

Learning with Attribute Costs

Haim Kaplan^{*}
Computer Science Dept.
Tel-Aviv university
haimk@post.tau.ac.il

Eyal Kushilevitz
Computer Science Dept.
Technion
eyalk@cs.technion.ac.il

Yishay Mansour[†]
Computer Science Dept.
Tel-Aviv university
mansour@cs.tau.ac.il

ABSTRACT

We study an extension of the “standard” learning models to settings where observing the value of an attribute has an associated cost (which might be different for different attributes). Our model assumes that the correct classification is given by some target function f from a class of functions \mathcal{F} ; most of our results discuss the ability to learn a clause (an OR function of a subset of the variables) in various settings:

Offline: We are given both the function f and the distribution D that is used to generate an input x . The goal is to design a strategy to decide what attribute of x to observe next so as to minimize the expected evaluation cost of $f(x)$. (In this setting there is no “learning” to be done but only an optimization problem to be solved; this problem is to be NP-hard and hence approximation algorithms are presented.)

Distributional online: We study two types of “learning” problems; one where the target function f is known to the learner but the distribution D is unknown (and the goal is to minimize the expected cost including the cost that stems from “learning” D), and the other where f is unknown (except that $f \in \mathcal{F}$) but D is known (and the goal is to minimize the expected cost while limiting the prediction error involved in “learning” f).

Adversarial online: We are given f , however the inputs are selected adversarially. The goal is to compare the learner’s cost to that of the best fixed evaluation order (i.e., we analyze the learner’s performance by a competitive analysis).

^{*}Work partially supported by Israel Science Foundation (ISF) Grant no. 548.

[†]This work was supported in part by the IST Programme of the European Community, under the PASCAL Network of Excellence, IST-2002-506778, by a grant no. 1079/04 from the Israel Science Foundation and an IBM faculty award. This publication only reflects the authors’ views.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

STOC’05, May 22-24, 2005, Baltimore, Maryland, USA.
Copyright 2005 ACM 1-58113-960-8/05/0005 ...\$5.00.

General Terms

Algorithms, Theory

Categories and Subject Descriptors

I.2.6 [Learning]; F.2.2 [Nonnumerical Algorithms]

Keywords

Learning, Online, Approximation Algorithms

1. INTRODUCTION

This paper is motivated by the observation that, in many cases, there is an inherent cost associated with observing each attribute of a given input x . Our goal is to model such scenarios and to offer algorithms that address this concern. The main variant of our model can be viewed as an extension of the classical PAC model [16], and has been proposed recently in the Machine Learning literature [14, 13]. To motivate the various variants of our model, we use a simple (hypothetical) medical setting. Each input x represents a patient where each attribute x_i represents the outcome of a medical test for this patient.

Consider a setting in which a doctor would like to diagnose if a patient has some disease. The doctor can run various tests on the patient in order to deduce the outcome; In fact, one would expect the order of tests to be determined so as to minimize the expected cost, an important criteria for any HMO (intuitively, one would expect that the doctor would run a simple blood test before using an expensive MRI; however, the probability of passing a certain test also influences its cost-effectiveness). This setting can be viewed as an *offline* setting, where the function f that the doctor evaluates (i.e., the final diagnosis based on the test results) is known, the distribution of the test outcomes (attributes) in the general public is also known, and the goal is to minimize the expected cost of evaluating $f(x)$.

Alternatively, consider a doctor conducting a research in order to formulate a diagnosis function f . She has the set of tests that she can run on each patient; after observing some of these tests she predicts the outcome. Since f is not known, the predictions might be incorrect. The doctor later learns whether her prediction was correct or not. Clearly, we would like to minimize those mistakes (for the insurance company sake) while maintaining a low overall cost, which is composed from learning the function and the evaluation process (for the HMO sake). We call this model the *distributional online* model, since the patients arrive in an online way; however, they are drawn from some fixed underlying

distribution D . The above describes the case where the target function f is unknown yet the distribution D might be known. An alternative case is when D is unknown but f is known; the goal is to minimize the expected cost of the tests, while collecting information regarding the distribution. (This scenario can be motivated in cases of changing the patients’ population; e.g., in the case of evaluating on humans a function learned using monkeys, where the function f is the same, but the distribution might be different.)

Finally, we consider the *adversarial online* model, where the target function is known but the adversary controls the arrival process. The hope is to bound the average online cost compared to a fixed evaluation order.

1.1 Our Model and Notation

We consider Boolean functions over the domain $\{0, 1\}^n$. We assume that each attribute x_i has an associated cost c_i . The evaluation of a given function f on input x is done by specifying a decision tree T of the attributes. This decision tree T is used to decide what attributes of x to observe next, given the outcome of the attributes that we observed. The value of f is written at the leaves of the tree. The cost of evaluating an input x using a tree T is the sum of the costs of the attributes along the path that x follows in T .

We mostly deal with simple functions, like clauses, where seeing an attribute which is one determines the value of the function. For such functions we can limit the evaluation model to a *leveled decision tree*. In a leveled decision tree, all the (internal) nodes at a same level (distance from the root) of the tree, test the same attribute of x . For such a tree we might encounter attributes which do not have an influence on the target function f . (Specifically, attribute x_i has influence on a function g if there exists two inputs x, x' that differ only on their i -th attribute and such that $g(x) \neq g(x')$.) In case we encounter such attribute at a node v , we assume that the subtree rooted at the left child of v and the subtree rooted at the right child of v are identical, and hence we can continue without evaluating the attribute at v . For this reason, the cost of evaluating an input x using a leveled decision tree T is the sum of the costs of the influential attributes along the path that x follows in T . A leveled decision tree is in fact equivalent to a *permutation* of the attributes, that determines the order in which we inspect them. While evaluating x we inspect the attributes one by one, in the order determined by the permutation, skipping attributes that have no influence, until the value of the function is determined. We pay the sum of the costs of the attributes that we inspected.

As mentioned, we consider target functions f that belong to certain simple classes of functions over $\{0, 1\}^n$. The target function f is used to give the correct classification for each input. A *Disjunctive Normal Form* (DNF) is an OR of terms, where each term is an AND of attributes. A Read-Once DNF has every attribute appearing in at most one term. A *Conjunctive Normal Form* (CNF) is an AND of clauses, where each clause is an OR of attributes.

The input generation is done as follows. In the *adversarial online model* the inputs are chosen by an adversary. In all other models, there is a fixed distribution D and the inputs are sampled i.i.d. using D . The description of a distribution (when it is known) is given as follows. A product distribution is specified by the probabilities p_i of each attribute x_i to get the value 1 (wlog, $0 < p_i < 1$). In the case of a general

distribution D , we assume a *conditional distribution oracle* that, given any assignment to a subset of the attributes, returns the conditional probability that the target function f is one, e.g., $\Pr_D[f = 1 | x_3 = 1, x_7 = 0]$.

The algorithm needs to compute (or to predict, in some settings) the value of the target function f for a given input x . By “computing the value”, we mean that the algorithm can “explain” this value; i.e., if $f(x) = 1$ then the algorithm outputs a term of the function which is true, and if $f(x) = 0$ the algorithm outputs a clause which is zero. (Namely, even if we are given a DNF which is always true, we still need, for every input, to exhibit a term which is true!) If $f(x) = 1$ we say that a term which equals 1 is a certificate for $f(x)$, and if $f(x) = 0$ we say that a clause which equals 0 is a certificate for $f(x)$. In case $f(x) = 1$ we denote the smallest cost of a term which is 1 by $\text{CER}(f, x)$, otherwise we denote the smallest cost of a clause which is 0 by $\text{CER}(f, x)$.

In general our algorithms would have two performance measures. The first is the cost of evaluating $f(x)$. The second, which is relevant when the target function is unknown, is the number of mistakes performed until the correct target function is identified.

1.2 Our Results

1. Known function and known distribution (Offline Computation): Given a function f and a distribution D , compute the optimal/approximate permutation that minimizes the average cost of computing f . (In this setting there is no learning and the main task is an optimization.)

For a general distribution D , the optimization problem is NP-Complete and we present a 4-approximation algorithm for the case that f is a clause. We generalize the results to the case where f can be represented as a read-once DNF. For monotone functions that have a small DNF and CNF (e.g., small decision trees), we give a logarithmic approximation algorithm when the cost and the distribution are both uniform.

2. Known function and unknown distribution (Distributional Online I): Given a target function $f \in \mathcal{F}$ (but no information about D), evaluate f while minimizing the average evaluation cost.

For the case where f is a clause, we show a learning algorithm for which, with probability $1 - \delta$, the cost of evaluating the first T examples is bounded by $O(T \cdot E[\text{OPT}_D] + \text{OPT}_D(S) \cdot \text{poly}(n))$, where $E[\text{OPT}_D]$ is the expected cost of the optimal permutation and $\text{OPT}_D(S)$ is the cost of OPT_D on a subset of the examples S , where $|S| = \text{poly}(n, \log 1/\delta)$.

