-h')^\top (\hat{\mu} - \hat{\mu}_{h,h',k})| + |(h-h')^\top (\hat{\mu}_{h,h',k} - \mu)| \\ &\leq 2c(\bar{u} + 2^{k/2}) \|h-h'\|_{A(t\lambda)^{-1}}. \end{aligned}$$

Step 4: Proving the inequality. Fix $\bar{h} \in \mathcal{G}$. Let $k(\bar{h}, l)$ be the smallest integer such that

$$d(\bar{h}, \mathcal{T}_{k(\bar{h}, l)}) \leq \frac{d(\bar{h}, \mathcal{T}_{k(\bar{h}, l-1)})}{2}.$$

Note that

$$d(\bar{h}, \mathcal{T}_{k(\bar{h}, l)}) \leq 2^{-l} \max_{h, h' \in \mathcal{G}} \|h-h'\|_{A(t\lambda)^{-1}}.$$

Let $\pi_l(\bar{h}) \in \mathcal{T}_{k(\bar{h}, l)}$ such that

$$\left\| \bar{h} - \pi_l(\bar{h}) \right\|_{A(t\lambda)^{-1}} = \min_{h \in \mathcal{T}_{k(\bar{h}, l)}} \left\| \bar{h} - h \right\|_{A(t\lambda)^{-1}}.$$

By the triangle inequality

$$\begin{aligned} |(\bar{h} - h_0)^\top (\hat{\mu} - \mu)| &\leq \sum_{l \geq 1} |(\pi_l(\bar{h}) - \pi_{l-1}(\bar{h}))^\top (\hat{\mu} - \mu)| \\ &\leq 2c \sum_{l \geq 1} (\bar{u} + 2^{k/2}) \left\| \pi_l(\bar{h}) - \pi_{l-1}(\bar{h}) \right\|_{A(t\lambda)^{-1}} \\ &\leq c' \sum_{l \geq 1} (\bar{u} + 2^{k/2}) \left\| \pi_l(\bar{h}) - \bar{h} \right\|_{A(t\lambda)^{-1}} \end{aligned} \quad (19)$$

where we use event \mathcal{E} and the triangle inequality. We have by construction that

$$\begin{aligned} \sum_{l \geq 1} \bar{u} \left\| \pi_l(\bar{h}) - \bar{h} \right\|_{A(t\lambda)^{-1}} &\leq \bar{u} \sum_{l \geq 1} 2^{-l} \max_{h, h' \in \mathcal{G}} \left\| h - h' \right\|_{A(t\lambda)^{-1}} \\ &\leq c' \bar{u} \max_{h, h' \in \mathcal{G}} \left\| h - h' \right\|_{A(t\lambda)^{-1}}. \end{aligned} \quad (20)$$

Furthermore,

$$\sum_{l \geq 1} 2^{k/2} \left\| \pi_l(\bar{h}) - \bar{h} \right\|_{A(t\lambda)^{-1}} \leq c' \gamma(\mathcal{G}, \left\| \cdot \right\|_{A(t\lambda)^{-1}}) \quad (21)$$

$$\leq c'' \mathbb{E}_{\zeta \sim N(0, I)} [\sup_{h \in \mathcal{G}} h^\top A(t\lambda)^{-1/2} \zeta] \quad (22)$$

where (21) follows by the definition of \mathcal{T}_k and (22) follows by Talagrand's majorizing measure theorem (Theorem 8.6.1 in (Vershynin, 2019)). Putting together (13), (19), (20), and (22), and noting that \bar{h} is arbitrary, the result follows. \square

Remark 3. We emphasize that the construction of this estimator does not use any knowledge of μ . The admissible sequence $(\mathcal{R}_k)_{k \in \mathbb{N}}$ does not require knowledge of μ to be chosen and the polyhedron P can be defined without knowledge of μ .

G. Fixed Confidence Algorithms

We restate Algorithm 1 in the language of combinatorial bandits.

Algorithm 3 ACED for Combinatorial Bandits.

Input: Confidence level $\delta \in (0, 1)$.

$\mathcal{H}_1 \leftarrow \mathcal{H}$, $k \leftarrow 1$, $\delta_k \leftarrow \delta/2k^2$.

while $|\mathcal{H}_k| > 1$ **do**

Let λ_k and τ_k be the solution and value of the following optimization problem

$$\inf_{\lambda \in \Delta_n} \mathbb{E}_{\zeta \sim N(0, I)} [\max_{h, h' \in \mathcal{H}_k} (h - h')^\top A(\lambda)^{-1/2} \zeta]^2 + 2 \log\left(\frac{1}{\delta_k}\right) \max_{h, h' \in \mathcal{H}_k} \left\| h - h' \right\|_{A(\lambda)^{-1}}^2$$

Set $N_k \leftarrow c\tau_k \left(\frac{2^{k+1}}{n}\right)^2$ where c is a universal constant.

Query $I_1, \dots, I_{N_k} \sim \lambda_k$ and receive rewards y_1, \dots, y_{N_k} .

Let $\hat{\mu}_k$ be the estimator defined in Theorem 4.

$\mathcal{H}_{k+1} \leftarrow \mathcal{H}_k \setminus \{h \in \mathcal{H}_k : \exists h' \text{ such that } (h' - h)^\top \hat{\mu}_k - \frac{n}{2^{k+1}} \geq 0\}$.

$k \leftarrow k + 1$

end while

Return: $\mathcal{H}_k = \{\hat{h}\}$.

We now restate Theorem 2 as Theorem 5 in the combinatorial bandits setting.

Theorem 5. Let $\delta \in (0, 1)$ and $\epsilon > 0$. With probability at least $1 - \delta$ Algorithm 1 returns $\hat{h} \in \mathcal{H}$ such that $\mu^\top \hat{h} + \epsilon \geq \mu^\top h_*$ and uses at most

$$c \log(1/\epsilon) [\log(1/\delta) \rho^*(\epsilon) + \gamma^*(\epsilon)]$$

samples where c is a positive universal constant.

The proof of Theorem 5 is essentially identical to the proof of Theorem 4 in (Katz-Samuels et al., 2020), and we therefore omit it. The key technical and conceptual technical hurdle in obtaining Theorem 5 is the estimator given in Theorem 4.

H. Fixed Budget Proof

We restate the fixed budget algorithm, Algorithm 2, and the Theorem 3 in the language of combinatorial bandits as Algorithm 4 and Theorem 6, respectively.

Algorithm 4 Fixed Budget ACED for Combinatorial Bandits.

Input: Budget T , tolerance $\epsilon > 0$.

$$\hat{\mu}_1 = \mathbf{0} \in \mathbb{R}^n, N \leftarrow \lceil T / \log_2(n\epsilon^{-1}) \rceil.$$

for $k = 1, 2, \dots, \lceil \log_2(n\epsilon^{-1}) \rceil$ **do**

$$\tilde{h}_k \leftarrow \arg \max_{h \in \mathcal{H}} \hat{\mu}_k^\top h$$

Let λ_k be the solution of the following optimization problem

$$\inf_{\lambda \in \Delta} \mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1}n + \hat{\mu}_k^\top (\tilde{h}_k - h)} \right]^2 \quad (23)$$

Sample $\{x_1, \dots, x_N\} \sim \lambda_k$.

Query x_1, \dots, x_N and receive rewards y_1, \dots, y_N .

Let $\hat{\mu}_{k+1}$ be the estimator defined in the proof of Theorem 6.

end for

Return: $\arg \max_{h \in \mathcal{H}} \hat{\mu}_{k+1}^\top h$.

Theorem 6. Let $T \in \mathbb{N}$ and $\epsilon > 0$. Let \hat{h} denote the $h \in \mathcal{H}$ returned by Algorithm 2. If $T \geq \log(n\epsilon^{-1})[\gamma^*(\epsilon) + \rho^*(\epsilon)]$, then

$$\mathbb{P}(\mu^\top \hat{h} + \epsilon < \mu^\top h_*) \leq \log(n\epsilon^{-1})^2 \exp\left(-\frac{T}{\log(n\epsilon^{-1})[\gamma^*(\epsilon) + \rho^*(\epsilon)]}\right)$$

The key technical challenge in the proof (see Step 1) is constructing an estimator that concentrates rapidly enough. To this end, we leverage the estimator in Theorem 4, applying it to various subsets of \mathcal{H} based on their estimated gaps and combine them into a single estimator by defining it to belong to a polyhedron (denoted P in step 1.2.2) characterizing estimators with the suitable concentration properties. We argue that on a good event, the true mean μ belongs to P , making it feasible and thus obtaining our estimator. The next challenge is bounding the probability of error of our algorithm. The estimator constructed at round $k + 1$ is chosen to have failure probability

$$\delta_{k+1} = \exp\left(-c \frac{N}{\mathbb{E}[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2} \zeta}{(\tilde{h}_k - h)^\top \hat{\mu}_{k+2^{-k+1}n}}]^2}\right),$$

for a suitably large universal constant $c > 0$, which we note is a function of (23). Thus, this step of the proof (step 3) shows that $\delta_{k+1} \leq \exp\left(-\frac{T}{\log(n\epsilon^{-1})[\gamma^*(\epsilon) + \rho^*(\epsilon)]}\right)$. We note that while this step of the proof and algorithm style are novel to our knowledge, the mechanics of bounding the various quantities appearing in this step are quite similar to arguments in the proof of Theorem 5 in (Katz-Samuels et al., 2020).

Proof of Theorem 6. Step 1: Construction of the Estimator. Define the sets

$$S_k = \{h \in \mathcal{H} : \Delta_h \leq n2^{-k+1}\}.$$

Let $\widehat{\mu}_k$ be some estimator formed from data collected in k th round. Define the event

$$\begin{aligned} \mathcal{E}_k(\widehat{\mu}_k) &= \{\forall h \in S_k^c, |(h_* - h)^\top (\widehat{\mu}_k - \mu)| \leq \frac{\Delta_h}{8}\} \\ &\cap \{\forall h \in S_k, |(h_* - h)^\top (\widehat{\mu}_k - \mu)| \leq \frac{2^{-k+1}n}{8}\} \end{aligned}$$

We show that at every round k , we can construct an estimator $\widehat{\mu}_k \in [-1, 1]^n$ such that $\Pr(\mathcal{E}_k^c(\widehat{\mu}_k) | \mathcal{E}_{k-1}(\widehat{\mu}_{k-1}), \dots, \mathcal{E}_1(\widehat{\mu}_1)) \leq \delta_k \log(n\epsilon^{-1})$ where $\delta_k > 0$ will be chosen at the k th round.

Step 1.1: Base Case

Let $\widehat{\mu}_1$ be the estimator from Theorem 1 applied with δ_1 to \mathcal{H} (to be chosen later). Then, we have that

$$\begin{aligned} \sup_{h, h' \in \mathcal{H}} |(h - h')^\top (\widehat{\mu}_1 - \mu)| &\leq c[\log(2/\delta_1) \max_{h, h' \in \mathcal{H}} \|h - h'\|_{A(N\lambda)^{-1}} + \mathbb{E}_{\zeta \sim N(0, I)}[\sup_{h \in \mathcal{H}} h^\top A(N\lambda)^{-1/2} \zeta]] \\ &\leq c[\frac{\pi}{2} \log(2/\delta_1) \mathbb{E}[\sup_{h \in \mathcal{H}} (h - h')^\top A(N\lambda)^{-1/2} \zeta] + \mathbb{E}_{\zeta \sim N(0, I)}[\sup_{h \in \mathcal{H}} h^\top A(N\lambda)^{-1/2} \zeta]] \quad (24) \\ &\leq c' \sqrt{\log(2/\delta_1)} \mathbb{E}_{\zeta \sim N(0, I)}[\sup_{h \in \mathcal{H}} h^\top A(N\lambda)^{-1/2} \zeta]. \end{aligned}$$

where in line (24) we used Lemma 2. Now, we have that

$$\begin{aligned} \frac{\sup_{h, h' \in \mathcal{H}} |(h - h')^\top (\widehat{\mu}_1 - \mu)|}{n} &\leq c \sqrt{\log(2/\delta_1)} \mathbb{E}[\frac{\sup_{h \in \mathcal{H}} h^\top A(N\lambda)^{-1/2} \zeta}{n}] \\ &\leq \frac{1}{8} \end{aligned}$$

where we chose

$$\delta_1 = 2 \exp(-c'' \frac{N}{\mathbb{E}[\frac{\sup_{h \in \mathcal{H}} h^\top A(\lambda)^{-1/2} \zeta}{n}]^2})$$

for a universal constant $c'' > 0$ large enough. This proves the base case for both $h \in S_1^c$ and $h \in S_1$.

