

Margin-Based Ranking Meets Boosting in the Middle^{*}

Cynthia Rudin¹, Corinna Cortes², Mehryar Mohri³, and Robert E. Schapire⁴

¹ Howard Hughes Medical Institute, New York University
⁴ Washington Place, Room 809, New York, NY 10003-6603
rudin@nyu.edu

² Google Research, 1440 Broadway, New York, NY 10018
corinna@google.com

³ Courant Institute, New York University, 719 Broadway, New York, NY 10003
mohri@cs.nyu.edu

⁴ Princeton University, Department of Computer Science
35 Olden St., Princeton NJ 08544
schapire@cs.princeton.edu

Abstract. We present a margin-based bound for ranking in a general setting, using the L_∞ covering number of the hypothesis space as our complexity measure. Our bound suggests that ranking algorithms that maximize the ranking margin will generalize well.

We produce a Smooth Margin Ranking algorithm, which is a modification of RankBoost analogous to Approximate Coordinate Ascent Boosting. We prove that this algorithm makes progress with respect to the ranking margin at every iteration and converges to a maximum margin solution. In the special case of bipartite ranking, the objective function of RankBoost is related to an exponentiated version of the AUC. In the empirical studies of Cortes and Mohri, and Caruana and Niculescu-Mizil, it has been observed that AdaBoost tends to maximize the AUC. In this paper, we give natural conditions such that AdaBoost maximizes the exponentiated loss associated with the AUC, i.e., conditions when AdaBoost and RankBoost will produce the same result, explaining the empirical observations.

1 Introduction

Consider the following supervised learning problem: Sylvia would like to get some recommendations for good movies before she goes to the theater. She would like a ranked list that agrees with her tastes as closely as possible, since she will see the movie closest to the top of the list that is playing at the local theater. For many pairs of movies she has seen, she will tell the learning algorithm whether she likes the first movie better than the second. This allows her to rank whichever pairs of movies she wishes, allowing for the possibility of ties between movies, and the possibility that certain movies cannot necessarily be compared by her, e.g, she

^{*} This work is partially supported by NSF grant CCR-0325463.

may not compare action movies with cartoons. Another advantage of this type of scoring over real-valued scoring is that Sylvia does not need to normalize her own scores in order to compare with the rankings of another person; she just compares rankings on pairs of movies. Sylvia does not need to be consistent, since she may rank $a > b > c > a$. Each pair of movies such that Sylvia ranks the first above the second is called a “crucial pair”.

The learning algorithm has access to a set of n individuals (“weak rankers”, or “ranking features”) who have also ranked pairs of movies. The learning algorithm must combine the views of the weak rankers in order to match Sylvia’s preferences, and generate a recommendation list that will generalize her views. This type of problem was studied in depth by Freund et al. [7], where the Rank-Boost algorithm was introduced.

In order to give some indication that an algorithm will generalize well (e.g., we want the ranking algorithm to predict movies that Sylvia will like), one often considers generalization bounds. Generalization bounds show that a small probability of error will most likely be achieved through a balance of the empirical error and the complexity of the hypothesis space. This complexity can be measured by an informative quantity, such as the VC dimension, covering number, Rademacher complexity, or a more specialized quantity, such as the bipartite rank shatter coefficient [1], which was used to derive a generalization bound specifically for the case of bipartite ranking. The “bipartite” ranking problem is a special case of the ranking problem where there are only two classes, a positive class, i.e., “good movies”, and a negative class, i.e., “bad movies”.

When deriving generalization bounds, it is important to consider the “separable” case, where all training instances are correctly handled by the learning algorithm so the empirical error is zero. Ranking in the separable case means that the algorithm’s chosen ranking is consistent with all crucial pairs; the algorithm ranks the first instance in each crucial pair above the second. In the bipartite ranking problem, the separable case means that all positive instances are ranked above all negative instances, and the AUC (Area Under the ROC Curve) is exactly 1.