3. Unknown function and known distribution (Distributional Online II): Given a class of functions \mathcal{F} , learn the target function $f \in \mathcal{F}$ while minimizing the average evaluation cost as well as the number of mistakes the algorithm performs. (There is a clear tradeoff between the two parameters, our aim would be to have a bounded number of mistakes and get our average cost to converge to the optimal average cost.)

We study product distributions and target functions that are clauses. We show how the online algorithm can guarantee at most $O(n)$ mistakes, while its expected cost is at most the optimal expected cost, in the case it is given the probability that $f = 1$, and otherwise it makes at most $O(n \log n)$ mistakes and has near optimal expected cost.

4. Adversarial online model: Given a function f , at each time step the online algorithm selects a permutation and the adversary selects an input.

For the lower bound, we consider a simple case where the target function is a clause, known to the learner, and all the costs are equal. We show that the adversary can force any deterministic algorithm to have average cost of $n/2$ while there always exists a fixed ordering of the attributes that has average cost at most 2. For the upper bound, we show that if we evaluate the attributes in order of increasing costs then the total cost is at most n times the optimal cost, for any input and any target function.

It is important to note that our results do not depend on the magnitude of the costs, or on the largest possible cost. This is a very important feature which makes both our algorithms and analysis much more delicate.

1.3 Related work

The offline version of our problem, when the distribution is known and the function is a clause containing all the variables, is in fact a generalization of a problem which has been studied recently, under different names, one of which is Min-Sum Set Cover (MSSC). The input to the MSSC problem (like the input to the classical Set Cover problem) consists of a collection of sets $\mathcal{F} = \{S_1, \dots, S_m\}$ containing elements from some universe \mathcal{U} . The goal is to find a permutation $S_{\pi(1)}, S_{\pi(2)}, \dots, S_{\pi(m)}$ that minimize the average covering times $f(e)$ for all $e \in \mathcal{U}$. The covering time of $e \in \mathcal{U}$ is the minimum j such that $e \in S_{\pi(j)}$. Bar-Noy et al. [3] show that the greedy algorithm for MSSC, which iteratively picks the set with the smallest ratio of cost to the number of uncovered elements which it covers, is a 4-approximation algorithm. More recently Feige, Lovasz, and Tetali [9] gave a simpler proof and also showed that it is NP-hard to approximate this problem to within a factor smaller than 4.

The MSSC problem had found several exciting applications recently. Cohen, Fiat, and Kaplan [6] considered a generalization of the problem in the context of exploiting associations to speed up search in P2P networks. Babu et al. [2, 15] used it in the context of recent work on data streams, where a stream had to go through a sequence of commutative *pipelined filters* that should be ordered so that average processing time is minimized. Some of our results may be interesting in the context of these applications as well.

The papers of Cohen et al. and of Babu et al. mentioned above, as well as this paper, generalize MSSC as defined in [9] in various ways: (1) sets are generalized to be predicates with certain success probability over some probability distribution, (2) predicates have non uniform costs, and (3) we may only have *estimates* on the success probability of each predicate. Theorem 2.4 unifies all these generalizations; The proof of Feige [9] with some changes suffices to establish it.

Etzioni et al. [7, 8] considered a related but different offline problem. In their work, attributes model information sources. An example x (a file, say) has $x_i = 1$ if the i^{th} information source stores x . They assume a product distribution, i.e., each source answers the query with fixed probability p_i . In their setup source i , in addition to cost c_i , also has a duration d_i which is the time that it takes source i to provide the answer. The basic problem is to find a schedule, that assigns to each source i a time t_i in which we query it, assuming we have not been able to get an answer for another

source before that time. The objective is to minimize the expected overall cost of the schedule, where the overall cost of the schedule is the sum of (1) the time it takes to perform the schedule, and (2) its cost (the sum of the costs of the sources we asked). They present a constant approximation algorithm for this basic problem.

Charikar et al. [5] also considered an offline problem similar to ours. They looked at functions defined by boolean AND/OR binary tree. Their goal was to find a decision tree with good *competitive ratio*. The *competitive ratio* of a decision tree T for a particular example x is the ratio between the cost of T on x , and $\text{CER}(f, x)$, the cost of the cheapest certificate for $f(x)$. The *competitive ratio* of T is the maximum of the competitive ratio of T over all examples x . Given a function f , Charikar et al. [5] produce a decision tree with competitive ratio bounded by $\max\{k, \ell\}$, where k is the size of the largest term in a minimal DNF representation of f , and ℓ is the size of the largest clause in a minimal CNF representation of f . Furthermore, they show how to design an algorithm specific to a function f and a cost vector c , achieving the best possible competitive ratio. In contrast our results in Section 2.2 give an algorithm for monotone functions f whose DNF and CNF representations are given, and uniform costs, that is $O(\log m)$ competitive, where m is the largest among the number of terms in the DNF representation, and the number of clauses in the CNF representation. Notice that $m \geq \max\{k, \ell\}$.

The most relevant work in the machine learning literature is the work of Madani et al. [14] which describes a fixed budget learning setting. In their setting, the learner needs to optimize its performance using a fixed budget. They study the coin problem – finding the coin with the highest bias; they show that in general the optimization problem is hard, and they consider specific priors and algorithms in the fixed budget model. In [13] the performance of the Naive Bayes estimator is evaluated using a Dirichlet prior.

Littlestone [12] formulated the online (mistake-bound) learning model which we use for presenting some of our algorithms. We also use the idea of “elimination” formulated by [12] (with roots in [16]) and used as a basic technique in other learning algorithms.

2. OFFLINE COMPUTATION

In this section, we discuss the offline problem; namely, both the target function f and the distribution D are known to the algorithm. For a target function which is a clause, as stated before, the optimal decision tree is a leveled decision tree that is equivalent to a permutation of the attributes. The following theorem characterizes the optimal offline permutation for a product distribution. (See, e.g., [10].)

THEOREM 2.1. *Given a clause $f = x_1 \vee \dots \vee x_k$ over the input $\{0, 1\}^n$, and a product distribution D , the optimal permutation is sorted according to c_i/p_i (from small to large). Hence, there is a polynomial time algorithm to compute the optimal permutation. \square*

For an arbitrary distribution given by an oracle¹ the following Theorem states that finding the best permutation

¹Recall that such oracle, given any assignment to a subset of the attributes, returns the conditional probability that the target function is one, e.g., $\Pr[f = 1 | x = 1, x_7 = 0]$

is hard. The proof is by a straightforward reduction from Min-Sum Set Cover (see Section 1.3).

THEOREM 2.2. *Deciding if the optimal average cost of a given clause f is at most K is NP-hard. \square*

Given that the problem is NP-Hard, it is natural to seek an approximation algorithm. A natural approximation algorithm is GREEDY, which selects the attributes in the permutation one after the other. The next attribute selected is the one that has the lowest marginal cost, i.e., $\mu_i = c_i/p_i$, where p_i is the probability that $x_i = 1$ given that all the attributes in the prefix are zero. Using the ideas of Feige et al. [9], we show that the GREEDY algorithm is an 4-approximation algorithm (The proof is a special case of Theorem 2.4).

THEOREM 2.3. *For any distribution D and any clause f , the expected cost of GREEDY is at most 4 times the expected cost of OPT. \square*

In the learning setting we might know the distribution only approximately. For this reason we consider a variation of algorithm GREEDY, called GREEDY'. At each iteration GREEDY' receives estimates of the conditional probabilities (rather than the exact conditional probabilities that GREEDY uses). Algorithm GREEDY' first determines $g(1)$, the index of its first attribute and, in general, it determines $g(k)$ after it has determined already $g(1), \dots, g(k-1)$. Let p_j^k be the probability that $x_j = 1$ given that the attributes $x_{g(1)}$ through $x_{g(k-1)}$ evaluate to 0. Algorithm GREEDY' is parameterized by two parameters ϵ and γ . We assume that when it determines $g(k)$, GREEDY' receives for x_i an estimate \hat{p}_i^k , such that: (1) if $\hat{p}_i^k \neq 0$ then $\hat{p}_i^k \in [(1-\epsilon)p_i^k, (1+\epsilon)p_i^k]$, and (2) for any $p_i^k > \gamma/n$ we have $\hat{p}_i^k > 0$. At phase k , GREEDY' computes $\hat{\mu}_\ell^k = c_\ell/\hat{p}_\ell^k$, and picks attribute x_ℓ with a minimal value $\hat{\mu}_\ell^k$, and places it in the k -th position in its permutation. To emphasize the dependency of GREEDY' on ϵ and γ we sometimes use the notation GREEDY'(ϵ, γ) instead of GREEDY'. We say that estimates \hat{p}_i^k which satisfy Conditions (1) and (2) above are *valid* for GREEDY'(ϵ, γ).

We prove the following upper bound on the cost of the permutation produced by GREEDY'.