Step 1.2: Inductive Step. Next, we show the inductive step. Suppose that at round k , the hypothesis is satisfied, i.e., the algorithm has constructed estimators $\widehat{\mu}_1, \dots, \widehat{\mu}_k$ such that $\mathcal{E}_k(\widehat{\mu}_k) \cap \dots \cap \mathcal{E}_1(\widehat{\mu}_1)$ holds. Now, we construct an estimator $\widehat{\mu}_{k+1}$ for round $k + 1$. Define for every $l \in [k] \cup \{0\}$, the set

$$\widehat{S}_l = \{h : (\tilde{h}_k - h)^\top \widehat{\mu}_k \leq 2^{-l+1}n\}.$$

We will construct an estimator each subset \widehat{S}_l and then combine these into a single estimator. \widehat{S}_l can be thought of as an estimate for S_l as suggested by the following claim.

Claim 1. $S_{l+1} \subset \widehat{S}_l \subset S_{l-1}$ for all $l \in [k]$.

Proof of Claim 1. Since \mathcal{E}_k holds, by Lemma 3, we have that

1. for all $h \in S_k^c$,

$$|(\tilde{h}_k - h)^\top \widehat{\mu}_k - (h_* - h)^\top \mu| \leq \frac{1}{2} \Delta_h. \quad (25)$$

2. for all $h \in S_k$,

$$|(\tilde{h}_k - h)^\top \widehat{\mu}_k - (h_* - h)^\top \mu| \leq \frac{1}{2} 2^{-k+1}n. \quad (26)$$

Suppose $h \in \widehat{S}_l$, that is, $(\tilde{h}_k - h)^\top \widehat{\mu}_k \leq 2^{-l+1}n$. We show that $h \in S_{l-1}$. If $h \in S_k$, then we automatically have that $h \in S_{l-1}$ since $l \in [k]$ and $S_k \subset S_{l-1}$. Thus, suppose $h \in S_k^c$. Then, (25) implies that $\Delta_h \leq (\tilde{h}_k - h)^\top \widehat{\mu}_k + \frac{\Delta_h}{2} \leq 2^{-l+1}n + \frac{\Delta_h}{2}$. Rearranging, we have that $\Delta_h \leq 2^{-l+2}n$, implying that $h \in S_{l-1}$. We conclude that $\widehat{S}_l \subset S_{l-1}$.

Now, suppose that $h \in S_{l+1}$. If $h \in S_k^c$, then we have that $(\tilde{h}_k - h)^\top \widehat{\mu}_k \leq 3/2\Delta_h \leq 2^{-l+1}n$ and hence $h \in \widehat{S}_l$. Suppose $h \in S_k$. (26) implies that

$$(\tilde{h}_k - h)^\top \widehat{\mu}_k \leq \Delta_h + 2^{-k}n \leq 2^{-k+1}n + 2^{-k}n \leq 2^{-k+2}n$$

where the second inequality follows since $h \in S_k$. Thus, $h \in \widehat{S}_{k-1}$. Showing that $S_{k+1} \subset \widehat{S}_k$ follows by a similar argument and we conclude that $S_{l+1} \subset \widehat{S}_l$. This shows the claim. \square

Step 1.2.1: Constructing the estimator for \widehat{S}_l . Let $l \in [k] \cup \{0\}$. We use Theorem 1 to construct an estimator $\widehat{\mu}_{k+1,l}(\delta_{k+1})$ for each \widehat{S}_l such that for all $l \in [k] \cup \{0\}$, the event

$$\begin{aligned} \mathcal{E}_{k+1,l} = \{ \sup_{h,h' \in \widehat{S}_l} |(h-h')^\top (\widehat{\mu}_{k+1,l}(\delta_{k+1}) - \mu)| \leq c[\log(2/\delta_{k+1}) \max_{h,h' \in \widehat{S}_l} \|h-h'\|_{A(N\lambda)^{-1}} \\ + \mathbb{E}_{\zeta \sim N(0,I)}[\sup_{h \in \widehat{S}_l} h^\top A(N\lambda)^{-1/2}\zeta]] \} \end{aligned}$$

holds with probability at least $1 - \delta_{k+1}$ (that will be chosen later). We assume that $\cap_{l \in [k] \cup \{0\}} \mathcal{E}_{k+1,l}$ holds for the remainder of the proof, which by the union bound holds with probability at least $1 - \delta \log(n\epsilon^{-1})$.

Fix $l \in [k]$. We have that

$$\sup_{h,h' \in S_{l+1}} |(h-h')^\top (\widehat{\mu}_{k+1,l}(\delta_{k+1}) - \mu)| \leq \sup_{h,h' \in \widehat{S}_l} |(h-h')^\top (\widehat{\mu}_{k+1,l}(\delta_{k+1}) - \mu)| \quad (27)$$

$$\begin{aligned} &\leq c[\log(2/\delta_{k+1}) \max_{h,h' \in \widehat{S}_l} \|h-h'\|_{A(N\lambda)^{-1}} + \mathbb{E}_{\zeta \sim N(0,I)}[\sup_{h \in \widehat{S}_l} h^\top A(N\lambda)^{-1/2}\zeta] \\ &\leq c \sqrt{\log(1/\delta_{k+1}) \mathbb{E}[\sup_{h \in \widehat{S}_l} h^\top A(N\lambda)^{-1/2}\zeta]^2} \quad (28) \end{aligned}$$

$$\leq c \sqrt{\log(1/\delta_{k+1}) \mathbb{E}[\sup_{h \in S_{l-1}} h^\top A(N\lambda)^{-1/2}\zeta]^2} \quad (29)$$

where inequality (27) follows by $S_{l+1} \subset \widehat{S}_l$ from Claim 1, inequality (28) follows by Lemma 2, and the inequality (29) follows from $\widehat{S}_l \subset S_{l-1}$ from Claim 1.

Now, we have that

$$\begin{aligned} \frac{\sup_{h,h' \in S_{l+1}} |(h-h')^\top (\widehat{\mu}_{k+1,l}(\delta_{k+1}) - \mu)|}{2^{-l}n} &\leq c \sqrt{\log(1/\delta_{k+1}) \mathbb{E}[\sup_{h \in S_{l-1}} \frac{h^\top A(N\lambda)^{-1/2}\zeta}{2^{-l}n}]^2} \\ &\leq c' \sqrt{\log(1/\delta_{k+1}) \frac{\mathbb{E}[\sup_{h \in S_{l-1}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2}\zeta}{2^{-l+2}n}]^2}{N}} \\ &\leq c'' \sqrt{\log(1/\delta_{k+1}) \frac{\mathbb{E}[\sup_{h \in S_{l-1}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2}\zeta}{\Delta_h + 2^{-k+1}n}]^2}{N}} \quad (30) \end{aligned}$$

$$\begin{aligned} &\leq c'' \sqrt{\log(1/\delta_{k+1}) \frac{\mathbb{E}[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2}\zeta}{\Delta_h + 2^{-k+1}n}]^2}{N}} \\ &\leq c''' \sqrt{\log(1/\delta_{k+1}) \frac{\mathbb{E}[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2}\zeta}{(\tilde{h}_k - h)^\top \widehat{\mu}_k + 2^{-k+1}n}]^2}{N}} \quad (31) \end{aligned}$$

Since $\tilde{h}_k \in S_{l-1}$ by Lemma 3 and for all $h \in S_{l-1}$, $2^{-l+2}n \geq \Delta_h + 2^{-k+1}n$, we may apply Lemma 4 to obtain line (30). (31) follows since we assumed \mathcal{E}_k holds and Lemma 3.

Now, we choose

$$\delta_{k+1} = \log(n\epsilon^{-1}) \exp\left(-c' \frac{N}{\mathbb{E}[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2} \zeta}{(\tilde{h}_k - h)^\top \hat{\mu}_{k+2^{-k+1}n}}]^2}\right)$$

for a universal constant $c' > 0$ large enough to guarantee that with probability at least $1 - \delta_{k+1}$,

$$\frac{\sup_{h, h' \in S_{l+1}} |(h - h')^\top (\hat{\mu}_{k+1, l}(\delta_{k+1}) - \mu)|}{2^{-l}n} \leq \frac{1}{32}.$$

Step 1.2.2: Combining the estimators into a single estimator $\hat{\mu}_{k+1}$. Now, define the polyhedron based on the estimators $\hat{\mu}_{k+1, l}(\delta_{k+1})$ for $l = 0, 1, \dots, k$:

$$P = \{y \in [-1, 1]^n : \forall l \in [k] \cup \{0\}, \forall h, h' \in \hat{S}_l : \frac{|(h - h')^\top (\hat{\mu}_{k+1, l}(\delta_{k+1}) - y)|}{2^{-l}n} \leq \frac{1}{32}\}.$$

We define the estimator $\hat{\mu}_{k+1}$ as follows: if P is nonempty, then let $\hat{\mu}_{k+1}$ be any point in P ; otherwise, let $\hat{\mu}_{k+1}$ be any point in \mathbb{R}^n .

On the event $\cap_{l \in [k] \cup \{0\}} \mathcal{E}_{k+1, l}$, P is nonempty since $\mu \in P$ and, thus, $\hat{\mu}_{k+1} \in P$.

Let $l \in [k] \cup \{0\}$ and $h, h' \in \hat{S}_l$. By the triangle inequality the event $\cap_{l \in [k] \cup \{0\}} \mathcal{E}_{k+1, l}$, and $\hat{\mu}_{k+1} \in P$, we have that

$$\frac{|(h - h')^\top (\hat{\mu}_{k+1} - \mu)|}{2^{-l}n} \leq \frac{|(h - h')^\top (\hat{\mu}_{k+1} - \hat{\mu}_{k+1, l})|}{2^{-l}n} + \frac{|(h - h')^\top (\mu - \hat{\mu}_{k+1, l})|}{2^{-l}n} \leq \frac{1}{16}.$$

By the union bound, with probability at least $1 - \delta_{k+1} \log(n\epsilon^{-1})$, for every $l \in [k] \cup \{0\}$,

$$\sup_{h, h' \in S_{l+1}} \frac{|(h - h')^\top (\hat{\mu}_{k+1} - \mu)|}{2^{-l}n} \leq \sup_{h, h' \in \hat{S}_l} \frac{|(h - h')^\top (\hat{\mu}_{k+1} - \mu)|}{2^{-l}n} \leq \frac{1}{16}.$$

where we used $S_{l+1} \subset \hat{S}_l$ from Claim 1. Furthermore, since $\hat{\mu}_k \in [-1, 1]^n$ note that $\hat{S}_0 = \mathcal{H}$ and so we also have that

$$\sup_{h, h' \in \mathcal{H}} \frac{|(h - h')^\top (\hat{\mu}_{k+1} - \mu)|}{n} \leq \frac{1}{16}.$$

Thus, we have shown that for all $l \in [k] \cup \{0\}$,

$$\sup_{h, h' \in S_l} \frac{|(h - h')^\top (\hat{\mu}_{k+1} - \mu)|}{2^{-l}n} \leq \frac{1}{16}. \quad (32)$$

Now, we are ready to finish the inductive step. Let $h \in S_{k+1}^c$. Let j be the largest integer such that $h \in S_j$. Then, $2^{-j}n \leq \Delta_h \leq 2^{-j+1}n$. Then, the inequality (32) implies that

$$\frac{|(h_* - h)^\top (\hat{\mu}_{k+1} - \mu)|}{\Delta_h} \leq \frac{|(h_* - h)^\top (\hat{\mu}_{k+1} - \mu)|}{2^{-j}n} \quad (33)$$

$$\leq \sup_{h, h' \in S_j} \frac{|(h' - h)^\top (\hat{\mu}_{k+1} - \mu)|}{2^{-j}n} \quad (34)$$

$$\leq \frac{1}{16}. \quad (35)$$

A similar argument shows that if $h \in S_{k+1}$,

$$|(h_* - h)^\top (\hat{\mu}_{k+1} - \mu)| \leq \frac{2^{-k}n}{16}. \quad (36)$$

Step 2: Correctness Consider the final round $\bar{k} = \lceil \log(n\epsilon^{-1}) \rceil$ and let h such that $\Delta_h \geq \epsilon$. Then, by the previous step, we have that

$$|(h_* - h)^\top (\hat{\mu}_{\bar{k}} - \mu)| \leq \frac{\Delta_h}{8}.$$

Therefore, $(h_* - h)^\top \hat{\mu}_{\bar{k}} \geq \frac{7}{8}\Delta_h > 0$. Thus, h cannot be the empirical maximizer in the final round and the algorithm outputs $\bar{h} \in \mathcal{H}$ such that $\mu^\top \bar{h} \geq \mu^\top h_* - \epsilon$.