In the separable case for classification, one important indicator of a classifier’s generalization ability is the “margin”, e.g., for boosting [15] and support vector machines. Although the empirical success of an algorithm depends on many factors (e.g., the type of data and how noisy it is), margin-based bounds often do provide a reasonable explanation (though not a complete explanation) of the success of many algorithms, both empirically and theoretically. Although there has been some work devoted to generalization bounds for ranking [7, 1], the bounds (that we are aware of) are not margin-based, and thus do not provide this useful type of discrimination between ranking algorithms in the separable case.

In Section 3, we provide a margin-based bound for ranking in a general setting. Our bound uses the L_∞ covering number as the complexity measure for the hypothesis space.

Since we have provided a general margin-based bound for ranking, we now try to derive algorithms which create large margins. For the classification problem, it was proved that AdaBoost does not always maximize the margin [12]. In fact, AdaBoost does not even necessarily make progress towards increasing the margin at every iteration. In analogy, RankBoost does not directly maximize the ranking margin, and it may not increase the margin at every iteration. In Section 4.1 we introduce a Smooth Margin Ranking algorithm, and prove that it makes progress towards increasing the “smooth” ranking margin at every iteration; this is the main step needed to prove convergence and convergence rates. This algorithm is analogous to Approximate Coordinate Ascent Boosting [14, 13] in its derivation, but the analogous proof that progress occurs at each iteration is much trickier; hence we present a sketch of this proof here, along with a proof that this algorithm converges to a maximum margin solution.

In the bipartite ranking problem, we want our recommendation list to minimize the misranking error, e.g., the probability that a bad movie is ranked above a good movie. The empirical version of this misranking error is closely related to the AUC. RankBoost [7] minimizes an exponentiated version of this misranking error, in analogy with the classification algorithm AdaBoost, which minimizes an exponentiated version of the margins of training instances. The bipartite ranking setting and the basic classification setting are different, since they use different loss functions. For most classification loss functions (e.g. in the SVM loss), a misclassified instance is given the same loss no matter how badly it is misclassified. In the ranking problem, if a negative instance is ranked more highly than only one positive instance, it hardly affects the loss, whereas if the negative instance is classified more highly than many positive instances, it affects the loss greatly.

Although AdaBoost and RankBoost were derived analogously for the two settings, the parallels between these algorithms are deeper than their derivations. Cortes and Mohri [5] and Caruana and Niculescu-Mizil [3] have noted that AdaBoost experimentally seems to be very good at the bipartite ranking problem, even though it was RankBoost that was explicitly designed to solve this problem, not AdaBoost. That is, AdaBoost often achieves a large AUC. In Section 5, we show an important reason for these observations. Namely, we show that if positive and negative instances contribute equally to AdaBoost’s final loss, then AdaBoost and RankBoost produce the same solution. We define a quantity called “F-skew”, an exponentiated version of the “skew” used in the expressions of Cortes and Mohri [5] and Agarwal et al. [1]. When the F-skew is zero, AdaBoost’s loss comes equally from positive and negative instances. We prove that if the F-skew vanishes, AdaBoost minimizes the exponentiated ranking loss, which is the same loss that RankBoost explicitly minimizes; the two algorithms will produce the same solution. Moreover, we show that if AdaBoost’s set of weak classifiers includes the constant classifier, then the F-skew always vanishes.

We proceed from the most general to the most specific. In Section 3 we provide a margin based bound for general ranking, which holds for each element of the hypothesis space. In Sections 4.1 and 4.2 we focus on the binary ranking

problem, where we fix the form of hypothesis space to match that of RankBoost. Here, we discuss coordinate-based ranking algorithms such as RankBoost, and introduce the Smooth Margin Ranking algorithm. In Section 5, we consider the bipartite ranking problem, which is a special case of the binary ranking problem. Here, we discuss conditions for AdaBoost to act as a bipartite ranking algorithm by minimizing the exponentiated loss associated with the AUC. Sections 3 and 4.2 focus on the separable case, and Sections 4.1 and 5 focus on the non-separable case.