THEOREM 2.4. *For any distribution D , the expected cost of GREEDY' is at most $4 \frac{(1+\epsilon)}{(1-\epsilon)(1-2\gamma)}$ times the expected cost of OPT $_D$.*

PROOF. We use the following notation for a permutation π . For an example e , let $I(e, \pi)$ be the minimum j such that $e_{\pi(j)} = 1$. We define: (1) $w_i^\pi = \Pr[I(x, \pi) = i]$ and $w_0 = 0$. (2) $\alpha_i^\pi = \sum_{j=1}^i w_j^\pi$, (3) $\beta_i^\pi = 1 - \alpha_i^\pi = \Pr[I(x, \pi) \geq i+1]$, (4) $C_i^\pi = \sum_{j=1}^i c_{\pi(j)}$ and (5) $D(S) = \Pr[x \in S]$, where (1)-(4) are defined for $0 \leq i \leq n$.

Let g denote the permutation produced by GREEDY' and let opt denote the permutation of OPT $_D$. For OPT $_D$ we define a step function h , such that $O(y) = C_k^{opt} = \sum_{j=1}^k c_j^{opt}$, where $y \in [\alpha_{k-1}^{opt}, \alpha_k^{opt}]$. Clearly the cost of OPT $_D$ is $\int_0^1 O(y)dy$. For GREEDY', we define a function G such that for $y \in [\alpha_{k-1}^g, \alpha_k^g]$, $G(y) = \beta_{k-1}^g \frac{c_{g(k)}}{w_k^g}$. Again, it is easy to see that the cost of GREEDY' is $\int_0^1 G(y)dy$. The crux of the proof is to show that for every k , $1 \leq k \leq n$,

$$O(1 - \beta_{k-1}^g/2) \geq G(1 - \beta_{k-1}^g) \frac{(1-2\gamma)(1-\epsilon)}{2(1+\epsilon)}.$$

This implies that if we shrink G by a factor of $\frac{2(1+\epsilon)}{(1-2\gamma)(1-\epsilon)}$ in the y dimension and by a factor of 2 in the x dimension it fits underneath O from which the theorem follows.

Consider OPT $_D$ after it covered 2 a $1 - \beta_{k-1}^g/2$ fraction of the space (and has cumulative cost $O(1 - \beta_{k-1}^g/2)$). Denote by S_1 the $1 - \beta_{k-1}^g/2$ fraction of the probability space that OPT $_D$ covers. Let S_2 be the subset of S_1 not covered by $\{x_{g(1)}, \dots, x_{g(k-1)}\}$. It follows that $D(S_2)$ is at least a $\beta_{k-1}^g/2$.

Denote by A the set of attributes x_j whose residual probability, p_j^k after GREEDY' picked its first $k-1$ attributes is smaller than γ/n . Let S_3 be the subset of S_2 which is covered by attributes from A . Since $|A| \leq n$, and each attribute in A covers at most γ/n of the remaining β_{k-1}^g fraction of the space that GREEDY' does not cover with its first $k-1$ attributes (i.e. those examples for which $x_{g(1)}, x_{g(2)}, \dots, x_{g(k-1)}$ are all zeros), we obtain that $D(S_3)$ is no larger than $\beta_{k-1}^g \cdot n(\gamma/n) = \beta_{k-1}^g \gamma$. Let $S_4 = S_2 - S_3$, we have that $D(S_4) \geq \beta_{k-1}^g(1/2 - \gamma)$. Obviously since $D(S_4) > 0$ there must be an attribute x_j whose residual probability $p_j^k \geq \gamma/n$. In particular $x_{g(k)}$ is such an attribute.

Let $q = \Pr[x_{g(1)} = 0, \dots, x_{g(k-1)} = 0]$. The probability $p_{g(k)}^k$ that $x_{g(k)} = 1$ given that $x_{g(1)}, x_{g(2)}, \dots, x_{g(k-1)}$ are all zeros is w_k^g/q by the definitions of w_k^g and q . So the residual cost of $x_{g(k)}$ after GREEDY' picked $x_{g(1)}, x_{g(2)}, \dots, x_{g(k-1)}$ is $c_{g(k)}/(w_k^g/q)$. The optimal algorithm OPT $_D$ must also have attributes $\notin \{x_{g(1)}, x_{g(2)}, \dots, x_{g(k-1)}\} \cup A$ to cover S_4 . Let x_o be one such attribute and let p_o^k the probability that x_o is one given that $x_{g(1)}, x_{g(2)}, \dots, x_{g(k-1)}$ are zero.

Since GREEDY' picks an attribute with smallest estimate of residual cost we know that the estimate that GREEDY' had for the residual cost of x_o was larger than the estimate GREEDY' had for the residual cost of $x_{g(k)}$ when it picked its k th attribute. If \hat{p}_o^k is the estimate that GREEDY' had for the residual probability of x_o and $\hat{p}_{g(k)}^k$ is the estimate that GREEDY' had for the residual probability of $x_{g(k)}$ then we have that $\frac{c_{g(k)}}{(1+\epsilon)w_k^g/q} \leq \frac{c_{g(k)}}{\hat{p}_{g(k)}^k} \leq \frac{c_o}{\hat{p}_o^k} \leq \frac{c_o}{(1-\epsilon)p_o^k}$

Rearranging the rightmost and leftmost sides we get that $\frac{(1-\epsilon)p_o^k c_{g(k)}}{(1+\epsilon)w_k^g} \leq c_o$.

Summing up over all attribute x_o that OPT $_D$ uses to cover S_4 we obtain that $O(1 - \beta_{k-1}^g/2) \geq \beta_{k-1}^g \frac{(1-2\gamma)(1-\epsilon)}{2(1+\epsilon)} \frac{c_{g(k)}}{w_k^g} = G(1 - \beta_{k-1}^g) \frac{(1-2\gamma)(1-\epsilon)}{2(1+\epsilon)}$, since $\sum_o p_o^k q \geq \beta_{k-1}^g \frac{1-2\gamma}{2}$. \square

Note that GREEDY is simply GREEDY' with $\gamma = 0$ and $\epsilon = 0$, and therefore it has an approximation ratio of 4.

2.1 Read-Once DNF

We first show that the optimal permutation for computing a Read-Once (R/O) DNF is a permutation on the terms. A *permutation of terms* is a permutation of the attributes that can be described as first a permutation of the terms and then for each term a permutation of its variables. (We will use this definition only for R/O DNF, in which case it is well defined.) Surprisingly, the following lemma shows that for any distribution and any set of costs there is an optimal permutation which is a permutation of terms.

²We use terminology related to set cover here. By saying that an algorithm "covers" some fraction p of the probability space with a set of attributes Y , we mean that for p fraction of the examples at least one attribute of Y is one.

LEMMA 2.5. *Let f be a R/O DNF and π a permutation. There exists a permutation π' which is a permutation of the terms, whose cost on any input is at most the cost of π .* \square

We can now extend Theorems 2.1 and 2.3 to R/O DNF.

THEOREM 2.6. *Let f be a function which is a R/O DNF. Then, (1) For a product distribution P , there is a polynomial time algorithm to compute the optimal permutation, and (2) for any distribution P , there is a polynomial time algorithm to compute a permutation whose expected cost is 4 times the optimal permutation.* \square

2.2 Small CNF and DNF

Assume that we are given a monotone function f that has at most m terms in its DNF representation and at most m clauses in its CNF representation. (We are given both the CNF and the DNF description of the function.) We assume uniform costs (i.e., $c_i = 1$ for every i), and uniform distribution D from which examples are drawn.

To formulate the upper bound, recall the definition in Section 1.1: $\text{CER}(f, x)$, when $f(x) = 1$ is the cost of the cheapest term of f which is one, and when $f(x) = 0$, $\text{CER}(f, x)$ is the cost of the cheapest clause of f which is zero.³ We show how to achieve an average cost which is bounded by the average (over D) of $\text{CER}(f, x)$ times $O(\ln m)$. Obviously the optimal cost for x is at least $\text{CER}(f, x)$ since the partial assignment it observes must always form a certificate.

THEOREM 2.7. *Let f be a monotone function that can be represented by both an m term DNF and by an m clause CNF. For uniform costs there is an algorithm that evaluates f with expected cost of $O(E_D[\text{CER}(f, x)] \ln m)$.*

PROOF. Since f is monotone we can assume, without loss of generality, that the DNF and CNF representations are both monotone (i.e., consisting only of non-negated variables). The idea is that we create two programs, one assuming that the function is one and the other assuming that the function is zero. We interleave the two programs, and show that the cost of the program with the correct assumption is at most a factor $O(\ln m)$ from the expected value of $\text{CER}(f, x)$. At each iteration we evaluate two variables. One that appears in the most terms (in the DNF) that were not yet eliminated and the other that appears in the most clauses (in the CNF) that were not yet eliminated.