Step 3: Bounding the probability of error Now, we need to bound the probability of error. We bound δ_k for all k . This argument is quite similar to the one given in the proof of Theorem 5 in (Katz-Samuels et al., 2020). Fix a round k . Lemma 3 yields

$$\mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1}n + \hat{\mu}_{k+1}^\top (\tilde{h}_k - h)} \right]^2 \leq c \mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1}n + \Delta_h} \right]^2 \quad (37)$$

$$\leq c' \left[\mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}} \frac{(h_* - h)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1}n + \Delta_h} \right]^2 \right] \quad (38)$$

$$+ \mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}} \frac{(h_* - \tilde{h}_k)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1}n + \Delta_h} \right]^2 \quad (39)$$

We start by bounding the first term. Fix $h_0 \in \mathcal{H} \setminus \{h_*\}$.

$$\begin{aligned} & \mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}} \frac{(h_* - h)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1}n + \Delta_h} \right]^2 \\ & \leq \mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H} \setminus \{h_*\}} \left| \frac{(h_* - h)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1}n + \Delta_h} \right| \right]^2 \\ & \leq 8 \mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H} \setminus \{h_*\}} \frac{(h_* - h)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1}n + \Delta_h} \right]^2 + 8 \frac{\|h_* - h_0\|_{A(\lambda)^{-1}}^2}{(2^{-k+1}n + \Delta_{h_0})^2} \end{aligned} \quad (40)$$

$$\begin{aligned} & \leq 8 \left[\mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H} \setminus \{h_*\}} \frac{(h_* - h)^\top A(\lambda)^{-1/2} \zeta}{\max(\epsilon, \Delta_h)} \right]^2 \right] \\ & \quad + \max_{h \neq h_*} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{\max(\epsilon, \Delta_h)^2} \end{aligned} \quad (41)$$

where line (40) is the consequence of exercise 7.6.9 in (Vershynin, 2019).

Now, we turn to the second term. We have that

$$\begin{aligned} \mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}} \frac{(h_* - \tilde{h}_k)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1}n + \Delta_h} \right]^2 & \leq \mathbb{E}_{\zeta \sim N(0, I)} \left[\max \left(\frac{(h_* - \tilde{h}_k)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1}n}, 0 \right) \right]^2 \\ & \leq c \frac{\|h_* - \tilde{h}_k\|_{A(\lambda)^{-1}}^2}{(2^{-k+1}n)^2} \\ & \leq c' \frac{\|h_* - \tilde{h}_k\|_{A(\lambda)^{-1}}^2}{\max(\epsilon, \Delta_{\tilde{h}_k})^2} \end{aligned} \quad (42)$$

$$\leq c' \max_{h \in \mathcal{H} \setminus \{h_*\}} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{\max(\epsilon, \Delta_h)^2} \quad (43)$$

where we obtain line (42) since $\tilde{h}_k \in S_{k+2}$ by Lemma 3.

Plugging in the design used at round k , λ_k , we have that

$$\mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda_k)^{-1/2} \zeta}{2^{-k+1}n + \hat{\mu}_{k+1}^\top (\tilde{h}_k - h)} \right]^2 = \min_{\lambda} \mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1}n + \hat{\mu}_{k+1}^\top (\tilde{h}_k - h)} \right]^2 \quad (44)$$

$$\leq c \min_{\lambda} [\rho^*(\epsilon; \lambda) + \gamma^*(\epsilon; \lambda)] \quad (45)$$

$$\leq c' [\rho^*(\epsilon) + \gamma^*(\epsilon)] \quad (46)$$

where

$$\rho^*(\epsilon; \lambda) := \sup_{h \in \mathcal{H} \setminus \{h_*\}} \frac{\|h_* - h\|_{A(\lambda)^{-1}}^2}{\max(\mu^\top (h_* - h), \epsilon)^2}$$

$$\gamma^*(\epsilon; \lambda) := \mathbb{E}_{\zeta \sim N(0, I)} \left[\sup_{h \in \mathcal{H} \setminus \{h_*\}} \frac{(h_* - h)^\top A(\lambda)^{-1/2} \zeta}{\max(\mu^\top (h_* - h), \epsilon)} \right]^2.$$

and (44) follows by definition of λ_k , (45) follows by (39), (41), and (43), and (46) follows by Lemma 13 of (Katz-Samuels et al., 2020).

Putting it together, for all k

$$\begin{aligned} \delta_{k+1} &= \exp\left(-c' \frac{N}{\mathbb{E}[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2} \zeta}{(\tilde{h}_k - h)^\top \hat{\mu}_k + 2^{-k+1}n}]^2}\right) \\ &\leq \exp\left(-c'' \frac{N}{[\gamma^*(\epsilon) + \rho^*(\epsilon)]}\right) \\ &\leq \exp\left(-c''' \frac{T}{\log(n\epsilon^{-1})[\gamma^*(\epsilon) + \rho^*(\epsilon)]}\right) \end{aligned}$$

This completes the proof. □

Remark 4. We stress that the construction of the estimators in the proof of Theorem 6 is based solely on data observed by the algorithm and \mathcal{H} and at no point is knowledge of μ used.

H.1. Technical Lemmas

The proof of Theorem 6 uses the following Lemmas, which appeared originally as Lemma 11, Lemma 1, and Lemma 13 in (Katz-Samuels et al., 2020).

Lemma 2. Let $\mathcal{G} \subset \mathcal{H}$. Then,

$$\mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h, h' \in \mathcal{H}} [A(\lambda)^{-1/2} (h - h')]^\top \zeta \right]^2 \geq \frac{2}{\pi} \max_{h, h' \in \mathcal{H}} \|h - h'\|_{A(\lambda)^{-1}}^2.$$

Lemma 3. Let $k \geq 1$. Consider the k th round of Algorithm 4. Suppose that

- if $h \in S_k^c$,

$$|(h_* - h)^\top (\hat{\mu}_k - \mu)| \leq \frac{\Delta h}{8} \quad (47)$$

- if $h \in S_k$,

$$|(h_* - h)^\top (\hat{\mu}_k - \mu)| \leq \frac{2^{-k+1}n}{8}. \quad (48)$$

Then, the following hold:

1.

$$\tilde{h}_k \in S_{k+2}, \quad (49)$$

2. if $h \in S_k^c$

$$|(\tilde{h}_k - h)^\top \hat{\mu}_k - (h_* - h)^\top \mu| \leq \frac{1}{2} \Delta_h. \quad (50)$$

3. if $h \in S_k$,

$$|(\tilde{h}_k - h)^\top \hat{\mu}_k - (h_* - h)^\top \mu| \leq \frac{1}{2} 2^{-k+1} n. \quad (51)$$

4. There exist universal constants $c, c' > 0$ such that

$$\begin{aligned} c \mathbb{E} \left[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2} \zeta}{\Delta_h + 2^{-k+1} n} \right]^2 &\leq \mathbb{E} \left[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2} \zeta}{(\tilde{h}_k - h)^\top \hat{\mu}_k + 2^{-k+1} n} \right]^2 \\ &\leq c' \mathbb{E} \left[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2} \zeta}{\Delta_h + 2^{-k+1} n} \right]^2 \end{aligned}$$

Lemma 4. Let $V = \{v_1, \dots, v_l\} \subset \mathbb{R}^d$ with $0 \in V$. Suppose $a_i \geq 1$ for all $i \in [l]$. Then,

$$\mathbb{E}_{\zeta \sim N(0, I)} \left[\sup_{v_i \in V} v_i^\top \zeta \right] \leq \mathbb{E}_{\zeta \sim N(0, I)} \left[\sup_{v_i \in V} a_i v_i^\top \zeta \right]$$

I. Efficient Fixed Budget Algorithm

Algorithm 5 Fixed Budget ACED for Combinatorial Bandits (Computationally Efficient).

Input: Budget T , tolerance $\epsilon > 0$.

$$\hat{\mu}_1 = \mathbf{0} \in \mathbb{R}^d, N \leftarrow \left\lceil \frac{T}{\log_2(n\epsilon^{-1})} \right\rceil.$$

for $k = 1, 2, \dots, \left\lceil \log_2(d\epsilon^{-1}) \right\rceil$ **do**

$$\tilde{h}_k \leftarrow \arg \max_{h \in \mathcal{H}} \hat{\mu}_k^\top h$$

Let $\lambda_k^{(1)}$ be the solution of the following optimization problem

$$\inf_{\lambda \in \Delta} \mathbb{E}_{\zeta \sim N(0, I)} \left[\max_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2} \zeta}{2^{-k+1} n + \hat{\mu}_k^\top (\tilde{h}_k - h)} \right]^2 \quad (52)$$

Let $\lambda_k^{(2)}$ be the solution to

$$\inf_{\lambda \in \Delta} \max_{h \in \mathcal{H}} \max_{i \in \tilde{h}_k \Delta_h} \frac{1}{\lambda_i} \frac{1}{2^{-k+1} n + \hat{\mu}_k^\top (\tilde{h}_k - h)} \quad (53)$$

$$\lambda_k \leftarrow \frac{1}{2} (\lambda_k^{(1)} + \lambda_k^{(2)})$$

Sample $\{x_1, \dots, x_N\} \sim \lambda_k$.

Query x_1, \dots, x_N and receive rewards y_1, \dots, y_N .

$$\text{Let } \hat{\mu}_{k+1} = \frac{1}{N} \sum_{s=1}^N A(\lambda)^{-1} x_s y_s.$$

end for

Return $\arg \max_{h \in \mathcal{H}} \hat{\mu}_{k+1}^\top h$

In this section, we present Algorithm 5, which can be implemented in practice. Algorithm 5 resembles Algorithm 4 with two differences. First, it uses the IPS estimator instead of the theoretical estimator derived in the proof of Theorem 6. Second, it mixes the design (given in used in (23)) with the design defined in (53) in order to control the worst-case deviations of the IPS estimator.

Now, we briefly discuss why Algorithm 5 can be efficiently implemented. As is standard in combinatorial bandits, we assume access to a linear maximization oracle: for any $v \in \mathbb{R}^d$, we can compute

$$\text{Oracle}(v) = \arg \max_{z \in \mathcal{Z}} v^\top z$$

efficiently.² (Katz-Samuels et al., 2020) gave a procedure for efficiently solving (52) up to a constant factor assuming a linear maximization oracle, and so we focus on (53). Define the function $g(\lambda) = \max_{i \in \tilde{h}_k \Delta h} \frac{\frac{1}{\lambda_i}}{2^{-k+1}n + \widehat{\mu}_k^\top (\tilde{h}_k - h)}$. g is convex since it is the maximum of a collection of convex functions. One can compute the subgradient of g using the linear maximization oracle. For each $i \in \tilde{h}_k$, one can compute

$$\max_{h: i \notin \tilde{h}_k} \widehat{\mu}_k^\top h$$

by defining

$$\tilde{\mu}_j^{(i)} = \begin{cases} \tilde{\mu}_j & j \neq i \\ -\infty & j = i \end{cases}$$

and use the linear maximization oracle to compute $\max_h h^\top \tilde{\mu}^{(i)}$. Using a similar technique, for each $i \notin \tilde{h}_k$, one can compute

$$\max_{h: i \in \tilde{h}_k} \widehat{\mu}_k^\top h.$$

Then, using a standard result about subgradients, we have that $\partial g(\lambda)$ consists of $\frac{\partial}{\partial \lambda} \frac{\frac{1}{\lambda_i}}{\widehat{\mu}_k^\top (\tilde{h}_k - h)}$ where \tilde{i} and h attain the maximum in

$$\max_{h \in \mathcal{H}} \max_{i \in \tilde{h}_k \Delta h} \frac{\frac{1}{\lambda_i}}{2^{-k+1}n + \widehat{\mu}_k^\top (\tilde{h}_k - h)}.$$

Thus, using this trick for computing the subgradient and any optimization procedure for nonsmooth convex optimization, we can optimize (53). In Theorem 7, for the sake of simplicity and focusing on the key ideas, we suppose that (52) and (53) are solved exactly, but using the optimization ideas laid out here would only affect the sample complexity up to constant factors.

The following complexity parameter is a key term in Theorem 7:

$$\psi^*(\epsilon) := \min_{\lambda \in \Delta} \max_{h \in \mathcal{H} \setminus \{h_*\}} \max_{i \in h_* \Delta h} \frac{\frac{1}{\lambda_i}}{\max(\epsilon, \mu^\top (h_* - h))}.$$

Theorem 7. *Let $T \in \mathbb{N}$ and $\epsilon > 0$. If $T \geq c_1 \log(n\epsilon^{-1})[\gamma^*(\epsilon) + \rho^*(\epsilon) + \psi^*(\epsilon) \log(|\mathcal{H}|) \log(\log(|\mathcal{H}|))]$, then Algorithm 5 satisfies*

$$\mathbb{P}(\mu^\top \widehat{h} + \epsilon < \mu^\top h_*) \leq c_2 \log(n\epsilon^{-1})^2 \exp\left(-\frac{T - \log(|\mathcal{H}|) \log(\log(|\mathcal{H}|)) \psi^*(\epsilon)}{\log(n\epsilon^{-1})[\gamma^*(\epsilon) + \rho^*(\epsilon) + \psi^*(\epsilon)]}\right)$$

where $c_1, c_2 > 0$ are universal constants.