The main contributions of this paper are: 1) a margin-based ranking bound, 2) a theorem stating that our Smooth Margin Ranking algorithm makes progress at every iteration towards increasing the smooth ranking margin, and 3) a set of conditions when AdaBoost acts as a bipartite ranking algorithm.

2 Notation

The training set, denoted by S , is $\{\mathbf{x}_i\}_{i=1,\dots,m}$, where $\mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^N$. The set \mathcal{X} may be finite or infinite. In the case of the movie ranking problem, the \mathbf{x}_i 's are the movies and \mathcal{X} is the database. The instances $\mathbf{x}_i \in \mathcal{X}$ are chosen independently and at random (iid) from a fixed but unknown probability distribution \mathcal{D} on \mathcal{X} . The notation $\mathbf{x} \sim \mathcal{D}$ means \mathbf{x} is chosen randomly according to \mathcal{D} , and $S \sim \mathcal{D}^m$ means the m elements of the training set S are chosen iid according to \mathcal{D} .

The values of the ‘‘truth’’ function $\pi : \mathcal{X} \times \mathcal{X} \rightarrow \{0, 1\}$, which is defined over pairs of instances, are analogous to the ‘‘labels’’ in classification. If $\pi(\bar{\mathbf{x}}, \tilde{\mathbf{x}}) = 1$, the pair $\bar{\mathbf{x}}, \tilde{\mathbf{x}}$ is a crucial pair, i.e., $\bar{\mathbf{x}}$ should be ranked more highly than $\tilde{\mathbf{x}}$. We require only that $\pi(\bar{\mathbf{x}}, \bar{\mathbf{x}}) = 0$, meaning $\bar{\mathbf{x}}$ cannot be ranked higher than itself, and also $\pi(\bar{\mathbf{x}}, \tilde{\mathbf{x}}) = 1$ implies $\pi(\tilde{\mathbf{x}}, \bar{\mathbf{x}}) = 0$, meaning that if $\bar{\mathbf{x}}$ is ranked higher than $\tilde{\mathbf{x}}$, that $\tilde{\mathbf{x}}$ cannot be ranked higher than $\bar{\mathbf{x}}$. (It is possible that these assumptions may be dropped.) It is possible to have $\pi(a, b) = 1$, $\pi(b, c) = 1$, and $\pi(c, a) = 1$; this forces us to be in the non-separable case. The quantity $E := \mathbb{E}_{\bar{\mathbf{x}}, \tilde{\mathbf{x}} \sim \mathcal{D}} \pi(\bar{\mathbf{x}}, \tilde{\mathbf{x}})$ is the expected fraction of pairs in the database that are crucial pairs, $0 \leq E \leq 1/2$. We assume that π is a deterministic (non-noisy) function, and that for each pair of training instances $\mathbf{x}_i, \mathbf{x}_k$, we receive $\pi(\mathbf{x}_i, \mathbf{x}_k)$.

Our goal is to construct a ranking function $f : \mathcal{X} \rightarrow \mathbb{R}$, which gives a real valued rank to each instance in \mathcal{X} . We do not care about the actual values of each instance, only the relative values; for crucial pair $\bar{\mathbf{x}}, \tilde{\mathbf{x}}$, we do not care if $f(\bar{\mathbf{x}}) = .4$ and $f(\tilde{\mathbf{x}}) = .1$, only that $f(\bar{\mathbf{x}}) > f(\tilde{\mathbf{x}})$. Also, $f \in L_\infty(\mathcal{X})$ (or if $|\mathcal{X}|$ is finite, $f \in \ell_\infty(\mathcal{X})$).

In the binary setting of boosting for classification, $|f(\mathbf{x})| \leq 1 \forall \mathbf{x}$ and the *margin of training instance i* (with respect to classifier f) is $y_i f(\mathbf{x}_i)$, where y_i is the classification label, $y_i \in \{-1, 1\}$ [15]. The *margin of classifier f* is the minimum margin over all training instances, $\min_i y_i f(\mathbf{x}_i)$. Intuitively, the margin tells us how much f can change before one of the training instances is misclassified; it gives us a notion of how stable the classifier is.