Given an input x , assume that $f(x) = 1$ (the case $f(x) = 0$ is similar). To simplify the analysis, we assume that the cost is 1 only when the variable evaluates to 1 and zero otherwise. We denote this new cost function by cost_1 . (We will overcome this assumption at the end of the proof.) Since $f(x) = 1$, some term of f is one according to x . Let K be the size of the smallest such term T ; i.e., $\text{CER}(f, x) = K$. This implies that the K variables in T intersect all the clauses of f . The argument proceeds in a way similar to the analysis of the greedy algorithm for set cover; that is, we first observe that there is a variable that appears in m/K clauses. If the variable evaluates to 1 then we have cost 1 and eliminate m/K clauses. If the variable is zero then we incur no cost (and might not eliminate any clauses). Therefore, after evaluating at most $K \ln m$ variables to 1 we eliminate

³For uniform costs, the cheapest term is the one containing the smallest number of attributes.

all the clauses. This shows that the cost which our algorithm pays according to cost_1 is no larger than $\text{CER}(f, x) \ln m$.

Generalizing to a uniform cost function, we would like intuitively to claim that about half of the variables that are evaluated, turn out to be ones. Consider the i -th variable that we evaluate. Given that the value of the function is 1 and any values for the other variables, we claim that the probability that $x_i = 1$ is at least $1/2$. To see why this holds, consider x_1 . We have two type of inputs for which $f = 1$: either $f(0y) = f(1y) = 1$ or $f(1y) = 1$ and $f(0y) = 0$ (here we use the fact that f is monotone). Since we assume that the distribution is uniform, it is clear that the probability that $x_1 = 1$ given that $f(x) = 1$ is at least $1/2$.

To complete the proof, note that since $\Pr[x_i = 1 \mid f(x) = 1] \geq 1/2$ we obtain that for any strategy A to query the variables and any termination condition, $E_D[\text{cost}_1(A) \mid f(x) = 1] \geq (1/2)E_D[\text{cost}(A) \mid f(x) = 1]$.

The case that the function is zero is similar (eliminating the terms of the DNF, rather than the clauses of the CNF, when an attribute is zero). \square

3. ADVERSARIAL ONLINE

In this section we consider learning in the adversarial online setting. We start by presenting a lower bound, which applies for the case where the target function f is a clause (specifically, the OR of all n variables). We then, show an almost matching upper bound that applies to *all* target functions.

THEOREM 3.1. *Let the target function be $f = x_1 \vee \dots \vee x_n$ and a uniform cost ($c_i = 1$, for all i). For any deterministic online algorithm A for f there is an input sequence such that the cost of A is at least $n/4$ times the cost of the best static permutation.*

Next, we show that sorting the variables according to their cost (and using this permutation for all inputs), almost achieves the lower bound.

THEOREM 3.2. *Let f be a target function and π be a permutation of the attributes according to the cost (from small to large). Then, for any input x , the cost of π satisfies $\text{cost}(\pi, x) \leq n \cdot \text{CER}(f, x)$.*

4. ONLINE DISTRIBUTION MODEL I: UNKNOWN DISTRIBUTION

In this section, we address the case where the distribution D is unknown and arbitrary while the function f , which is a clause, is given. We assume that $f(x) = 1$ since if $f(x) = 0$ every algorithm has to observe all attributes to prove it. The basic approach is to sample from D , estimate the residual costs of the attributes, and run GREEDY' using the estimates. However, since some of the attributes may be much more expensive than others, we have to be careful not to pay too much for learning parts of D which are insignificant.

Consider an attribute x_i . By sampling from the unknown distribution D , we can compute an estimate, \hat{p}_i , of the probability $p_i \stackrel{\text{def}}{=} \Pr_D(x_i = 1)$, as follows. Draw L examples from D and check for each example if $x_i = 1$. Let b be the number of examples for which $x_i = 1$. If $b > \theta$, for some fixed threshold θ , we let $\hat{p}_i = b/L$. Otherwise we let $\hat{p}_i = 0$.

More specifically, let ϵ and δ' be two parameters (to be fixed). Set $\theta = \frac{6 \ln(3/\delta')}{(\epsilon)^2}$ and $L = 2n\theta/\gamma$; then, by Lemma 6.2, with probability at least $1 - \delta'$ we have that for all i : (1) if $p_i < \frac{\gamma}{4n}$ then our estimate \hat{p}_i is zero, (2) if $p_i \geq \frac{\gamma}{4n}$ then \hat{p}_i is not zero and is within $(1 \pm \epsilon)$ from p_i , and (3) if $\frac{\gamma}{4n} \leq p_i < \frac{\gamma}{n}$ and $\hat{p}_i \neq 0$, then \hat{p}_i is within $(1 \pm \epsilon)$ from p_i .⁴

Recall the offline algorithm $\text{GREEDY}'(\epsilon, \gamma)$ from Section 2. After fixing the first $k - 1$ attributes $x_{g(1)}, x_{g(2)}, \dots, x_{g(k-1)}$ in its permutation, g , the algorithm uses estimates \hat{p}_i^k on the probabilities p_i that $x_{g(1)}, x_{g(2)}, \dots, x_{g(k-1)}$ are all zeros and x_i is one, to choose its k -th attributes $x_{g(k)}$. We now describe a straightforward algorithm that answers examples drawn from D while computing estimates \hat{p}_i^k , as required by GREEDY' , and converging to its permutation g . We call this algorithm the Naive Learning algorithm (NL).

Algorithm NL maintains a sequence of attributes which is a *prefix* of the permutation that it gradually constructs. This prefix is constructed such that it is also a prefix that $\text{GREEDY}'(\epsilon, \gamma)$ would have constructed, for some valid estimates of the probabilities p_i^k .

The prefix is initially empty. In general assume that the prefix of NL consists of $k - 1$ attributes after answering some sequence of examples. For the next example, e , NL first observes the attributes in the prefix. If NL hits an attribute in the prefix which is one then NL stops there and reports that $f(e) = 1$. Otherwise, if all attributes in the prefix are zero then NL uses e to estimate p_i^k for every attribute x_i not in the prefix. To do that (and find out the value of $f(e)$ at the same time) NL observes the value of each of the remaining attributes in e . After L examples for which NL observed all attributes not in its prefix, NL has estimates for p_i^k that valid for $\text{GREEDY}'(\epsilon, \gamma)$ with probability $1 - \delta'$. That is, each estimate \hat{p}_i^k , with probability $1 - \delta'$, is either 0 if $p_i \leq \gamma/n$, or it is accurate to within $(1 \pm \epsilon)$ relative error. NL then picks the attribute with smallest residual cost according to its estimates and adds it to its prefix.

Taking $\delta' = \delta/n^3 2^n$ and using the union bound we obtain that with probability $1 - \delta$, NL's estimates for all p_i^k 's are valid for $\text{GREEDY}'(\epsilon, \delta)$. The following theorem summarizes the properties of NL, its proof is similar to the proof of Theorem 4.2 and hence omitted.

THEOREM 4.1. *Let $S = e_1, \dots, e_T$ be a sequence of examples to which we apply the algorithm NL. Let S' be the subset of S for which NL observes the attributes that are not in the prefix (notice that, by the definition of NL, we have $|S'| \leq nL$). Suppose that we run NL on S with $\delta' \leq \delta/n^3 2^n$. Then, with probability at least $1 - \delta$, the following hold.*

- (1) *Let $c^{\text{pre}}(e_j)$ be the cost paid by NL for observing attributes in the prefix at the time in which it processed e_j . Then, $\sum_{j=1}^T c^{\text{pre}}(e_j) = T \cdot \frac{4(1+\epsilon)}{(1-2\gamma)(1-\epsilon)} \cdot (1+\epsilon)E_D(\text{OPT}_D)$, where $E_D(\text{OPT}_D)$ is the expected cost of OPT_D on an example drawn from D .*
- (2) *The cost NL pays for observing additional attributes, for examples in S' , is at most $O(nL)$ times $\sum_{i=1}^n c_i = O(nc_{\max})$, where c_{\max} is the cost of the most expensive attribute.*

The weakness of this straightforward approach is that the cost that NL pays on S' (and therefore on the entire sequence) may be very large compared to the cost of OPT_D

⁴In Lemma 6.2 use $q = \frac{\gamma}{2n}$ and $F = \theta$. The number of samples m is then equal to L .

on the particular sample or even compared to the expected cost of OPT_D on the distribution D . This is because we may be observing an expensive attribute for about nL times while OPT_D may not observe it at all. To be able to relate the learning cost to the cost of OPT_D on the sample, we describe an algorithm which we call the *Greedy Learning (GL)* algorithm. The GL algorithm also uses a permutation which starts with a prefix, that is also a prefix of GREEDY' . However, in contrast with NL, when GL cannot resolve an example with the attributes in its prefix, it does not observe all other attributes. Instead, GL carefully observes attributes in its suffix one by one, in a particular order, to avoid large cost. Even before defining GL precisely, to contrast its performance with the performance of NL we state our main theorem.