Proof of Theorem 7. Let $\alpha > 0$ be a universal constant that will be specified later. Let $\delta > 0$ such that

$$T = \alpha \log(n\epsilon^{-1})[\log(1/\delta)(\gamma^*(\epsilon) + \rho^*(\epsilon) + \psi^*(\epsilon)) + \psi^*(\epsilon) \log(|\mathcal{H}|) \log(\log(|\mathcal{H}|))],$$

which exists since by assumption $T \geq c \log(n\epsilon^{-1})[\gamma^*(\epsilon) + \rho^*(\epsilon) + \psi^*(\epsilon) \log(|\mathcal{H}|) \log(\log(|\mathcal{H}|))]$. Note that

$$N = \log(1/\delta) \alpha (\gamma^*(\epsilon) + \rho^*(\epsilon) + \psi^*(\epsilon)) + \psi^*(\epsilon) \log(|\mathcal{H}|) \log(\log(|\mathcal{H}|)). \quad (54)$$

²In the setting of active classification, this is equivalent to assuming a weighted 0/1 loss minimization oracle.

Recall $\widehat{\mu}_k = \frac{1}{N} A(\lambda)^{-1} \sum_{s=1}^N x_{I_s} y_s$ and

$$S_k = \{h \in \mathcal{H} : \Delta_h \leq n2^{-k+1}\}.$$

Note that $S_0 = \mathcal{H}$ since $\mu \in [-1, 1]$ by assumption. Define the event at round k for $l \in [k] \cup \{0\}$,

$$\begin{aligned} \mathcal{E}_{k,l} = & \left\{ \sup_{h, h' \in S_l} |(h - h')^\top (\widehat{\mu}_k - \mu)| \leq \right. \\ & c[\mathbb{E}[\sup_{h \in S_l} h^\top A(N\lambda)^{-1/2} \zeta] + (\log(\log(|S_l|)) \log(|S_l|) + \log(1/\delta)) \frac{\max_{h, h' \in S_l} \max_{i \in h \Delta h'} \frac{1}{\lambda_i}}{N}] \\ & \left. + \left\| h - h' \right\|_{A(N\lambda)^{-1}} \sqrt{\log(1/\delta)} \right\} \end{aligned}$$

Define $\mathcal{E}_k = \cap_{l=0}^k \mathcal{E}_{k,l}$ and $\mathcal{E} = \cap_{k=1}^{\log_2(n\epsilon^{-1})} \mathcal{E}_k$. Now, using the law of total probability, the independence of samples between rounds, and Lemma 5

$$\mathbb{P}(\mathcal{E}^c) \leq \sum_{k \geq 1} \sum_{l=0}^k \mathbb{P}(\mathcal{E}_{k,l}^c | \cap_{j=1}^{k-1} \mathcal{E}_j) \leq \log(n\epsilon^{-1})^2 \delta.$$

Suppose that \mathcal{E} occurs for the remainder of the proof. Note that using the same series of inequalities as in (27)-(29), we have that

$$\begin{aligned} & \sup_{h, h' \in S_l} |(h - h')^\top (\widehat{\mu}_k - \mu)| \leq \\ & c'[\sqrt{\log(1/\delta)} \mathbb{E}[\sup_{h \in S_l} h^\top A(N\lambda)^{-1/2} \zeta] + (\log(\log(|S_l|)) \log(|S_l|) + \log(1/\delta)) \frac{\max_{h, h' \in S_l} \max_{i \in h \Delta h'} \frac{1}{\lambda_i}}{N}] \end{aligned}$$

We argue inductively that at every round $k \geq 2$

1. for all $h \in S_k^c$,

$$|(h_* - h)^\top (\widehat{\mu}_k - \mu)| \leq \frac{\Delta_h}{8}$$

2. for all $h \in S_k$,

$$|(h_* - h)^\top (\widehat{\mu}_k - \mu)| \leq \frac{2^{-k+1}n}{8}.$$

Claim: Base Case. Let $k = 2$. Then, the event $\mathcal{E}_{2,0}$ implies that

$$\begin{aligned} & \sup_{h \in \mathcal{H}} \frac{|(h_* - h)^\top (\widehat{\mu}_k - \mu)|}{2^{-1}n} \\ & \leq \sup_{h, h' \in \mathcal{H}} \frac{|(h - h')^\top (\widehat{\mu}_k - \mu)|}{2^{-1}n} \\ & \leq \frac{1}{2^{-1}n} c'[\sqrt{\log(1/\delta)} \mathbb{E}[\sup_{h \in \mathcal{H}} h^\top A(N\lambda)^{-1/2} \zeta] + (\log(|\mathcal{H}|) + \log(1/\delta)) \frac{\max_{h, h' \in \mathcal{H}} \max_{i \in h \Delta h'} \frac{1}{\lambda_i}}{N}] \\ & = \frac{1}{2^{-1}n} c' \left[\sqrt{\frac{\log(1/\delta) \mathbb{E}[\sup_{h \in \mathcal{H}} h^\top A(\lambda)^{-1/2} \zeta]^2}{\log(1/\delta) \alpha[\gamma^*(\epsilon) + \rho^*(\epsilon)]}} + (\log(|\mathcal{H}|) + \log(1/\delta)) \frac{\max_{h, h' \in \mathcal{H}} \max_{i \in h \Delta h'} \frac{1}{\lambda_i}}{\alpha \psi^*(\epsilon) (\log(\log(\mathcal{H})) \log(\mathcal{H}) + \log(1/\delta))} \right]. \end{aligned} \tag{55}$$

where the last line follows by (54). Recall that $\lambda_1 = \frac{1}{2}(\lambda_1^{(1)} + \lambda_2^{(1)})$ where $\lambda_1^{(1)}$ is the minimizer of (52) and $\lambda_1^{(1)}$ is the minimizer of (53). We have that

$$\mathbb{E}[\sup_{h \in \mathcal{H}} \frac{h^\top A(\lambda_1)^{-1/2} \zeta}{2^{-1}n}]^2 \leq 4\mathbb{E}[\sup_{h \in \mathcal{H}} \frac{h^\top A(\lambda_1^{(1)})^{-1/2} \zeta}{2^{-1}n}]^2 \quad (56)$$

$$= 4 \min_{\lambda} \mathbb{E}[\sup_{h \in \mathcal{H}} \frac{h^\top A(\lambda)^{-1/2} \zeta}{2^{-1}n}]^2 \quad (57)$$

$$\leq c[\gamma^*(\epsilon) + \rho^*(\epsilon)] \quad (58)$$

(56) follows by sudakov-fernieque inequality since $A(\lambda_1) \succeq \frac{1}{2}A(\lambda_1^{(1)})$ implies that $\sqrt{2}A(\lambda_1^{(1)})^{-1/2} \succeq A(\lambda_1)^{-1/2}$. Line (57) follows by definition of $\lambda_1^{(1)}$. Finally, (58) follows by the same series of inequalities that established (46).

Fix $h, h' \in \mathcal{H}$. Then, if $i \in h\Delta h'$, then either $i \in h\Delta h_*$ or $i \in h'\Delta h_*$. Thus,

$$\begin{aligned} \max_{h, h' \in \mathcal{H}} \max_{i \in h\Delta h'} \frac{1}{\lambda_{1,i}} &\leq 2 \max_{h \in \mathcal{H} \setminus \{h_*\}} \max_{i \in h\Delta h_*} \frac{1}{\lambda_{1,i}} \\ &\leq 4 \max_{h \in \mathcal{H} \setminus \{h_*\}} \max_{i \in h\Delta h_*} \frac{1}{\lambda_{1,i}^{(2)}} \end{aligned} \quad (59)$$

$$= 8 \min_{\lambda} \max_{h \in \mathcal{H} \setminus \{h_*\}} \max_{i \in h\Delta h_*} \frac{1}{\lambda_i} \quad (60)$$

$$\leq c\psi^*(\epsilon) \quad (61)$$

where (59) follows by definition of λ_1 , (60) follows by definition of $\lambda_1^{(2)}$, and (61) follows by definition of $\psi^*(\epsilon)$ and since for all $h \in \mathcal{H} \setminus \{h_*\}$, $\Delta_h \leq 2n$.

Then, putting together (55), (58), and (61), we have that if the universal constant α is large enough, then

$$\sup_{h, h' \in \mathcal{H}} \frac{|(h - h')^\top (\hat{\mu}_k - \mu)|}{2^{-1}n} \leq \frac{1}{8}$$

this proves the base case.

Claim: Inductive Step. Suppose that at round $k \geq 2$

1. for all $h \in S_k^c$,

$$|(h_* - h)^\top (\hat{\mu}_k - \mu)| \leq \frac{\Delta_h}{8}$$

2. for all $h \in S_k$,

$$|(h_* - h)^\top (\hat{\mu}_k - \mu)| \leq \frac{2^{-k+1}n}{8}.$$

Now, we prove the statement for round $k + 1$. Let $l \in [k + 1]$. Using $\mathcal{E}_{k+1,l}$, we have that

$$\frac{\sup_{h, h' \in S_l} |(h - h')^\top (\hat{\mu}_{k+1} - \mu)|}{2^{-l}n} \leq c \sqrt{\log(1/\delta) \mathbb{E}[\sup_{h \in S_l} \frac{h^\top A(N\lambda_k)^{-1/2} \zeta}{2^{-l}n}]^2} \quad (62)$$

$$+(\log(\log(|S_l|)) \log(|S_l|) + \log(1/\delta)) \frac{\max_{h, h' \in S_l} \max_{i \in h\Delta h'} \frac{1/\lambda_{k,i}}{2^{-l}n}}{N} \quad (63)$$

Bounding the first term, we have that

$$\begin{aligned} \sqrt{\log(1/\delta)\mathbb{E}\left[\sup_{h \in S_l} \frac{h^\top A(N\lambda_k)^{-1/2}\zeta}{2^{-l}n}\right]^2} &\leq c' \sqrt{\log(1/\delta) \frac{\mathbb{E}\left[\sup_{h \in S_l} \frac{(\tilde{h}_k - h)^\top A(\lambda_k)^{-1/2}\zeta}{2^{-l+2}n}\right]^2}{N}} \\ &\leq c'' \sqrt{\log(1/\delta) \frac{\mathbb{E}\left[\sup_{h \in S_l} \frac{(\tilde{h}_k - h)^\top A(\lambda_k)^{-1/2}\zeta}{\Delta_h + 2^{-k+1}n}\right]^2}{N}} \end{aligned} \quad (64)$$

$$\leq c'' \sqrt{\log(1/\delta) \frac{\mathbb{E}\left[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda_k)^{-1/2}\zeta}{\Delta_h + 2^{-k+1}n}\right]^2}{N}} \quad (65)$$

$$\leq c''' \sqrt{\log(1/\delta) \frac{\mathbb{E}\left[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda_k)^{-1/2}\zeta}{(\tilde{h}_k - h)^\top \hat{\mu}_k + 2^{-k+1}n}\right]^2}{N}} \quad (66)$$

Since $\tilde{h}_k \in S_l$ by Lemma 3 and for all $h \in S_{l-1}$, $2^{-l+2}n \geq \Delta_h + 2^{-k+1}n$, we may apply Lemma 4 to obtain line (64). (66) follows since we assumed \mathcal{E}_k holds and Lemma 3. Using a similar series of inequalities to Step 3 in the proof of Theorem 3, we have that

$$\begin{aligned} \sqrt{\log(1/\delta)\mathbb{E}\left[\sup_{h \in S_l} \frac{h^\top A(N\lambda_k)^{-1/2}\zeta}{2^{-l}n}\right]^2} &\leq c''' \sqrt{\log(1/\delta) \frac{\mathbb{E}\left[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda_k)^{-1/2}\zeta}{(\tilde{h}_k - h)^\top \hat{\mu}_k + 2^{-k+1}n}\right]^2}{N}} \\ &\leq c'''' \sqrt{\log(1/\delta) \frac{\mathbb{E}\left[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda_k^{(1)})^{-1/2}\zeta}{(\tilde{h}_k - h)^\top \hat{\mu}_k + 2^{-k+1}n}\right]^2}{N}} \end{aligned} \quad (67)$$

$$= c'''' \min_{\lambda} \sqrt{\log(1/\delta) \frac{\mathbb{E}\left[\sup_{h \in \mathcal{H}} \frac{(\tilde{h}_k - h)^\top A(\lambda)^{-1/2}\zeta}{(\tilde{h}_k - h)^\top \hat{\mu}_k + 2^{-k+1}n}\right]^2}{N}} \quad (68)$$

$$\leq c'''' \sqrt{\log(1/\delta) \frac{\rho^*(\epsilon) + \gamma^*(\epsilon)}{N}} \quad (69)$$

$$\leq \frac{1}{32} \quad (70)$$

where (67) follows by Sudakov-Fernique inequality and the definition of λ_k , (68) follows by the definition of $\lambda_k^{(1)}$, (69) follows by the same series of inequalities used to establish (46), and the last line follows by plugging in (54) and letting α be a large enough universal constant.