For the ranking setting, we define an analogous notion of margin. Here, we can normalize f so that $|f| \leq 1$. The *margin of crucial pair i, k* with respect

to ranking function f will be defined as $f(\mathbf{x}_i) - f(\mathbf{x}_k)$. The *margin of ranking function* f , is the minimum margin over all crucial pairs,

$$\mu_f := \min_{\{i,k \mid \pi(\mathbf{x}_i, \mathbf{x}_k)=1\}} f(\mathbf{x}_i) - f(\mathbf{x}_k).$$

Intuitively, the margin tells us how much the ranking function can change before one of the crucial pairs is misranked. As with classification, we are in the separable case whenever the margin of f is positive.

3 A Margin-Based Bound for Ranking

In this section, we provide a bound which gives us an intuition for separable-case ranking and yields theoretical encouragement for margin-based ranking algorithms. The quantity we hope to minimize is analogous to the misclassification probability for classification; for two randomly chosen instances, if they are a crucial pair, we want to minimize the probability that these instances will be misranked. That is, we want to minimize:

$$\mathbb{P}_D\{\text{misrank}_f\} := \mathbb{P}_D[f(\tilde{\mathbf{x}}) \leq f(\tilde{\mathbf{x}}) \mid \pi(\tilde{\mathbf{x}}, \tilde{\mathbf{x}})=1] = \frac{\mathbb{E}_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}} \sim D} [\mathbf{1}_{[f(\tilde{\mathbf{x}}) \leq f(\tilde{\mathbf{x}})]} \pi(\tilde{\mathbf{x}}, \tilde{\mathbf{x}})]}{E}. \quad (1)$$

The numerator of (1) is the fraction of crucial pairs that are correctly ranked by f , and the denominator, $E := \mathbb{E}_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}} \sim D} \pi(\tilde{\mathbf{x}}, \tilde{\mathbf{x}})$ is the fraction of pairs that are crucial pairs. Thus, $\mathbb{P}_D\{\text{misrank}_f\}$ is the proportion of crucial pairs that are correctly ranked by f .

Since we do not know D , we may use only empirical quantities that rely on our training sample. An empirical quantity analogous to $\mathbb{P}_D\{\text{misrank}_f\}$ is:

$$\mathbb{P}_S\{\text{misrank}_f\} := \mathbb{P}_S\{\text{margin}_f \leq 0\} := \mathbb{E}_S[\mathbf{1}_{(f(\mathbf{x}_i) \leq f(\mathbf{x}_k))} \pi(\mathbf{x}_i, \mathbf{x}_k)].$$

We make this definition more general, by allowing it to include a margin of $\theta \geq 0$:

$$\begin{aligned} \mathbb{P}_S\{\text{margin}_f \leq \theta\} &:= \mathbb{P}_S[\mathbf{1}_{(f(\mathbf{x}_i) - f(\mathbf{x}_k) \leq \theta)} \mid \pi(\mathbf{x}_i, \mathbf{x}_k) = 1] \\ &= \frac{\sum_{i=1}^m \sum_{k=1}^m [\mathbf{1}_{(f(\mathbf{x}_i) - f(\mathbf{x}_k) \leq \theta)} \pi(\mathbf{x}_i, \mathbf{x}_k)]}{\sum_{i=1}^m \sum_{k=1}^m \pi(\mathbf{x}_i, \mathbf{x}_k)}, \end{aligned}$$

i.e., $\mathbb{P}_S\{\text{margin}_f \leq \theta\}$ is the fraction of crucial pairs with margin no larger than θ .