THEOREM 4.2. *Let $S = e_1, \dots, e_T$ be a sequence of examples when a round of GL ends, and let S' be the subset of S for which GL observed attributes that are not in the prefix (as we shall see, from the definition of GL would follow that $|S'| \leq n^2 L$). Suppose we run GL on S with $\delta' = \delta/n^3 2^n$. Then, with probability at least $1 - \delta$, the followings hold:*

- (1) *Let $c^{\text{pre}}(e_j)$ be the cost paid by GL for observing attributes in its prefix while processing e_j . Then, $\sum_{j=1}^T c^{\text{pre}}(e_j) \leq T \cdot (1+\epsilon)E_D(\text{GREEDY}'(\epsilon, \gamma)) \leq T \cdot \frac{4(1+\epsilon)}{(1-2\gamma)(1-\epsilon)} \cdot (1+\epsilon)E_D(\text{OPT}_D)$.*
- (2) *The cost GL pays for observing additional attributes in the examples of S' is at most $O(n^2)$ times the cost paid by OPT_D on the examples in S' .*

It is worthwhile to compare the theorem regarding NL (Theorem 4.1) and that regarding GL (Theorem 4.2). The main and most important difference is in the second guarantee. While NL has a guarantee of $O(n^2 c_{\max} L)$ on its cost for processing the set S' , algorithm GL has a guarantee of $O(n^2)$ times the cost of OPT_D to process S' . This is a crucial difference, showing that GL is as good as OPT_D , up to a factor which is independent of the costs.

4.1 The Greedy Learning algorithm: GL

The *Greedy Learning (GL)* algorithm maintains a permutation of the attributes partitioned into two subsequences called *prefix* and *suffix* and one additional attribute in between the prefix and the suffix which is the *currently explored attribute*. The order of attributes in the prefix is determined by the algorithm. The attributes in the suffix are sorted by their costs. Each attribute in the suffix is marked either as “explored” or “unexplored”.

The GL algorithm progresses in *phases*, where each phase is divided to *rounds*. In phase k , let \mathcal{P}^k be the attributes in the prefix, and let \mathcal{S}^k denote the set of the remaining attributes. At round j of this phase, the explored attributes are $\mathcal{E}^{k,j}$, and the currently explored attribute is $a^{k,j}$, where $a^{k,j} \notin \mathcal{E}^{k,j}$ (the currently explored attribute does not change during a round). The permutation $\pi^{k,j}$ in round j of phase k has first the attributes in the prefix \mathcal{P}^k , in the order that they were added to the prefix, followed by $a^{k,j}$ followed but the attributes in $\mathcal{S}^k \setminus \{a^{k,j}\}$, sorted according to their cost.

In phase 0, we have an empty prefix, i.e., $\mathcal{P}^0 = \emptyset$, and all the attributes are marked “unexplored”. In the first round (of phase 0) the currently explored attribute, $a^{0,1}$, is the cheapest attribute. The permutation $\pi^{0,1}$ starts with $a^{0,1}$ followed by the remaining attributes sorted by their cost.

In general, in round j of phase k , we use permutation $\pi^{k,j}$ until we observe the currently explored attribute $a^{k,j}$ in L examples; (where L , and θ are defined by fixing δ' , ϵ , and γ) in which case round j terminates, $a^{k,j}$ is marked “explored”, and $\mathcal{E}^{k,j+1} = \mathcal{E}^{k,j} \cup \{a^{k,j}\}$. When round j terminates, we first check the number of times b , that $a^{k,j}$ was true (out of the L times it was observed). If $b < \theta$, we estimate $\hat{p}_i^k = 0$; otherwise, we estimate $\hat{p}_i^k = b/L$. We also set $\hat{\mu}_i^k = c_i/\hat{p}_i^k$ to be our estimate of the residual cost of attribute $a^{k,j}$. This completes a round.

Now we test if the phase terminates, or we should continue to the next round of the same phase. We calculate the minimum among the residual costs of the explored attributes in the phase, i.e., in $\mathcal{E}^{k,j+1}$, denote it by $\hat{w}^{k,j} = \min_{i \in \mathcal{E}^{k,j+1}} \hat{\mu}_i^k$. We also calculate the minimum cost among the costs of all unexplored attributes in the suffix, denote it by w' .

If $\hat{w}^{k,j} > w'$ then we set $a^{k,j+1}$ to be the unexplored attribute in the suffix with the cheapest cost, i.e., let $\ell = \arg \min_{i \in \mathcal{S}^{k,j+1}} c_i$, and $a^{k,j+1} = x_\ell$. We continue to round $j+1$ of phase k .

If $\hat{w}^{k,j} \leq w'$ we add to the prefix the explored attribute with the smallest residual cost in the suffix. I.e., $\mathcal{P}^{k+1} = \mathcal{P}^k \cup \{x_a\}$, where $a = \arg \min_{i \in \mathcal{E}^{k,j+1}} \hat{\mu}_i^k$, and $\mathcal{S}^{k+1} = \mathcal{S}^k \setminus \{x_a\}$. This completes phase k , and we move to phase $k+1$. At the beginning of phase $k+1$, we mark all elements that are not in \mathcal{P}^{k+1} as “unexplored” and make the cheapest one the currently explored attribute.

Algorithm GL simulates the algorithm GREEDY', described in Section 2, where at each iteration we try to avoid estimating the p_i 's of expensive attributes, when knowing they cannot change the decision of GREEDY'. Specifically, at each iteration we have estimates for p_i 's of the explored attributes in the suffix. If there are unexplored attributes in the suffix when we add the next attribute to the prefix, then it is guaranteed that, whatever the \hat{p}_i of those attributes are, GREEDY' could have taken the same decision.

4.2 The analysis of GL

We assume that we run GL on a sequence $S = e_1, \dots, e_T$ and after e_T a round of GL ends. To establish the first part of Theorem 4.2, we show that the prefix of GL is in fact a prefix that could have obtained by GREEDY'(ϵ, γ), and we also show that the cost of such prefix on S is close to T times its expected cost on an example from D . Then, to establish the second part of Theorem 4.2, we show that GL does not observe attributes which are too expensive even for those examples for which the attributes in the prefix are all zero. Immediately from Lemma 6.2 we obtain that:

LEMMA 4.3. *Let S' be a subset of S such that $|S'| \geq L$. Fix a set of attributes Y , and an attribute $x_i \notin Y$. Let p be the probability that all attributes in Y are zero and $x_i = 1$ for an example drawn from D . Let \hat{p} be the fraction of examples in S' for which all attributes in Y are zero and $x_i = 1$. Then with probability $1 - \delta'$, (1) if $\hat{p} \neq 0$ then $\hat{p} \in [(1 - \epsilon)p, (1 + \epsilon)p]$, and (2) for any $p > \gamma/n$ we have $\hat{p} > 0$. \square*

Combining this lemma with the union bound we obtain the following lemma and its corollary.

LEMMA 4.4. *Fix $\delta' = \frac{\delta}{n^3 2^n}$. Then with probability $1 - \delta$ in every round j of phase k , for every subset of attributes Y , and an attribute $x_i \notin Y$, the following hold: Let p be the probability that all attributes in Y are zero and $x_i = 1$ for*

an example drawn from D . Let \hat{p} be the fraction of examples in round j of phase k for which all attributes in Y are zero and $x_i = 1$. Then with probability $1 - \delta'$, (1) if $\hat{p} \neq 0$ then $\hat{p} \in [(1 - \epsilon)p, (1 + \epsilon)p]$, and (2) for any $p > \gamma/n$ we have $\hat{p} > 0$.

Let $c(e_j)$, $1 \leq j \leq T$, be the cost of GL for answering example e_j . We partition $c(e_j)$ into three components $c(e_j) = c^{pre}(e_j) + c^{cur}(e_j) + c^{suf}(e_j)$. The first part $c^{pre}(e_j)$ is the cost paid by GL for observing attributes in the prefix while processing e_j . The second part $c^{cur}(e_j)$ is the cost paid by GL if it observed the currently explored attribute while processing example e_j , otherwise $c^{cur}(e_j) = 0$. The third part $c^{suf}(e_j)$ is the cost incurred by GL for observing attributes in the suffix while processing e_j . Let $E^{pre} = \sum_j c^{pre}(e_j)$, $E^{cur} = \sum_j c^{cur}(e_j)$, and $E^{suf} = \sum_j c^{suf}(e_j)$. We bound E^{pre} , E^{cur} , and E^{suf} in the following. The next lemma bounds E^{pre} . By combining it with Theorem 2.4, we establish the first part of Theorem 4.2.

LEMMA 4.5. *If we use GL with $\delta' = \delta/n^3 2^n$ then, with probability $1 - \delta$ over the choice of S , we have $E^{pre} \leq T(1 + \epsilon)E_D(\text{GREEDY}')$.*

Our last two lemmas establish the second part of Theorem 4.2. Let $E^{cur} = \sum_k E_k^{cur}$ where E_k^{cur} is the contribution to E^{cur} of the cost of observing the currently explored attribute for examples during phase k . Assume we explore $j(k) > 1$ attributes during phase k . Then, the number of examples for which we observe the currently explored attribute during phase k is $b(k) = jL$. Let $x_{e(k)}$ be the $j(k)$ -th and therefore the most expensive attribute that GL explores at phase k . Then $E_k^{cur} \leq bc_{e(k)}$. The following lemma shows that with probability at least $1 - \delta$, for every phase k , OPT_D also pays at least $(1 - \epsilon)^2 bc_{e(k)}/n^2$ to serve the examples in phase k .