Now, we bound the second term. We have that

$$\begin{aligned} \max_{h, h' \in S_l} \max_{i \in h \Delta h'} \frac{1/\lambda_{k,i}}{2^{-l}n} &\leq 2 \max_{h \in S_l} \max_{i \in \tilde{h}_k \Delta h} \frac{1/\lambda_{k,i}}{2^{-l}n} \\ &\leq 4 \max_{h \in S_l} \max_{i \in \tilde{h}_k \Delta h} \frac{1/\lambda_{k,i}}{2^{-k+1}n + \Delta_h} \end{aligned} \quad (71)$$

$$\leq c \max_{h \in S_l} \max_{i \in \tilde{h}_k \Delta h} \frac{1/\lambda_{k,i}}{2^{-k+1}n + \hat{\mu}_k^\top(\tilde{h}_k - h)} \quad (72)$$

$$\leq c \max_{h \in \mathcal{H}} \max_{i \in \tilde{h}_k \Delta h} \frac{1/\lambda_{k,i}}{2^{-k+1}n + \hat{\mu}_k^\top(\tilde{h}_k - h)}$$

$$\begin{aligned} &\leq c' \max_{h \in \mathcal{H}} \max_{i \in \tilde{h}_k \Delta h} \frac{1/\lambda_{k,i}^{(2)}}{2^{-k+1}n + \hat{\mu}_k^\top(\tilde{h}_k - h)} \\ &= c' \min_{\lambda} \max_{h \in \mathcal{H}} \max_{i \in \tilde{h}_k \Delta h} \frac{1/\lambda_i}{2^{-k+1}n + \hat{\mu}_k^\top(\tilde{h}_k - h)} \end{aligned} \quad (73)$$

$$\leq c'' \min_{\lambda} \max_{h \in \mathcal{H}} \max_{i \in \tilde{h}_k \Delta h} \frac{1/\lambda_i}{2^{-k+1}n + \Delta_h} \quad (74)$$

$$\leq c'' \min_{\lambda} \max_{i \in \tilde{h}_k \Delta h_*} \frac{1/\lambda_i}{2^{-k+1}n + \Delta_{\tilde{h}_k}} + \max_{h \in \mathcal{H}} \max_{i \in h_* \Delta h} \frac{1/\lambda_i}{2^{-k+1}n + \Delta_h} \quad (75)$$

$$\begin{aligned} &\leq c''' \min_{\lambda} \max_{h \in \mathcal{H}} \max_{i \in h_* \Delta h} \frac{1/\lambda_i}{2^{-k+1}n + \Delta_h} \\ &\leq c''' \psi^*(\epsilon) \end{aligned} \quad (76)$$

where in (71) we used that since for all $h \in S_l$, $2^{-l}n \geq \frac{2^{-k+1}n + \Delta_h}{2}$ and in (72) we used the inductive hypothesis and Lemma 3 and in (73), we used the definition of $\lambda_k^{(2)}$, and in (74) we used the inductive hypothesis and Lemma 3, and finally in (75) we used if $i \in \tilde{h}_k \Delta h$, then either $i \in h_* \Delta h$ or $i \in \tilde{h}_k \Delta h_*$.

Thus,

$$(\log(\log(|S_l|)) \log(|S_l|) + \log(1/\delta)) \frac{\max_{h, h' \in S_l} \max_{i \in h \Delta h'} \frac{1/\lambda_{k,i}}{2^{-l}n}}{N} \leq c''' (\log(\log(|S_l|)) \log(|S_l|) + \log(1/\delta)) \frac{\psi^*(\epsilon)}{N} \quad (77)$$

$$\leq \frac{1}{32} \quad (78)$$

where the last line follows by plugging in (54) and letting α be a large enough universal constant. Then, putting together (63), (70), and (78) yields for every $l \in [k+1] \cup \{0\}$

$$\frac{\sup_{h, h' \in S_l} |(h - h')^\top (\hat{\mu}_{k+1} - \mu)|}{2^{-l}n} \leq \frac{1}{16}$$

Using the same series of inequalities used to establish (35) and (36) completes the inductive step.

Correctness. This argument is the same as in the correctness step in the proof of Theorem 6 □

I.1. Lemmas

In this Section, we prove the main concentration inequality, Lemma 5, that is used in the proof of Theorem 7. We need the following definition for the proof of Lemma 5.

Definition 2. Fix a set $T \subset \mathbb{R}^n$. Let d be a metric on \mathbb{R}^n . The γ_1 -functional of T is defined as

$$\gamma_1(T, d) = \inf_{(T_k)} \sup_{t \in T} \sum_{k=0}^{\infty} 2^k d(t, T_k)$$

where the infimum is taken over all admissible sequences of T .

Lemma 5. *Let $\mathcal{G} \subset \mathcal{H}$. Fix some $\lambda \in \Delta_n$ and draw $\{x_s\}_{i=1}^t \sim \lambda$ and then observe y_s with mean $x_s^\top \mu$ and $|y_i| \leq 1$ with probability 1 for $s = 1, \dots, t$. Then, there exists a universal constant $c > 0$ such that for any $u > 0$ with probability at most $\exp(-u)$*

$$\begin{aligned} \sup_{h, h' \in \mathcal{G}} |(h - h')^\top (\frac{1}{t} A(\lambda)^{-1} \sum_{s=1}^t x_{I_s} y_s - \mu)| &\leq \\ c[\mathbb{E}[\sup_{h \in \mathcal{G}} h^\top A(t\lambda)^{-1/2} \zeta] + (\log(\log(|\mathcal{G}|)) \log(|\mathcal{G}|) + u) \frac{\max_{h, h' \in \mathcal{G}} \max_{i \in h \Delta h'} \frac{1}{\lambda_i}}{t} & \\ + \left\| h - h' \right\|_{A(t\lambda)^{-1}} \sqrt{u}] & \end{aligned}$$

Proof. Using Corollary 7.9 of (Ledoux, 2001), we have that with probability at least $1 - \exp(-u)$,

$$\begin{aligned} \sup_{h, h' \in \mathcal{G}} |(h - h')^\top (\frac{1}{t} A(\lambda)^{-1} \sum_{s=1}^t x_{I_s} y_s - \mu)| &\leq \\ c[\mathbb{E}[\sup_{h, h' \in \mathcal{G}} |(h - h')^\top (\frac{1}{t} A(\lambda)^{-1} \sum_{s=1}^t x_{I_s} y_s - \mu)|] + u \frac{\max_{h, h' \in \mathcal{G}} \max_{i \in h \Delta h'} \frac{1}{\lambda_i}}{t} & \\ + \left\| h - h' \right\|_{A(t\lambda)^{-1}} \sqrt{u}. & \end{aligned} \quad (79)$$

Using the standard technique of symmetrization, we have that

$$\mathbb{E}[\sup_{h, h' \in \mathcal{G}} |(h - h')^\top (\frac{1}{t} A(\lambda)^{-1} \sum_{s=1}^t x_{I_s} y_s - \mu)|] \leq 2\mathbb{E}_{\epsilon_1, \dots, \epsilon_t}[\sup_{h, h' \in \mathcal{G}} |(h - h')^\top (\frac{1}{t} A(\lambda)^{-1} \sum_{s=1}^t x_{I_s} \epsilon_s)|] \quad (80)$$

where $\epsilon_1, \dots, \epsilon_t$ are Rademacher random variables. By Bernstein's inequality, we have that

$$\mathbb{P}(|(h - h')^\top (\frac{1}{t} A(\lambda)^{-1} \sum_{s=1}^t x_{I_s} \epsilon_s| \geq u) \leq 2 \exp(-c \frac{u^2}{\left\| h - h' \right\|_{A(t\lambda)^{-1}}^2}, \frac{u}{\max_{i \in h \Delta h'} \frac{1}{\lambda_i}}). \quad (81)$$

Define the set $\tilde{\mathcal{G}} = \{A(\lambda)^{-1} h : h \in \mathcal{G}\}$. Then, using (81) and applying Theorem 2.2.23 of (Talagrand, 2014), we have that

$$\mathbb{E}[\sup_{h, h' \in \mathcal{G}} |(h - h')^\top (\frac{1}{t} A(\lambda)^{-1} \sum_{s=1}^t x_{I_s} \epsilon_s)|] \leq c[\gamma_2(\tilde{\mathcal{G}}, d_2) + \gamma_1(\tilde{\mathcal{G}}, d_1)] \quad (82)$$

where d_2 is the metric induced by $\left\| \cdot \right\|_2$ and d_1 is the metric induced $\left\| \cdot \right\|_\infty$. Using Talagrand's majorizing measure theorem (Theorem 8.6.1 in (Vershynin, 2019)), we have that

$$\gamma_2(\tilde{\mathcal{G}}, d_2) \leq c\mathbb{E}_{\zeta \sim N(0, I)}[\sup_{h \in \mathcal{H}} h^\top A(\lambda)^{-1} \zeta]. \quad (83)$$

From the definition of γ_1 , it follows trivially that

$$\gamma_1(\tilde{\mathcal{G}}, d_1) \leq \log(\log(|\mathcal{G}|)) \log(|\mathcal{G}|) \frac{\max_{h, h' \in \mathcal{G}} \max_{i \in h \Delta h'} \frac{1}{\lambda_i}}{t}. \quad (84)$$

Finally, putting together (79), (80), (82), (83), and (84), we obtain the result. \square

J. Thresholds

Let n be a power of 2. Let $\mathcal{H} = \{e_{[k]} : k \in [n]\}$ where $[e_A]_i = \mathbf{1}\{i \in A\}$. Assume that $\mu_* \in \epsilon(2e_{[k_*]} - 1)$ for some $k_* \in [n]$. Then

$$\begin{aligned} \gamma^* &:= \inf_{\lambda} \mathbb{E} \left[\sup_{h \in \mathcal{H}} \frac{\langle h - h_*, \widehat{\mu} - \mu_* \rangle}{\langle h - h_*, \mu_* \rangle} \right]^2 \\ &= \inf_{\lambda} \mathbb{E} \left[\sup_{h \in \mathcal{H}} \frac{\langle h - h_*, A(\lambda)^{-1/2} \zeta \rangle}{\langle h - h_*, \mu_* \rangle} \right]^2 \\ &= \inf_{\lambda} \mathbb{E} \left[\sup_{k > k_*} \frac{\sum_{i=k_*+1}^k \frac{\zeta_i}{\sqrt{\lambda_i}}}{(k - k_*)h} \vee \sup_{k < k_*} \frac{\sum_{i=k+1}^{k_*} \frac{\zeta_i}{\sqrt{\lambda_i}}}{(k_* - k)h} \right]^2 \\ &\leq \inf_{\lambda} \frac{4}{\epsilon^2} \mathbb{E} \left[\max_{k=1, \dots, n} \frac{1}{k} \sum_{i=1}^k \frac{\zeta_i}{\sqrt{\lambda_i}} \right]^2. \end{aligned}$$

and

$$\begin{aligned} \max_{k=1, \dots, n} \frac{1}{k} \sum_{i=1}^k \frac{\zeta_i}{\sqrt{\lambda_i}} &= \max_{j=0, \dots, \log_2(n)} \max_{k=[2^j, 2^{j+1})} \frac{1}{k} \sum_{i=1}^k \frac{\zeta_i}{\sqrt{\lambda_i}} \\ &\leq \max_{j=0, \dots, \log_2(n)} \max_{k=[2^j, 2^{j+1})} \frac{1}{2^j} \sum_{i=1}^k \frac{\zeta_i}{\sqrt{\lambda_i}} \\ &= \max_{j=0, \dots, \log_2(n)} \left(\frac{1}{2^j} \sum_{i=1}^{2^j-1} \frac{\zeta_i}{\sqrt{\lambda_i}} + \max_{k=[2^j, 2^{j+1})} \frac{1}{2^j} \sum_{i=2^j}^k \frac{\zeta_i}{\sqrt{\lambda_i}} \right) \\ &\leq \max_{j=0, \dots, \log_2(n)} \left(\frac{1}{2^j} \sum_{i=1}^{2^j-1} \frac{\zeta_i}{\sqrt{\lambda_i}} \right) + \max_{j=0, \dots, \log_2(n)} \left(\max_{k=[2^j, 2^{j+1})} \frac{1}{2^j} \sum_{i=2^j}^k \frac{\zeta_i}{\sqrt{\lambda_i}} \right). \end{aligned}$$

We now take $\lambda_i = \frac{1}{i \log_2(n)}$. Note that $\sum_{i=1}^n \lambda_i \geq 1/2$ and for any $k > 1$ we have $\sum_{i=1}^{k-1} \frac{1}{\lambda_i} \leq k^2 \log_2(n)/2$.