We want to bound $\mathbb{P}_D\{\text{misrank}_f\}$ in terms of an empirical, margin-based term and a complexity term. The type of complexity we choose is the L_∞ covering number of the hypothesis space \mathcal{F} , $\mathcal{F} \subset L_\infty(\mathcal{X})$. (Upper bounds on the covering number can be calculated; see [6]). The covering number $\mathcal{N}(\mathcal{F}, \sigma)$ is defined as the minimum number of balls of radius σ needed to cover \mathcal{F} , using the L_∞ metric. Now we state our theorem (see Appendix A for the proof):

Theorem 1. For $\epsilon > 0$, $\theta \geq 0$, for all $f \in \mathcal{F}$,

$$\mathbb{P}_{S \sim D^m} [\mathbb{P}_D\{\text{misrank}_f\} \leq \mathbb{P}_S\{\text{margin}_f \leq \theta\} + \epsilon] \geq 1 - \mathcal{N}\left(\mathcal{F}, \frac{\epsilon\theta}{4}\right) 4 \exp\left[-2m \left(\frac{\epsilon E}{4 + \epsilon}\right)^2\right].$$

That is, with probability depending on m , E , θ , ϵ , and \mathcal{F} , the misranking probability is less than the fraction of instances with margin above θ , plus ϵ .

We have chosen to write our bound in terms of E , but we could equally well have used an analogous empirical quantity, namely $\mathbb{E}_{\mathbf{x}_i, \mathbf{x}_k \sim \mathcal{S}} \pi(\mathbf{x}_i, \mathbf{x}_k) = \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{k=1}^m \pi(\mathbf{x}_i, \mathbf{x}_k)$. This is an arbitrary decision; we cannot maximize E in practice because the data is random. Either way, the bound tells us that the margin should be an important quantity to consider in the design of algorithms.

As a special case of the theorem, we consider the case of a finite hypothesis space \mathcal{F} , where the covering number achieves its largest value (for any θ), i.e., $\mathcal{N}(\mathcal{F}, \frac{\epsilon\theta}{4}) = |\mathcal{F}|$. Also take ϵ small so $4 + \epsilon \approx 4$. Now we can solve for ϵ :

$$\delta := |\mathcal{F}| 4 \exp \left[-2m \left(\frac{\epsilon E}{4} \right)^2 \right] \implies \epsilon = \frac{1}{\sqrt{m}} \sqrt{\frac{8}{E^2} (\ln 4 |\mathcal{F}| + \ln(1/\delta))}$$

This could be compared directly with Theorem 1 of Schapire et al [15]. Cucker and Smale [6] have reduced the factor ϵ^2 to a factor of ϵ in certain cases; it is possible that this bound may be tightened, especially in the case of a convex combination of weak rankers. An interesting open problem is to prove generalization bounds using Rademacher complexity or a more specialized bound analogous to those in Koltchinskii and Panchenko [10]; here the trick would be to find a symmetrization step analogous to the usual symmetrization for classification. In any case, our bound indicates that the margin is an important quantity for generalization.

4 Coordinate-Based Ranking Algorithms

We have just presented a uniform bound that holds for all $f \in \mathcal{F}$. Now we will discuss how a learning algorithm might pick one of those functions, in order to make $\mathbb{P}_D\{\text{misrank}_f\}$ as small as possible, based on intuition gained from the bound of Theorem 1; the bound reveals the margin to be a useful quantity in the learning process, so it deserves consideration in our design of algorithms.

In Section 4.1, we discuss RankBoost’s objective function \tilde{F} . Then, we describe a coordinate descent algorithm on this objective. In Section 4.2 we define the smooth ranking margin \tilde{G} . Then we present the Smooth Margin Ranking algorithm, and prove that it makes progress towards increasing \tilde{G} at each iteration and converges to a maximum margin solution.