LEMMA 4.6. *Let $B(k)$ denote the set of examples for which we observe the currently explored attribute during phase k , and let $b(k) = |B(k)|$. Let $x_{e(k)}$ be the most expensive attribute which GL explores during phase k . Then, with probability at least $1 - \delta$, for every phase k , the cost of OPT_D during phase k is at least $(1 - \epsilon)^2 b(k)c_{e(k)}/n^2$.*

PROOF. Recall that \mathcal{P}^k denotes the prefix of the permutation of GL in phase k . Let D^k be the distribution D conditioned on examples for which all attributes in \mathcal{P}^k are zero. Consider using OPT_D on D^k . Let $x_{opt(k)}$ be the cheapest attribute that is one with probability at least $\frac{1}{n}$ over D^k .

By Lemma 4.4, with probability $1 - \delta$, for every k , the fraction of the examples in $B(k)$ that are answered by $x_{opt(k)}$ is at least $(1 - \epsilon)\frac{1}{n}$. So the cost of OPT_D , just on the examples in $B(k)$, is at least $b(k)(1 - \epsilon)c_{opt(k)}/n$.

If $c_{e(k)} \leq c_{opt(k)}$ the lemma follows. Otherwise we claim that $c_{e(k)} \leq nc_{opt(k)}/(1 - \epsilon)$. By the definition of our algorithm if $c_{e(k)} > c_{opt(k)}$ then $x_{opt(k)}$ is explored before $x_{e(k)}$. Since $x_{opt(k)}$ is true for a fraction $\geq 1/n$ of D^k then, by Lemma 4.4, when our algorithm explores it estimates its residual probability to be at least $(1 - \epsilon)/n$. Therefore the residual cost of $x_{opt(k)}$ after exploring it is at most $nc_{opt(k)}/(1 - \epsilon)$. It follows from the definition of our algorithm that $c_{e(k)} \leq nc_{opt(k)}/(1 - \epsilon)$ as otherwise we would not explore it. \square

Last we bound E^{suf} . Notice that for each example j , $1 \leq j \leq T$, such that $c^{suf}(e_j) > 0$, OPT_D must also observe an attribute in the suffix. This is since, by the definition of GL, all attributes in the prefix and the currently explored attribute at time j are zero. By Theorem 3.2, for each j such that $c^{suf}(e_j) > 0$, the cost of OPT_D to answer example j is at least $c^{suf}(e_j)/n$. This establishes the following lemma.

LEMMA 4.7. *The cost E^{suf} is at most n times the cost of OPT_D on S' . \square*

5. ONLINE DISTRIBUTION MODEL II: UNKNOWN TARGET FUNCTION

In this section, we assume that the target function f is a monotone clause which is not known to the algorithm, and that the distribution from which examples are drawn is a product distribution D , which is given to the algorithm. At each time step t , our online algorithm specifies a hypothesis h^t and a way to evaluate it on the input z^t (which is drawn from D). At the end of each step, the algorithm observes the correct value of the target function, i.e., $f(z^t)$. When evaluating the performance of an online algorithm, there are two different measures: the first measure is the number of mistakes, i.e., $|\{t \mid h^t(z^t) \neq f(z^t)\}|$, and the second measure is the average cost per example. (In our setting h^t will also be a monotone clause.)

There is a clear tradeoff between the two measures. An extreme trivial solution, is not to observe *any* attribute and output a random guess! The cost is zero but the error rate is high (half). In the other extreme, one can use the best learning algorithm for the class in question (monotone clauses, in our case) regardless of the cost of observing the examples. In this section, we aim at getting a tradeoff between the two extremes. Moreover, we would like our number of mistakes to be independent of the number of time steps (specifically it will be $O(n)$, where n is the number of variables, which is optimal) while keeping the cost near optimal.

Before describing our algorithm, we describe a standard procedure which is widely used in online learning of clauses. The procedure $\text{ELIM}(z^t, h^t)$ is performed when $f(z^t) = 0$ and $h^t(z^t) = 1$. In such a case $\text{ELIM}(z^t, h^t)$ returns the set of variables in the clause h^t which are TRUE, i.e., $\text{ELIM}(z^t, h^t) = \{x_i \in h^t \mid z_i^t = 1\}$ (and we say that $\text{ELIM}(z^t, h^t)$ are not in the target function f . This implies that we can “eliminate” any such variable. Clearly, the maximum number of times we can use ELIM is n .

We assume first that our online algorithm $\text{ONLINE}(q, D)$ receives as an input, in addition to the product distribution $D = (p_1, \dots, p_n)$, also the probability q such that $q \stackrel{\text{def}}{=} \Pr[f = 1]$. Let X^t be the set of variables that have not been eliminated before time t . (Initially X^1 includes all n variables.) To specify h^t , let $x_{i_1}, \dots, x_{i_\ell}$ be the variables in X^t sorted by increasing value of $\mu_{i_j} = c_{i_j}/p_{i_j}$ (break ties arbitrarily)⁵. Let $H_{k,r}^t \stackrel{\text{def}}{=} \bigvee_{j=1}^{k-1} x_{i_j} \vee (x_{i_k} \wedge \text{Br}(r))$, where $\text{Br}(r)$ is a Bernoulli random variable whose probability of 1 is r (and the probability of 0 is $1 - r$). The evaluation of $H_{k,r}^t$ is done by observing the variables x_{i_1} to $x_{i_{k-1}}$ in order. If any of these variables is 1 then $H_{k,r}^t$ terminates

⁵It is assumed, wlog, that $p_{i_j} > 0$; otherwise, the variable x_{i_j} can be ignored.

with output 1. If they are all zero then it flips a coin with bias r . If the coin comes up 1 it observes x_{i_k} and returns its value and otherwise it returns zero. The hypothesis h^t is the unique $H_{k,r}^t$ such that $\Pr[H_{k,r}^t = 1] = q$.

LEMMA 5.1. *For any t , there is a unique pair (k, r) such that $\Pr[H_{k,r}^t = 1] = q$.*

An online algorithm is *balanced* if, at every time step t , we have $\Pr[f = 1] = \Pr[h^t = 1]$. (Algorithm $\text{ONLINE}(q, D)$ is balanced by definition.) For a balanced online algorithm, given that at time t a mistake occurred, i.e., $f(z^t) \neq h^t(z^t)$, the probability that $\text{ELIM}(z^t, h^t)$ is applicable is exactly $1/2$. This implies that the expected number of mistakes is $O(n)$.

CLAIM 5.2. *The expected number of mistakes of a balanced online algorithm is at most $O(n)$.*

Let $\text{OPT}(f, D)$ be the optimal evaluation of f , with respect to a distribution D . (Recall that, by Lemma 2.1, when D is a product distribution then the optimal evaluation sorts the variables in f according to $\mu_i = c_i/p_i$.) The next lemma bounds the expected cost of $\text{ONLINE}(q, D)$ at any time t .

LEMMA 5.3. *Let D be a product distribution, and let $h^t = H_{k,r}^t$ be such that $\Pr[f = 1] = \Pr[h^t = 1]$. Then, $E[\text{cost}(\text{OPT}(f, D))] \geq E[\text{cost}(h^t)]$.*

PROOF. To simplify notation, assume that the variables $x_{i_1}, \dots, x_{i_\ell}$ in X^t when sorted by increasing μ_{i_j} , are in fact x_1, \dots, x_ℓ . That is $h^t = H_{k,r}^t = \bigvee_{j=1}^{k-1} x_j \vee (x_k \wedge \text{Br}(r))$. We can write the expected cost of h^t as,

$$\begin{aligned} E[\text{cost}(h^t)] &= \sum_{i=1}^{k-1} c_i \prod_{j=1}^{i-1} (1 - p_j) + r c_k \prod_{j=1}^{k-1} (1 - p_j) \\ &= \sum_{i=1}^{k-1} \mu_i p_i \prod_{j=1}^{i-1} (1 - p_j) + r \mu_k p_k \prod_{j=1}^{k-1} (1 - p_j). \end{aligned}$$

We can visualize $E[\text{cost}(h^t)]$ as an integral of a step function $\eta(x)$, whose values are μ_i . More precisely, let $q_0 = 0$, and let $q_m = \sum_{i=1}^m p_i \prod_{j=1}^{i-1} (1 - p_j)$, for $1 \leq m \leq k$. Note that q_m is the probability that one of the variables x_1, \dots, x_m is one. For $x \in (q_{i-1}, q_i]$, where $1 \leq i \leq k-1$, we have $\eta(x) = \mu_i$; for $x \in (q_{k-1}, q_{k-1} + r(q_k - q_{k-1})]$, we have $\eta(x) = \mu_k$; and for $x \in (q_{k-1} + r(q_k - q_{k-1}), 1]$, we have $\eta(x) = 0$. This implies that $\int_0^1 \eta(x) dx = E[\text{cost}(h^t)]$.