For any ν we have

$$\mathbb{E} \left[\exp \left(\nu \frac{1}{2^j} \sum_{i=1}^{2^j-1} \frac{\zeta_i}{\sqrt{\lambda_i}} \right) \right] = \exp \left(\nu^2 \frac{1}{2^{2j}} \sum_{i=1}^{2^j-1} \frac{1}{\lambda_i} \right) \leq \exp(\nu^2 \log_2(n)/2)$$

so we apply Proposition 1 with $\sigma_t^2 = \log_2(n)$ and $\mathcal{T} = \{0, 1, \dots, \log_2(n)\}$ to obtain

$$\mathbb{E} \left[\max_{j=0, \dots, \log_2(n)} \left(\frac{1}{2^j} \sum_{i=1}^{2^j-1} \frac{\zeta_i}{\sqrt{\lambda_i}} \right) \right] \leq \sqrt{2 \log_2(n) \log(\log_2(2n))}.$$

The second term deserves a bit more care. First, note that

$$\begin{aligned} \max_{j=0, \dots, \log_2(n)} \left(\max_{k=[2^j, 2^{j+1})} \frac{1}{2^j} \sum_{i=2^j}^k \frac{\zeta_i}{\sqrt{\lambda_i}} \right) &\leq \max_{j=0, \dots, \log_2(n)} \left(\max_{k=[2^j, 2^{j+1})} \sqrt{\frac{2 \log_2(n)}{2^j}} \sum_{i=2^j}^k \zeta_i \right) \\ &= \max_{j=0, \dots, \log_2(n)} \left(\max_{k=1, \dots, 2^j} \sqrt{\frac{2 \log_2(n)}{2^j}} \sum_{i=1}^k \zeta_i^{(j)} \right) \end{aligned}$$

where each $\{\zeta_i^{(j)}\}_{i=1}^{2^j}$ are i.i.d. sequences of $N(0, 1)$. Now

$$\begin{aligned} & \mathbb{P} \left(\max_{j=0, \dots, \log_2(n)} \left(\max_{k=1, \dots, 2^j} \sqrt{\frac{2 \log_2(n)}{2^j}} \sum_{i=1}^k \zeta_i^{(j)} \right) > t \right) \\ & \leq \sum_{j=0, \dots, \log_2(n)} \mathbb{P} \left(\max_{k=1, \dots, 2^j} \sqrt{\frac{2 \log_2(n)}{2^j}} \sum_{i=1}^k \zeta_i^{(j)} > t \right) \\ & \leq \sum_{j=0, \dots, \log_2(n)} \exp(-t^2/4 \log_2(n)) \\ & = \log_2(2n) \exp(-t^2/4 \log_2(n)) \end{aligned}$$

where the last inequality follows from Doob's maximal inequality. Thus,

$$\begin{aligned} & \mathbb{E} \left[\max_{j=0, \dots, \log_2(n)} \left(\max_{k=1, \dots, 2^j} \sqrt{\frac{2 \log_2(n)}{2^j}} \sum_{i=1}^k \zeta_i^{(j)} \right) \right] \\ & \leq \int_{t=0}^{\infty} \mathbb{P} \left(\max_{j=0, \dots, \log_2(n)} \left(\max_{k=1, \dots, 2^j} \sqrt{\frac{2 \log_2(n)}{2^j}} \sum_{i=1}^k \zeta_i^{(j)} \right) > t \right) dt \\ & \leq \int_{t=0}^{\infty} \min\{1, \log_2(2n) \exp(-t^2/4 \log_2(n))\} dt \\ & \leq a + \log_2(2n) \sqrt{4\pi \log_2(n)} \int_{t=a}^{\infty} \frac{1}{\sqrt{4\pi \log_2(n)}} \exp(-t^2/4 \log_2(n)) dt \\ & \leq a + \log_2(2n) \sqrt{4\pi \log_2(n)} \exp(-a^2/4 \log_2(n)) \\ & \leq \sqrt{4 \log_2(n) \log(\log_2(2n))} + \sqrt{4\pi \log_2(n)} \end{aligned}$$

for $a = \sqrt{4 \log_2(n) \log(\log_2(2n))}$.

Putting it all together we have

$$\begin{aligned} \gamma^* & \leq \frac{4}{\epsilon^2} \left(\sqrt{2 \log_2(n) \log(\log_2(2n))} + \sqrt{4 \log_2(n) \log(\log_2(2n))} + \sqrt{4\pi \log_2(n)} \right)^2 \\ & \leq \frac{194 \log_2(n) \log(\log_2(2n))}{\epsilon^2} \end{aligned}$$

K. Implementation Details

In this section we discuss the modifications and implementation details of Algorithm 2. We consider the slightly modified version present in Algorithm 9.

K.1. Optimization

First, we reiterate the following definition from Section 4.2:

$$f(\lambda; h; \zeta) := \frac{\sum_{i \in [n]} (\tilde{h}_k(x_i) - h(x_i)) \frac{\zeta_i}{n \lambda_i^{1/2}}}{2^{-k+1} + \text{err}(h, \hat{\eta}_{k-1}) - \text{err}(\tilde{h}_k, \hat{\eta}_{k-1})}.$$

K.1.1. MIRROR DESCENT

In Algorithm 2, we need to solve the optimization problem in (5), namely,

$$\inf_{\lambda \in \Delta_n} \mathbb{E}_{\zeta \sim N(0, I)} [\max_{h \in \mathcal{H}} f(\lambda; h; \zeta)].$$

We use stochastic mirror descent method. Specifically, given a current iterate $\hat{\lambda}$ and an i.i.d. sample ζ_1, \dots, ζ_B we computed an unbiased estimate of the sub-gradient

$$g(\hat{\lambda}) = \frac{1}{B} \sum_{i=1}^B \nabla_{\lambda} \max_{h \in H} f(\lambda; h; \zeta_i) \Big|_{\lambda=\hat{\lambda}},$$

and then move to $\hat{\lambda}_+ = \Pi(\exp(\log(\hat{\lambda}) - \eta g(\hat{\lambda})))$ where $\Pi(x) = x/\|x\|_1$ and \log, \exp are applied element-wise. The step size η is chosen by back-tracking line search where two step sizes are equivalent if the difference in estimated function values is dominated by the the square root of the empirical variance, with respect to the finite batch size estimates.

The batch size was grown adaptively. Let $\bar{f}(\lambda) = \mathbb{E}_{\zeta \sim N(0, I)}[\max_{h \in \mathcal{H}} f(\lambda; h; \zeta)]$, $\lambda_* \in \arg \min_{\lambda} \bar{f}(\lambda)$ and suppose the algorithm is at some current configuration $\hat{\lambda}$. By convexity of \bar{f} we have

$$\begin{aligned} \bar{f}(\hat{\lambda}) - \bar{f}(\lambda_*) &\leq \langle \nabla \bar{f}(\hat{\lambda}), \hat{\lambda} - \lambda_* \rangle \\ &= \langle \nabla \bar{f}(\hat{\lambda}) - g(\hat{\lambda}), \hat{\lambda} - \lambda_* \rangle + \langle g(\hat{\lambda}), \hat{\lambda} - \lambda_* \rangle \\ &\leq 2 \max_k |[\nabla \bar{f}(\hat{\lambda}) - g(\hat{\lambda})]_k| + \max_k \langle g(\hat{\lambda}), \hat{\lambda} - \mathbf{e}_k \rangle \\ &\approx 2 \max_k \hat{\sigma}_k + \max_k \langle g(\hat{\lambda}), \hat{\lambda} - \mathbf{e}_k \rangle \end{aligned}$$

where $\hat{\sigma}_k^2$ is the empirical variance of $[g(\hat{\lambda})]_k$, namely, the sample variance of $\{[\nabla_{\lambda} \max_{h \in H} f(\lambda; h; \zeta_i)]_k\}_{i=1}^B$ divided by B . Note that $\hat{\sigma}_k = O(1/\sqrt{B})$. Thus, we double the batch size $B \mapsto 2B$ whenever $2 \max_k \hat{\sigma}_k \geq \max_k \langle g(\hat{\lambda}), \hat{\lambda} - \mathbf{e}_k \rangle$. This also motivates our stopping condition: for input ϵ , terminate when $2 \max_k \hat{\sigma}_k + \max_k \langle g(\hat{\lambda}), \hat{\lambda} - \mathbf{e}_k \rangle \leq \epsilon$. Because B is increasing over time, this stopping condition will always, eventually, be met.

For a fixed ζ_i , to compute the corresponding gradient, we computed $\bar{h} = \arg \max_{h \in \mathcal{H}} f(\lambda; h; \zeta)$ and used the fact that $\nabla_{\lambda} \max_{h \in \mathcal{H}} f(\lambda; h; \zeta) = \nabla_{\lambda} f(\lambda; \bar{h}; \zeta)$. As described in the next section, finding \bar{h} is difficult due to the use of surrogate loss functions. We provide a line search method, $\text{LineSearch}(\lambda, \zeta)$, that finds an approximate value for it. Together, the full optimization is as follows:

Algorithm 6 SMD($\tilde{h}_k, \hat{\eta}_{k-1}$).

Goal: Solve for $\max_{\lambda} \mathbb{E}_{\zeta \sim N(0, I)}[f(\lambda; h; \zeta)]$, where f inherently depends on \tilde{h}_k and $\hat{\eta}_{k-1}$.

Input: Tolerance ϵ for stopping criteria.

Initialize: $\hat{\lambda} \leftarrow \frac{1}{d} \mathbf{1}$.

repeat

 Sample $\zeta_1, \dots, \zeta_B \sim N(0, I)$.

 Compute $\bar{h}_i = \text{LineSearch}(\hat{\lambda}, \zeta_i)$ for all $i \in [B]$ and their corresponding gradient estimates $\nabla_{\lambda} f(\hat{\lambda}; \bar{h}_i, \zeta_i)$.

$g(\hat{\lambda}) \leftarrow \frac{1}{B} \sum_{i=1}^B \nabla_{\lambda} f(\hat{\lambda}; \bar{h}_i, \zeta_i)$.

$\hat{\lambda} \leftarrow \Pi(\exp(\log(\hat{\lambda}) - \eta g(\hat{\lambda})))$ where η is chosen as described above.

if $2 \max_k \hat{\sigma}_k \geq \max_k \langle g(\hat{\lambda}), \hat{\lambda} - \mathbf{e}_k \rangle$ **then**

$B \leftarrow 2B$.

end if

until $2 \max_k \hat{\sigma}_k + \max_k \langle g(\hat{\lambda}), \hat{\lambda} - \mathbf{e}_k \rangle \leq \epsilon$

K.1.2. LINE SEARCH

Then, computing $\max_{h \in \mathcal{H}} f(\lambda; h; \zeta)$ is equivalent to

$$\min_{r \in \mathbb{R}^+} r \quad \text{subject to } g(r) \leq 0$$

where $a = -2^{-k+1} - \sum_{i \in [n]} (1 - 2\hat{\eta}_{k,i}) \tilde{h}_k(x_i)$, $b = \sum_{i \in [n]} \frac{\zeta_i}{n\lambda_i^{1/2}} \tilde{h}_k(x_i)$, $c_i = 1 - 2\hat{\eta}_{k-1,i}$, $d_i = -\frac{\zeta_i}{n\lambda_i^{1/2}}$ and

$$g(r) = ar + b + \max_{h \in \mathcal{H}} \sum_{i \in [n]} (c_i r + d_i) h(x_i).$$

In particular at the optimal value of r , called r^* , $\arg \max_{h \in \mathcal{H}} f(\lambda; h; \zeta) = \arg \max_{h \in \mathcal{H}} g(r^*)$.

As shown in Section 4.2, given access to a weighted classification oracle, computing $\max_{h \in \mathcal{H}} \sum_{i \in [n]} (c_i r + d_i) h(x_i)$ is equivalent to a 0/1-loss minimization problem that is solvable using a weighted classification oracle. If we had such an oracle then since $g(r)$ is monotonically decreasing as a function of r we can use a binary search procedure to solve this optimization problem. Indeed, for any given range $[r_{\min}, r_{\max}]$ where the optimal r lies in, we are checking if $g(\frac{r_{\min} + r_{\max}}{2}) \leq 0$. If that is the case, we just set $r_{\min} \leftarrow \frac{r_{\min} + r_{\max}}{2}$, otherwise, we set $r_{\max} \leftarrow \frac{r_{\min} + r_{\max}}{2}$. We then check the sign of $g(\frac{r_{\min} + r_{\max}}{2})$ again and repeat this procedure until a sufficient tolerance is met.

However, in practice, we do not have access to a weighted 0/1 loss oracle and must employ a convex surrogate loss that may not correctly solve the weighted classification problem. To be concrete, in all of our experiments we used Scikit-learn's `LogisticRegression` classifier. Given such a surrogate, which we denote as $\widetilde{\max}_{h \in H}$ to acknowledge that it may not find the optimal h , the resulting function

$$\tilde{g}(r) = ar + b + \widetilde{\max}_{h \in \mathcal{H}} \sum_{i \in [n]} (c_i r + d_i) h(x_i)$$

may no longer be monotonically decreasing in r hence a binary search procedure would fail. However intuitively it suffices to look at a large enough set of r 's near a zero $\tilde{g}(r)$. To overcome this issue, we used the procedure in Algorithm 8.