4.1 Coordinate Descent on RankBoost’s Objective

We are in the binary setting of RankBoost, which we now describe. The hypothesis space \mathcal{F} is the class of convex combinations of “weak” rankers $\{h_j\}_{j=1, \dots, n}$, where $h_j : X \rightarrow \{0, 1\}$. We assume that if h_j is a weak ranker, that $1 - h_j$ is not chosen as a weak ranker; this assumption avoids the complicated seminorm notation in earlier work [13]. The function f is constructed as a normalized linear combination of the h_j ’s: $f = \sum_j \lambda_j h_j / \|\boldsymbol{\lambda}\|_1$, where $\|\boldsymbol{\lambda}\|_1 = \sum_j \lambda_j$.

We construct a structure \mathbf{M} , which tells us how each individual weak ranker j ranks each pair of instances i, k . We define \mathbf{M} element-wise as: $M_{ikj} := h_j(\mathbf{x}_i) - h_j(\mathbf{x}_k)$. Thus, $M_{ikj} \in \{-1, 0, 1\}$. Since \mathbf{M} has three indices, we need to define right multiplication: $(\mathbf{M}\boldsymbol{\lambda})_{ik} := \sum_{j=1}^n \lambda_j h_j(\mathbf{x}_i) - \lambda_j h_j(\mathbf{x}_k)$, and left multiplication: $(\mathbf{d}^T \mathbf{M})_j := \sum_{i=1}^m \sum_{k=1}^m d_{ik} M_{ikj}$.

Just as AdaBoost can be represented as a coordinate descent algorithm on a specific objective function of $\boldsymbol{\lambda}$ (see [9]), so can RankBoost. The objective function for RankBoost is:

$$\tilde{F}(\boldsymbol{\lambda}) := \sum_{\{i,k | [\pi(\mathbf{x}_i, \mathbf{x}_k)=1]\}} e^{-(\mathbf{M}\boldsymbol{\lambda})_{ik}}.$$

We now perform standard coordinate descent on \tilde{F} to derive ‘‘Coordinate Descent RankBoost’’. The direction chosen at iteration t (i.e., the choice of weak ranker j_t) in the ‘‘optimal’’ case (where the best weak ranker is chosen at each iteration) is given by: $j_t \in \underset{j}{\operatorname{argmax}} \left[-\frac{d\tilde{F}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_j)}{d\alpha} \Big|_{\alpha=0} \right] = \underset{j}{\operatorname{argmax}} (\mathbf{d}_t^T \mathbf{M})_j$, where the ‘‘weights’’ $d_{t,ik}$ are defined over pairs of instances by: $d_{t,ik} = 0$ for non-crucial pairs, and for crucial pair i, k : $d_{t,ik} := e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_{ik}} / \tilde{F}(\boldsymbol{\lambda}_t)$. One can see that the chosen weak ranker is a natural choice, namely, j_t is the most accurate weak ranker with respect to the weighted crucial training pairs.

Define $\mathcal{I}_+ := \{i, k | M_{ikj_t} = 1, \pi(\mathbf{x}_i, \mathbf{x}_k) = 1\}$ (although \mathcal{I}_+ is different for each t , we eliminate the subscript), also $\mathcal{I}_- := \{i, k | M_{ikj_t} = -1, \pi(\mathbf{x}_i, \mathbf{x}_k) = 1\}$. Also define $d_+ := \sum_{\mathcal{I}_+} d_{t,ik}$ and $d_- := \sum_{\mathcal{I}_-} d_{t,ik}$. The step size at iteration t is α_t , where α_t satisfies the equation for the line search along direction j_t :

$$0 = -\frac{d\tilde{F}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_{j_t})}{d\alpha} \Big|_{\alpha=\alpha_t} \Rightarrow \alpha_t = \frac{1}{2} \ln \frac{d_+}{d_-}. \quad (2)$$

Thus, we have derived the first algorithm, Coordinate Descent RankBoost.

RankBoost, as it is described by Freund et al. [7], is similar, but differs by the ordering of steps: α_t is calculated first (via (2)), and afterwards j_t is determined. For RankBoost, there may not be a natural interpretation of this type of weak learning algorithm as there is for Coordinate Descent RankBoost.