Similarly, we can write the expected value of the optimal evaluation of f as an integral of $\eta_f(x)$, which is a step function defined as $\eta(x)$ but with respect to $\text{OPT}(f, D)$. Since $\Pr[h^t = 1] = \Pr[f = 1]$, we have $\eta_f(x) \neq 0$ for any x such that $\eta(x) \neq 0$. Furthermore, since X^t contains all the variables in f , and both h^t and the optimal evaluation of f order their variables by increasing value of μ_i , it follows that for every x such that $\eta(x) \neq 0$ we have $\eta(x) \leq \eta_f(x)$. Thus, $\eta(x) \leq \eta_f(x)$ for every x and $E[\text{cost}(h^t)] \leq E[\text{cost}(\text{OPT}(f, D))]$. \square

The following theorem now follows from Claim 5.2 and Lemma 5.3.

THEOREM 5.4. *Let D be a product distribution. The expected number of mistakes of Algorithm $\text{ONLINE}(q, D)$ is $O(n)$ and at any time t we have $E(\text{cost}(h^t)) \leq E[\text{cost}(\text{OPT}(f, D))]$, where h^t is the hypothesis of $\text{ONLINE}(q, D)$ at time t (and its evaluation order). (The expectation is over an example drawn from the product distribution D .) \square*

Now, we extend the result to the case where $q = \Pr[f = 1]$, the probability that the target function is one, is unknown to the algorithm. We modify the algorithm and its analysis and show that, the same goal can be achieved with a slight increase in the cost, and $O(n \log n)$ mistakes.

The algorithm $\text{ON_APPROX}(\bar{\tau}, \delta, D)$ receives as a parameter $\bar{\tau} \in (0, 1/4]$, a confidence parameter $0 < \delta < 1$, and a product distribution D . The algorithm will try to learn a monotone clause f from examples drawn from D . The parameter $\bar{\tau}$ controls the tradeoff between the number of mistakes that $\text{ON_APPROX}(\bar{\tau}, \delta, D)$ makes to its expected cost. For larger $\bar{\tau}$ we will make less mistakes but the cost would be larger relative to the optimal cost. The parameter δ is the probability that $\text{ON_APPROX}(\bar{\tau}, \delta, D)$ would not behave as we expect (“fail”). Specifically, with probability $1 - \delta$ $\text{ON_APPROX}(\bar{\tau}, \delta, D)$ will indeed converge to the function f if given enough examples from D , while making at most $O(\frac{n}{\bar{\tau}} \log(\frac{n}{\delta}))$ mistakes with expected cost which is larger than the optimal expected cost by at most a factor of $(1 + \bar{\tau})$.

The algorithm $\text{ON_APPROX}(\bar{\tau}, \delta, D)$ runs in phases and maintains a set X of variables that have not been eliminated. At phase k , we have a hypothesis h^k which is used throughout this phase, and the evaluation of h^k is done by sorting the variables that it contains according to $\mu_i = c_i/p_i$. We start the first phase with $h^1 = \emptyset$. Phase k terminates in one of two ways:

(1) There is a mistake on an example x where $f(x) = 0$ and $h^k(x) = 1$. We call such a mistake an *eliminating mistake* since $\text{ELIM}(x, h^k)$ will eliminate at least one variable from h^k . We call a phase that terminates due to an eliminating mistake an *eliminating phase*.

When an eliminating phase k terminates with an eliminating mistake x then we set $X = X \setminus \text{ELIM}(x, h^k)$ and $h^{k+1} = h^k \cap X$.

(2) During phase k , there has been $\alpha = (\gamma/\epsilon^2) \log 4n/\delta$ mistakes where $f(x) = 1$ and $h^k(x) = 0$ (where $\epsilon = \bar{\tau}/16$ and $\gamma \geq (1 + \epsilon)$ is a fixed constant). We call such a phase a *growing phase*.

We then proceed to phase $k + 1$. Consider a growing phase k that terminates. Let L_k be the total number of steps in phase k and $\hat{q}_k = \alpha/L_k$. Let $Z_k = \{x_i \in X \mid p_i \leq (1 + \epsilon)\hat{q}_k / \Pr[h = 0]\} \setminus h^k$. If Z_k is empty then the algorithm continues with $h^{k+1} = h^k$ (we will show that this occurs with negligible probability). If Z_k is not empty then let $x_{i_{k+1}} = \arg \min_{x_i \in Z_k} \mu_i = c_i/p_i$. We set $h^{k+1} = h^k \vee x_{i_{k+1}}$ and proceed to phase $k + 1$.

Clearly there are at most n eliminating phases since, after each such phase, we eliminate at least one variable that cannot get back into h . We say that a growing phase k *succeeds* if $Z_k \neq \emptyset$ and thereby $h^k \subsetneq h^{k+1}$. Clearly, there are at most n growing phases that succeed. We will show that with probability $1 - \delta$ all growing phases succeed and therefore there are at most $2n$ phases altogether (Theorem 5.10).

Once we establish that with probability $1 - \delta$ there are only $2n$ phases, by the definition of the algorithm, it makes at most $O(n\alpha) = O(\frac{n}{\bar{\tau}} \log(\frac{n}{\delta}))$ mistakes.

We focus now on a particular phase k . Denote by M^k the event that we obtain a mistake; that is, the event that for an example x drawn from D , $f(x) \neq h^k(x)$. Denote by M_e^k the event that $f(x) = 0$ and $h^k(x) = 1$, in which case we say that x is an eliminating example (recall that if this event happens then the phase ends with an elimination). Denote by M_{ne}^k the event that $f(x) = 1$ and $h^k(x) = 0$. When

we focus on a particular phase k , we sometimes drop the superscript k when it is clear from the context and use M , M_e , and M_{ne} to denote the three events.

Our next lemma shows that if there is a non-negligible probability of getting an eliminating mistake in a phase, then indeed the phase is likely to end with an elimination.

LEMMA 5.5. *If $\Pr[M_e^k \mid M^k] \geq \epsilon^2$ then with probability at least $1 - \delta/4n$ phase k ends with an elimination. \square*

We say that phase k is *heavy* if $\Pr[M_e^k \mid M^k] \geq \epsilon^2$; otherwise, phase k is *light*. Lemma 5.5 asserts that a heavy phase is likely to be eliminating. Specifically,

LEMMA 5.6. *With probability $1 - \delta/2$, all heavy phases, in the first $2n$ phases, are eliminating phases. \square*

We focus now on a light growing phase and show that it is likely to succeed. Since in a light phase $\Pr[M_e \mid M] \leq \epsilon^2$, it follows that the probability of a mistake (i.e., the event M) in such a phase is similar to the probability of the event M_{ne} .

Specifically, let k be a light phase and let $p = \Pr[M^k]$, $p_{ne} = \Pr[M_{ne}^k]$, and $p_e = \Pr[M_e^k]$. Since the phase is light, we have $(1 - \epsilon^2)p \leq p_{ne} \leq p$ and $p_e \leq \epsilon^2 p \leq \epsilon^2$. Let $q = \Pr[M_{ne} \mid \overline{M_e}]$, where $\overline{M_e}$ is the event that the example is not eliminating. Since $M_{ne} \subseteq \overline{M_e}$, we get $q = \Pr[M_{ne}] / \Pr[\overline{M_e}] = p_{ne} / (1 - p_e)$. Combining the last inequalities,

$$(1 - \epsilon^2)p \leq p_{ne} \leq q = \frac{p_{ne}}{1 - p_e} \leq \frac{p_{ne}}{(1 - \epsilon^2)} \leq \frac{p}{(1 - \epsilon^2)}, \quad (1)$$

and so

$$(1 - \epsilon^2)q \leq p \leq \frac{q}{1 - \epsilon^2} \leq (1 + 2\epsilon^2)q. \quad (2)$$

Recall that $\hat{q}_k = \alpha/L_k$, where L_k is the length of the phase. The following lemma shows that q (and therefore also p and p_{ne}) are close to \hat{q}_k with high probability.

LEMMA 5.7. *Let k be a growing phase then with probability $1 - \delta/2n$, $|q - \hat{q}_k| \leq \epsilon \hat{q}_k$. \square*

Equations (1) and (2), and Lemma 5.7 imply the following corollary

COROLLARY 5.8. *Let k be a light growing phase. Then, with probability $1 - \delta/2n$, for $\epsilon \leq 1/4$, $|p - \hat{q}_k| \leq (1 - 2\epsilon)$ and $|p_{ne} - \hat{q}_k| \leq (1 - 2\epsilon)$.*

The next lemma guarantees that, with probability $1 - \delta/2n$, each light growing phase indeed succeeds.