Algorithm 8 overcomes this issue by considering a large set of potential r values and for each value computing the corresponding $h_r = \arg \max_h \tilde{g}(r)$ and stores these in the array S . It then returns $\arg \max_{h \in S} f(\lambda; h; \zeta)$. The set of r values considered is chosen by a multi-scale procedure that looks at r values on a finer and finer geometric grid given a budget N_{\max} .

Algorithm 7 Oracle(r, S).

Goal: Solve for $\widetilde{\max}_{h \in \mathcal{H}} ar + b + \sum_{i \in [n]} (c_i r + d_i) h(x_i)$.
 Compute $\bar{h} = \widetilde{\max}_{h \in \mathcal{H}} \sum_{i \in [n]} (c_i r + d_i) h(x_i)$ by the relaxed weighted classification oracle.
 Set $\hat{g} \leftarrow ar + b + \sum_{i \in [n]} (c_i r + d_i) \bar{h}(x_i)$.
 $S \leftarrow S \cup \{\bar{h}\}$.
Return: \hat{g}, S .

Algorithm 8 LineSearch(λ, ζ).

Goal: Solve for $\max_{h \in \mathcal{H}} f(\lambda; h; \zeta)$.
Input: fixed λ and ζ ; maximum number of iterations N_{\max} ; tolerance ϵ .
Initialize: $S \leftarrow \{\}$, $r \leftarrow 100$, $\gamma \leftarrow 10$, $\delta \leftarrow \sqrt{2}$, and $t \leftarrow 0$.
 $\hat{g}, S \leftarrow \text{Oracle}(r, S)$.
while $\hat{g} < 0$ and $t < N_{\max}$ **do**
 $r \leftarrow r/2$. { r is too large}
 $\hat{g}, S \leftarrow \text{Oracle}(r, S)$.
 $t \leftarrow t + 1$.
end while
for $j = t, \dots, N_{\max} - 1$ **do**
 $\hat{g}, S \leftarrow \text{Oracle}(r, S)$.
 if $\hat{g} > 0$ **then**
 $r \leftarrow \gamma \cdot r$. { r is too small, need to increase r so that \hat{g} decreases.}
 else
 $r \leftarrow r/\gamma^2$. {Reaches a point where r is too large, scale back.}
 $\gamma \leftarrow \gamma/\delta$. {Start searching with a finer grid scale.}
 end if
end for
Return: $\arg \max_{h \in S} f(\lambda; h; \zeta)$.

K.2. Sampling

The last portion of our algorithm is the sampling scheme. The algorithm described in Algorithm 2 does not reuse samples between rounds to compute an estimate for $\hat{\eta}_k$. In practice, this can be very wasteful and we instead want an estimator

based on all samples up to and including those taken in round k . In each round the algorithm computes λ_k with the goal of $\lambda_k \approx \lambda_*$ for large enough values of k to ensure that we are sampling from the optimal distribution in that round. Hence, since we are recycling samples, we need to ensure that the distribution of *all* samples taken by the end of round k , including those in previous rounds, match λ_k .

To ensure those, we use a waterfilling technique. We set $p_1 = \lambda_1$ and then for each $k \geq 1$ we set

$$p_k = \arg \min_{q \in \Delta_n} \max_{j \leq n} \max \{0, k \cdot \lambda_{k,j} - (\sum_{i=1}^{k-1} p_{i,j}) - q_j\}. \quad (85)$$

Our algorithm being implemented is shown below:

Algorithm 9 Fixed Budget ACED with Waterfilling.

Input: Budget T , tolerance $\epsilon > 0$, batch size N (default 250).

$\hat{\eta}_0 \leftarrow 0$

for $k = 1, 2, \dots, \lfloor \log_2(\epsilon^{-1}) \rfloor$ **do**

$\tilde{h}_k \leftarrow \arg \min_{h \in \mathcal{H}} \widetilde{\text{err}}(h, \hat{\eta}_{k-1})$.

Optimization

$\lambda_k \leftarrow \text{SMD}(\tilde{h}_k, \hat{\eta}_{k-1})$.

Sampling

Sample till observing N unique $\{x_1^{(k)}, \dots, x_N^{(k)}\} \sim p_k$ that has not been queried before (from rounds 1, ..., $k-1$).

Query x_1, \dots, x_N and observe y_1, \dots, y_N .

Compute an estimate $\hat{\eta}_k$ with the naive estimator $\hat{\eta}_k^{(\text{Naive})}$ as defined in Section 5.

end for

Return: $\arg \min_{h \in \mathcal{H}} \widetilde{\text{err}}^{(k)}(h)$

K.3. Batched IWAL

In this section we explain our implementation of the IWAL algorithm (Beygelzimer et al., 2010) (and variants such as IWAL1 and oracular variants) for streaming active algorithms. In round k we assume access to a (labeled) dataset $S_{k-1} = \{(x_i, y_i, p_i)_{i=1}^{n_k}\} \subset \mathcal{X} \times \{\pm 1\} \times [0, 1]$, with $n_k \leq k$. Given a new point from a stream, x_k , the decision to label x_k is made by computing two hypothesis. Firstly we compute

$$h_k = \arg \min_{h \in \mathcal{H}} \sum_{i=1}^{n_k} \frac{\mathbf{1}\{h(x_i) \neq y_i\}}{p_i}$$

and second we compute

$$h'_k = \arg \min_{\substack{h \in \mathcal{H} \\ h_k(x_k) \neq h'_k(x_k)}} \sum_{i=1}^{n_k} \frac{\mathbf{1}\{h(x_i) \neq y_i\}}{p_i}.$$

Thus to compute h'_k we need access to a weighted classification oracle for 0/1 loss that can handle a single constraint.

In general including the constraint $h_k(x_k) \neq h'_k(x_k)$ is not easy for arbitrary classes and past methods have considered the optimization

$$\arg \min_{h \in \mathcal{H}} \sum_{i=1}^{n_k} \frac{\mathbf{1}\{h(x_i) \neq y_i\}}{p_i} + q \mathbf{1}\{h(x_k) \neq h'_k(x_k)\}$$

for a sufficiently large weight $q \in \mathbb{R}$ to ensure the constraint (Karampatziakis & Langford, 2010).

However in the case of linear classes under the surrogate convex logistic loss, the main focus of experiments in this paper, we take a different approach. Assume that $\mathcal{X} \subset \mathbb{R}^p$ and that $\mathcal{H} = \{h(x) = w^\top x + b : w \in \mathbb{R}^p, b \in \mathbb{R}\}$. W.l.o.g. assume that $h_k(x_k) = -1$. Under this convex relaxation, we seek a classifier where $w^\top x_k + b = \epsilon$, where $\epsilon \geq 0$ i.e. the linear

predictor flips the predicted sign of x_k and has margin ϵ . Thus we learn (an approximate) h'_k by solving

$$\begin{aligned} \min_{\substack{w \in \mathbb{R}^p, b \in \mathbb{R} \\ w^\top x + b = \epsilon}} \sum_{i=1}^{n_k} \frac{1}{p_i} \log(1 + \exp(-y_i(w^\top x + b))) &= \min_{w \in \mathbb{R}^p} \sum_{i=1}^{n_k} \frac{1}{p_i} \log(1 + \exp(-y_i(w^\top x_i - w^\top x_k + \epsilon))) \\ &= \min_{w \in \mathbb{R}^p} \sum_{i=1}^{n_k} \frac{1}{p_i} \log(1 + \exp(-y_i(w^\top (x_i - x_k) + \epsilon))) \end{aligned}$$

This is a convex optimization problem with no constraints that is easily solvable using a procedure for logistic regression where we assume that the intercept is some ϵ with really small magnitude.