It is interesting that for AdaBoost’s objective function, the plain coordinate descent algorithm and the variation (choosing the coordinate with knowledge of the step size) actually turn out to both yield the same algorithm, i.e., AdaBoost.

4.2 Smooth Margin Ranking

The value of \tilde{F} does not directly tell us anything about the margin, only whether the margin is positive. Using any $\boldsymbol{\lambda}$ that yields a positive margin, we can actually make \tilde{F} arbitrarily small by multiplying $\boldsymbol{\lambda}$ by a large positive constant, so the objective is arbitrarily small, yet the margin may not be maximized. Actually, the same problem occurs for AdaBoost. It has been proven [12] that for certain \mathbf{M} , AdaBoost does not converge to a maximum margin solution, nor does it even make progress towards increasing the margin at every iteration. Since the

calculations are identical for RankBoost, there are certain cases in which we can expect RankBoost not to converge to a maximum margin solution.

In earlier work, we have proposed a smooth margin function which one can maximize in order to achieve a maximum margin solution for the classification problem [13]. We have proposed a coordinate ascent algorithm on this function which makes progress towards increasing the smooth margin at every iteration. Here, we present the analogous smooth ranking function and the Smooth Margin Ranking algorithm. The smooth ranking function \tilde{G} is defined as follows:

$$\tilde{G}(\boldsymbol{\lambda}) := \frac{-\ln \tilde{F}(\boldsymbol{\lambda})}{\|\boldsymbol{\lambda}\|}.$$

With proofs identical to those of Rudin et al. [13], one can show that:

$$\tilde{G}(\boldsymbol{\lambda}) < \mu(\boldsymbol{\lambda}) \leq \rho, \quad \text{where} \quad (3)$$

$\rho = \min_{\{\mathbf{d} \mid \sum_{ik} d_{ik} = 1, d_{ik} \geq 0\}} \max_j (\mathbf{d}^T \mathbf{M})_j = \max_{\{\bar{\boldsymbol{\lambda}} \mid \sum_j \bar{\lambda}_j = 1, \bar{\lambda}_j \geq 0\}} \min_i (\mathbf{M} \bar{\boldsymbol{\lambda}})_i$, i.e., the smooth ranking margin is less than the true margin, and the true margin is no greater than ρ , the min-max value of the game defined by \mathbf{M} (see [8]).

We now define the Smooth Margin Ranking algorithm, which is approximately coordinate ascent on \tilde{G} . As usual, the input to the algorithm is matrix \mathbf{M} . Also, we will only define this algorithm when $\tilde{G}(\boldsymbol{\lambda})$ is positive, so that we only use this algorithm once the data has become separable; we can use RankBoost or Coordinate Descent RankBoost to get us to this point.

We will define iteration $t + 1$ in terms of the quantities known at iteration t , namely, the following quantities can be calculated from $\boldsymbol{\lambda}_t$: the current value of the objective $g_t := \tilde{G}(\boldsymbol{\lambda}_t)$, the weights $d_{t,ik} := e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_{ik}} / \tilde{F}(\boldsymbol{\lambda}_t)$, the direction $\mathbf{j}_t = \underset{j}{\operatorname{argmax}} (\mathbf{d}_t^T \mathbf{M})_j$, and the edge $r_t := (\mathbf{d}_t^T \mathbf{M})_{\mathbf{j}_t}$. The choice of \mathbf{j}_t is the same as for Coordinate Descent RankBoost, also see [13]. The step size α_t is chosen to obey (6) below, but we need more definitions before we state its value. We define recursive equations for \tilde{F} and \tilde{G} , and then use these to build up to (6). We also have $s_t = \|\boldsymbol{\lambda}_t\|_1$ and $s_{t+1} = s_t + \alpha_t$, and $g_{t+1} = \tilde{G}(\boldsymbol{\lambda}_t + \alpha_t \mathbf{e}_{\mathbf{j}_t})$.