LEMMA 5.9. *With probability $1 - \frac{\delta}{2n}$, in every light growing phase k we have $V \subset Z_k$, where V is the set of variables in f and not in h^k .*

PROOF. Let g be the clause of the variables in V and $x_i \in V$. We have that

$$\begin{aligned} p_{ne} &= \Pr[f(x) = 1, h^k(x) = 0] = \Pr[g(x) = 1, h^k(x) = 0] \\ &= \Pr[g(x) = 1] \Pr[h^k(x) = 0] \geq p_i \Pr[h^k(x) = 0]. \end{aligned}$$

By Corollary 5.8, $p_{ne} \leq (1 + \epsilon)\hat{q}_k$, with probability $1 - \frac{\delta}{2n}$. (Note that this holds for a phase, independently of any particular $x_i \in V$.) So, for every $x_i \in V$, we have

$$p_i \leq (1 + \epsilon)\hat{q}_k / \Pr[h^k(x) = 0],$$

with probability $1 - \frac{\delta}{2n}$, and therefore $x_i \in Z_k$, with probability $1 - \frac{\delta}{2n}$. \square

The following theorem states that indeed, with probability $1 - \delta$, our algorithm has at most $2n$ phases.

THEOREM 5.10. *For any product distribution D , algorithm $\text{ON_APPROX}(\bar{\epsilon}, \delta, D)$, with probability at least $1 - \delta$, terminates after at most $2n$ phases. \square*

Last we prove a bound on the cost of the algorithm.

LEMMA 5.11. *With probability $1 - \delta$, for any phase k , the expected cost of h^k is at most $(1 + 16\epsilon)\text{OPT}(f, D)$.*

PROOF. The statement clearly holds for h^1 . We prove that, at the end of a light growing phase k , the statement holds for h^{k+1} with probability $1 - \delta/2n$.⁶ Since after a terminating phase the cost only decreases, and with probability $1 - \delta/2$ we have only terminating phases and light growing phases the lemma would follow.

Note that $p_e = \Pr[M_e] = \Pr[M_e|M] \Pr[M]$. Since the phase is light then, by Corollary 5.8, with probability $1 - \delta/2n$, we have $p = \Pr[M] \leq (1 + 2\epsilon)\hat{q}_k$, and $p_{ne} \geq (1 - 2\epsilon)\hat{q}_k$. So $p_e = \Pr[M_e|M] \Pr[M] \leq \epsilon^2(1 + 2\epsilon)\hat{q}_k \leq \epsilon\hat{q}_k$, for $\epsilon \leq 1/4$. Therefore,

$$\begin{aligned} \Pr[f(x) = 1] - \Pr[h^k(x) = 1] &= p_{ne} - p_e \\ &\geq ((1 - 2\epsilon) - \epsilon)\hat{q}_k \geq (1 - 3\epsilon)\hat{q}_k. \end{aligned}$$

Since we select $x_{i_{k+1}}$ such that $\Pr[x_{i_{k+1}} = 1] \leq \frac{(1+\epsilon)\hat{q}_k}{\Pr[h^k(x)=0]}$, we have $\Pr[h^{k+1}(x) = 1] \leq \Pr[h^k(x) = 1] + (1 + \epsilon)\hat{q}_k$. Consider

$$h' = \frac{1 - 3\epsilon}{1 + \epsilon}h^{k+1} + \frac{4\epsilon}{1 + \epsilon}h^k.$$

For h' we have that $\Pr[h'(x) = 1] \leq \Pr[f(x) = 1]$. Consider the set of variables $h^k \cup Z_k$, then by Lemma 5.9 both h' and f are subsets of it, and h' is in fact $H_{k+1, \frac{1-3\epsilon}{1+\epsilon}}$ on that set. Therefore, by an argument as in the proof of Lemma 5.3, we have that $E[\text{cost}(h')] \leq E[\text{cost}(f)]$. On the other hand, $E[\text{cost}(h')] \geq \frac{1-3\epsilon}{1+\epsilon}E[\text{cost}(h^{k+1})]$. Therefore, $E[\text{cost}(h^{k+1})] \leq (1 + 16\epsilon)E[\text{cost}(f)]$, for $\epsilon \in (0, 1/4)$. \square

The following theorem summarizes this section.

THEOREM 5.12. *For any product distribution D and any clause f , with probability $1 - \delta$, the number of mistakes of $\text{ON_APPROX}(\bar{\epsilon}, \delta, D)$ is at most $O(n\alpha) = O(\frac{n}{\bar{\epsilon}} \log(n/\delta))$ and, at any time t we have $E[\text{cost}(h^t)] \leq (1 + \bar{\epsilon})E[\text{OPT}(f, D)]$.*

6. REFERENCES

- [1] D. Angluin and L. G. Valiant. Fast probabilistic algorithms for Hamiltonian circuits and matchings. *JCSS*, 18(2):155–193, April 1979.
- [2] S. Babu, R. Motwani, K. Munagala, I. Nishizawa, and J. Widom. Adaptive ordering of pipelined stream filters. In *SIGMOD*, 2004.
- [3] A. Bar-Noy, M. Bellare, M. M. Halldórsson, H. Shachnai, and T. Tamir. On chromatic sums and distributed resource allocation. *Inf. Comput.*, 140(2):183–202, 1998.
- [4] A. Bar-Noy, M. M. Halldórsson, and G. Kortsarz. A matched approximation bound for the sum of a greedy coloring. *IPL*, 71(3-4):135–140, 1999.

⁶The dependency is on the same event that guarantees that the phase succeed.

- [5] M. Charikar, R. Fagin, V. Guruswami, J. Kleinberg, P. Raghavan, and A. Sahai. Query strategies for priced information. In *32nd STOC*, pages 582–591, 2000.
- [6] E. Cohen, A. Fiat, and H. Kaplan. Associative search in peer to peer networks: Harnessing latent semantics. In *Proceedings IEEE INFOCOM*, 2003.
- [7] O. Etzioni, S. Hanks, T. Jiang, R. M. Karp, O. Madani, and O. Waarts. Efficient information gathering on the internet. In *FOCS*, 234–243, 1996.
- [8] O. Etzioni, S. Hanks, T. Jiang, and O. Madani. Optimal information gathering on the internet with time and cost constraints. *SIAM J. Comput.*, 29(5):1596–1620, 2000.
- [9] U. Feige, L. Lovász, and P. Tetali. Approximating min-sum set cover. In *APPROX*, 94–107, 2002.
- [10] J. M. Hellerstein and M. Stonebraker. Predicate migration: optimizing queries with expensive predicates. In *SIGMOD*, 267–276. ACM Press, 1993.
- [11] A. T. Kalai and S. Vempala. Efficient algorithms for online decision problems. In *COLT*, 26–40, 2003.
- [12] N. Littlestone. Learning when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, 2:285–318, 1988.
- [13] D. J. Lizotte, O. Madani, and R. Greiner. Budgeted learning of naive-bayes classifiers. In *UAI*, 378–385, 2003.
- [14] O. Madani, D. J. Lizotte, and R. Greiner. Active model selection. *ICML*, 2004.
- [15] K. Munagala, S. Babu, R. Motwani, and J. Widom. The pipelined set cover problem. *ICDT*, 2005.
- [16] L. G. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.

Appendix: Concentration Bounds

For our analysis in Section 4 and 5 we need the following version of the Chernoff bound [1].

THEOREM 6.1. *Let X_1, \dots, X_m be i.i.d. $\{0, 1\}$ random variables, where $E[X_i] = p$. Let $b = \sum_{i=1}^m X_i$. Then,*

1. For all $\beta \in (0, 1]$ we have $\Pr[b \leq (1 - \beta)pm] \leq e^{-\beta^2 pm/2}$
2. For all $\beta \in (0, 1]$ we have $\Pr[b \geq (1 + \beta)pm] \leq e^{-\beta^2 pm/3}$
3. For all $\beta \geq 1$ we have $\Pr[b \geq (1 + \beta)pm] \leq e^{-\beta^2 pm/(2 + \beta)}$
4. (combining the above 2 and 3) For all $\beta > 0$ we have $\Pr[b \geq (1 + \beta)pm] \leq e^{-\beta pm/3}$

An immediate consequence we obtain the following lemma.

LEMMA 6.2. *For $q \in [0, 1]$ and $F \geq 0$ let $m = F/q$. Let X_1, \dots, X_m be i.i.d. $\{0, 1\}$ random variables, where $E[X_i] = p$. Let $b = \sum_{i=1}^m X_i$ and $\hat{p} = b/m$. Then, for any $\epsilon \in (0, 1)$, with probability $1 - 4e^{-\epsilon^2 F/6}$,*

1. If $p \geq 2q$ then $b \geq F$.
2. If $p \leq q/2$ then $b < F$.
3. If $p \geq q/2$ then $\hat{p} \in [(1 - \epsilon)p, (1 + \epsilon)p]$.