L. Full Scale Plots for Performances on the Pool

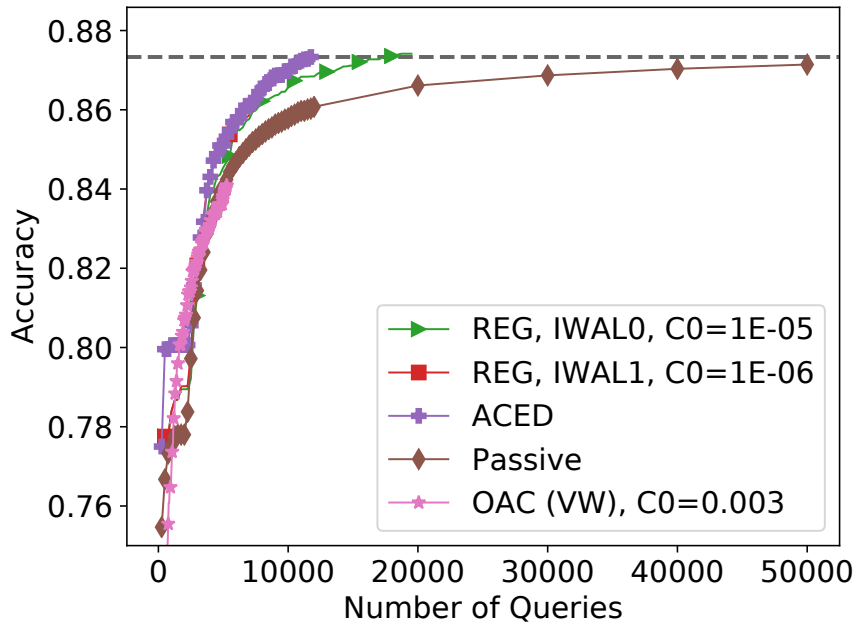


Figure 5. Full scale of Figure 1

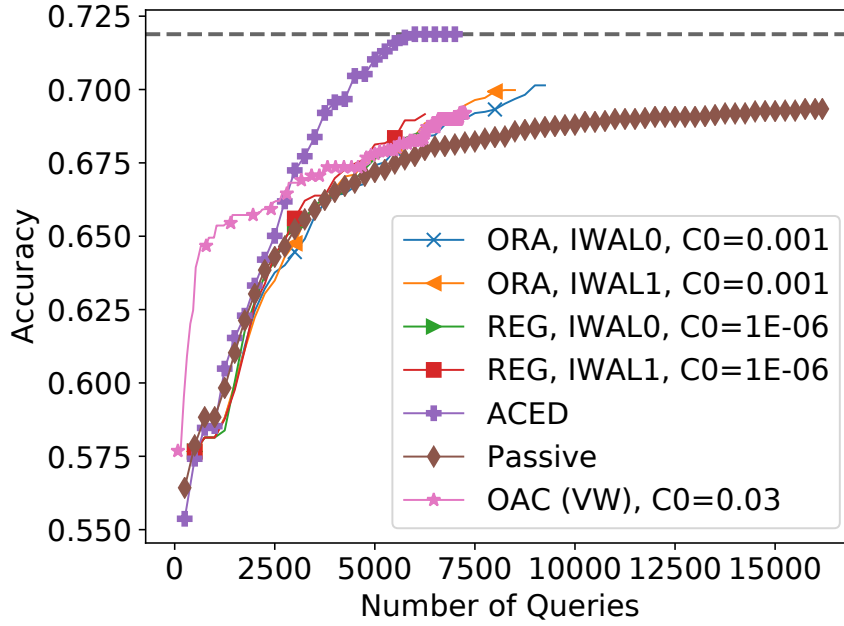


Figure 6. Full scale of Figure 2

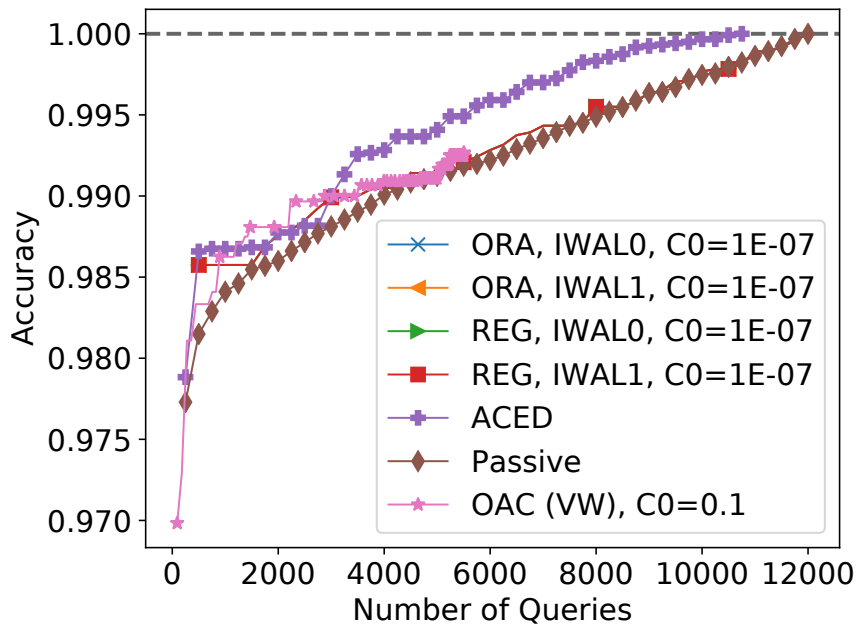


Figure 7. Full scale of Figure 4

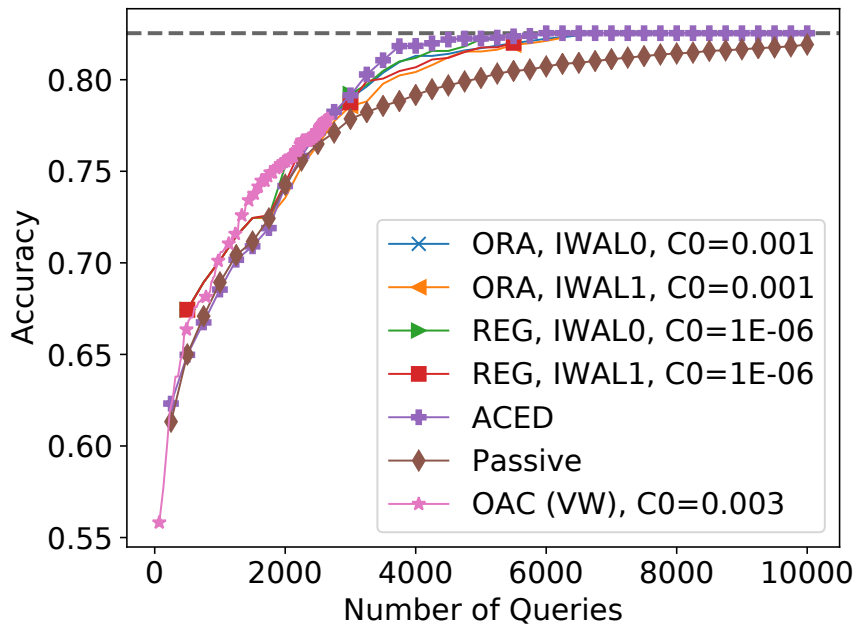


Figure 8. Full scale of Figure 3

M. Hyperparameters for Baselines

We searched in the following C_0 that uses the same grid fineness as (Huang et al., 2015):

Baseline	MNIST	SVHN	FashionMNIST	CIFAR
IWAL-0&1	$10^{-7}, 10^{-6}, \dots, 1$	$10^{-7}, 10^{-6}, \dots, 1$	$10^{-8}, 10^{-7}, \dots, 1$	$10^{-7}, 10^{-6}, \dots, 1$
ORA-IWAL-0&1	N/A	$10^{-4}, 10^{-3}, \dots, 10^{-1}$	$10^{-8}, 10^{-7}, \dots, 1$	$10^{-4}, 10^{-3}, \dots, 10^{-1}$
OAC	.001, .003, .01, .03	.01, .03, .1, .3	.001, .003, ..., 1	.001, .003, .01, .03

 Table 1. Ranges of C_0 searched for different experiments.

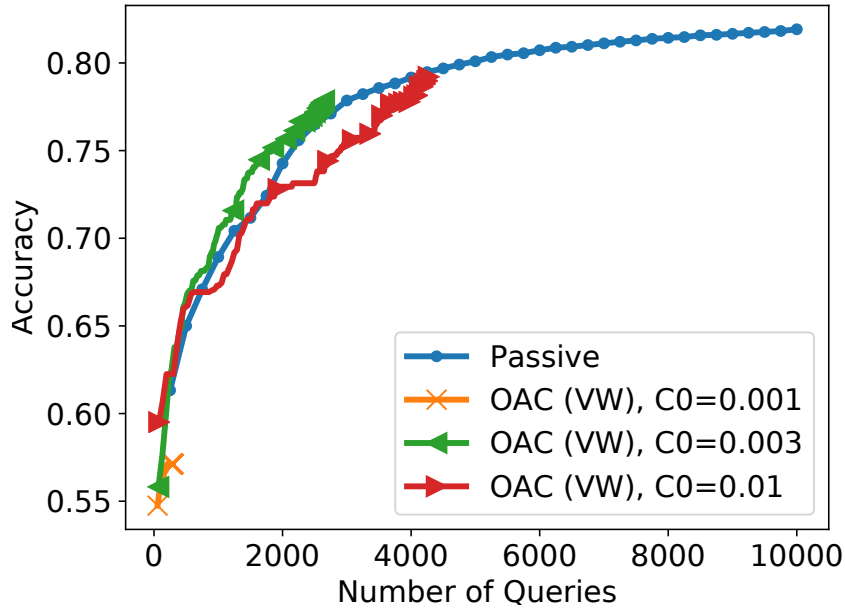


Figure 9. Sensitivity of OAC to C_0 on the CIFAR 2 vs 7 dataset (training accuracy).

We used the following amount of passes over dataset for the following baselines. For OAC, we made sure the number of passes is sufficient enough so that the algorithm is no longer taking more queries.

Baseline	MNIST	SVHN	FashionMNIST	CIFAR
IWAL-0&1 and ORA-IWAL-0&1	1 (N/A for ORA variants)	2	2	2
OAC	5	10	10	10

Table 2. Ranges of C_0 searched for different experiments.

For Vowpal Wabbit, we used an initial learning rate of 0.5 for CIFAR, and 1 for every other experiments.

N. Generalization Performance on Holdout Set

In the following figures, we show performances of the algorithms on a holdout test set. We note that there’s no algorithm that is consistently the best among all four experiments, but ACED is consistently among the top two algorithms, and is the best in two of the four experiments. We also reiterate that the OAC curves have stopped taking more queries at the end, so their final accuracies are inferior than other algorithms in most cases.

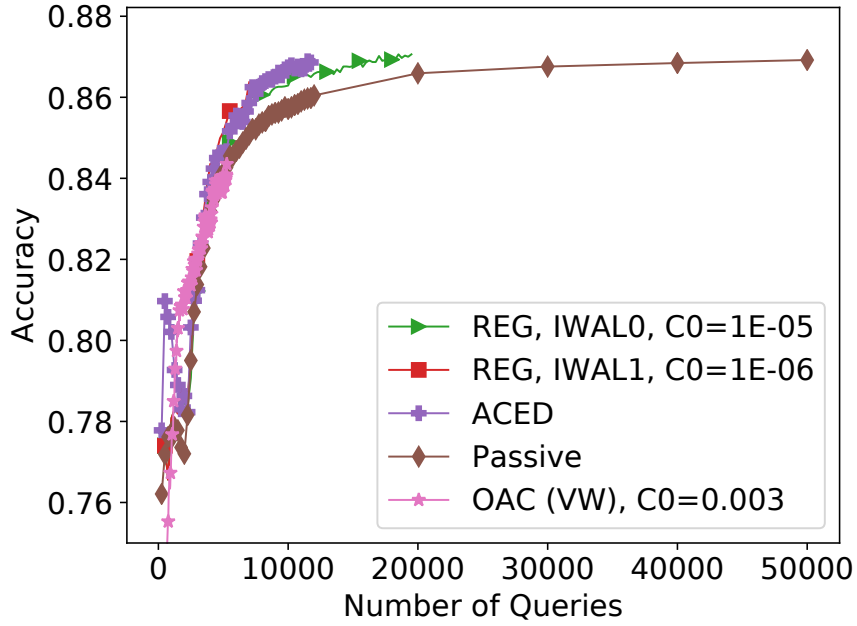


Figure 10. MNIST performance on test set

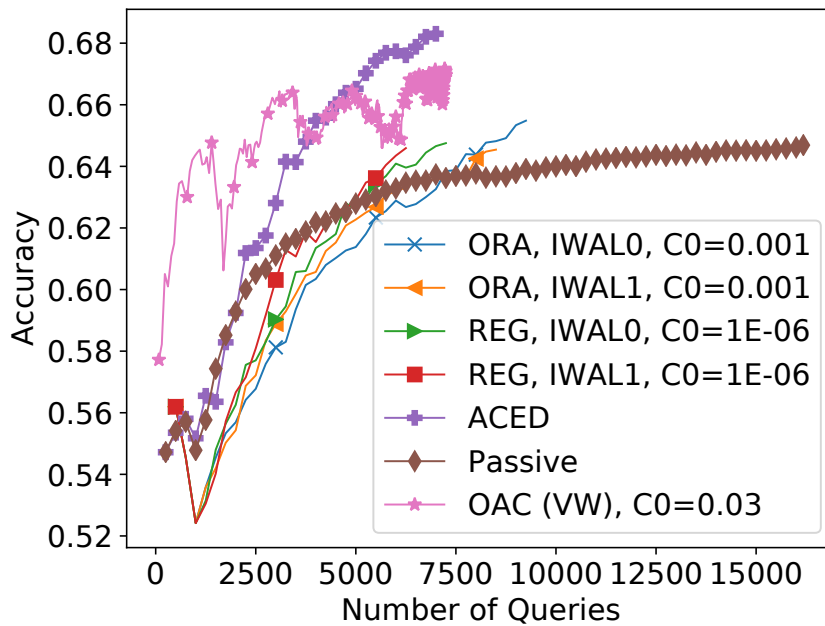


Figure 11. SVHN performance on test set

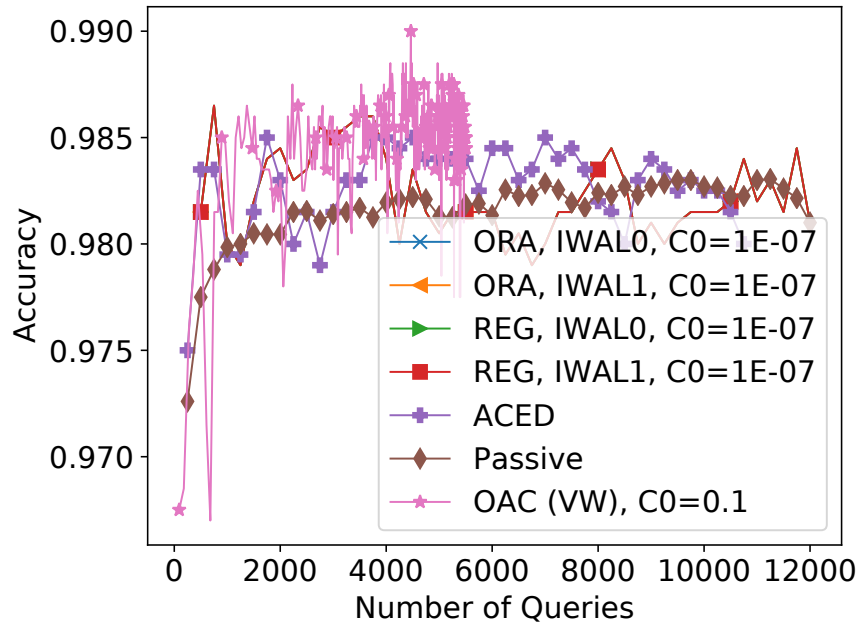


Figure 12. FashionMNIST performance on test set

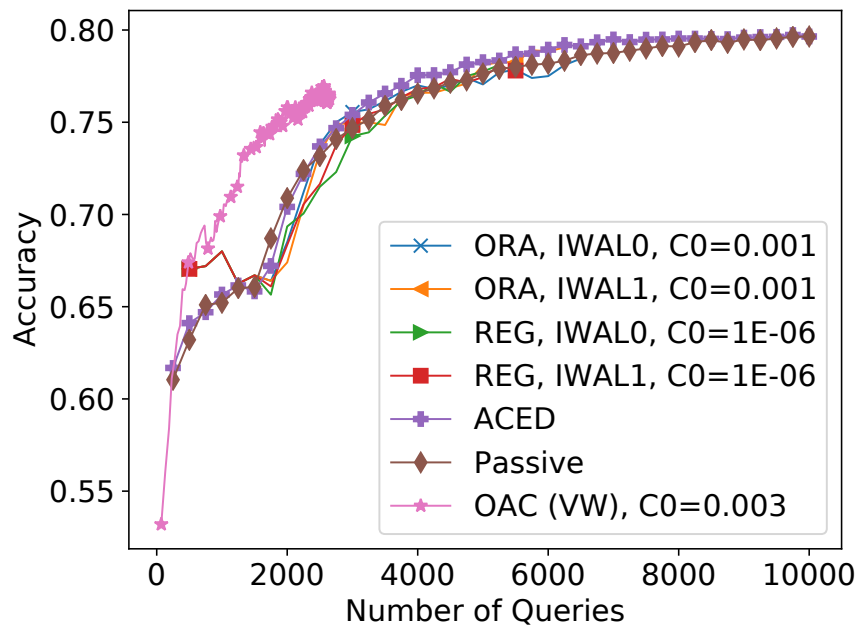


Figure 13. CIFAR performance on test set