As before, $\mathcal{I}_+ := \{i, k \mid M_{ik\mathbf{j}_t} = 1, \pi(\mathbf{x}_i, \mathbf{x}_k) = 1\}$, $\mathcal{I}_- := \{i, k \mid M_{ik\mathbf{j}_t} = -1, \pi(\mathbf{x}_i, \mathbf{x}_k) = 1\}$, and now, $\mathcal{I}_0 := \{i, k \mid M_{ik\mathbf{j}_t} = 0, \pi(\mathbf{x}_i, \mathbf{x}_k) = 1\}$. Also $d_+ := \sum_{\mathcal{I}_+} d_{t,ik}$, $d_- := \sum_{\mathcal{I}_-} d_{t,ik}$, and $d_0 := \sum_{\mathcal{I}_0} d_{t,ik}$. So, by definition, $d_+ + d_- + d_0 = 1$. Now, r_t becomes $r_t = d_+ - d_-$. Define the factor τ_t and its ‘‘derivative’’ τ'_t :

$$\tau_t := d_+ e^{-\alpha_t} + d_- e^{\alpha_t} + d_0, \quad \text{and} \quad \tau'_t := -d_+ e^{-\alpha_t} + d_- e^{\alpha_t}.$$

We now derive a recursive equation for \tilde{F} , true for any α :

$$\tilde{F}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_{\mathbf{j}_t}) = \sum_{\{i,k \mid \pi(\mathbf{x}_i, \mathbf{x}_k) = 1\}} e^{(-\mathbf{M}\boldsymbol{\lambda}_t)_{ik}} e^{-M_{ik\mathbf{j}_t} \alpha} = \tilde{F}(\boldsymbol{\lambda}_t) (d_+ e^{-\alpha} + d_- e^{\alpha} + d_0).$$

Thus, we have defined τ_t so that $\tilde{F}(\boldsymbol{\lambda}_{t+1}) = \tilde{F}(\boldsymbol{\lambda}_t + \alpha_t \mathbf{e}_{\mathbf{j}_t}) = \tilde{F}(\boldsymbol{\lambda}_t) \tau_t$. We use this to write a recursive equation for \tilde{G} :

$$\tilde{G}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_{\mathbf{j}_t}) = \frac{-\ln(\tilde{F}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_{\mathbf{j}_t}))}{s_t + \alpha} = g_t \frac{s_t}{s_t + \alpha} - \frac{\ln(d_+ e^{-\alpha} + d_- e^{\alpha} + d_0)}{s_t + \alpha}.$$

For our algorithm, we set $\alpha = \alpha_t$ in the above expression:

$$g_{t+1} = g_t \frac{s_t}{s_t + \alpha_t} - \frac{\ln \tau_t}{s_t + \alpha_t} \Rightarrow g_{t+1} - g_t = -\frac{1}{s_{t+1}} [g_t \alpha_t + \ln \tau_t]. \quad (4)$$

Now we have gathered enough notation to write the equation for α_t for Smooth Margin Ranking. For plain coordinate ascent, the update α^* solves:

$$0 = \frac{d\tilde{G}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_{j_t})}{d\alpha} \Big|_{\alpha=\alpha^*} = \frac{1}{s_t + \alpha^*} \left[-\tilde{G}(\boldsymbol{\lambda}_t + \alpha^* \mathbf{e}_{j_t}) + \left[\frac{-d\tilde{F}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_{j_t})/d\alpha \Big|_{\alpha=\alpha^*}}{\tilde{F}(\boldsymbol{\lambda}_t)} \right] \right].$$

We could solve this equation numerically for α^* to get a smooth margin coordinate ascent algorithm, however, we avoid this line search. To get the update rule for Smooth Margin Ranking, we set α_t to solve:

$$0 = \frac{1}{s_t + \alpha_t} \left[-\tilde{G}(\boldsymbol{\lambda}_t) + \left[\frac{-d\tilde{F}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_{j_t})/d\alpha \Big|_{\alpha=\alpha_t}}{\tilde{F}(\boldsymbol{\lambda}_t)} \right] \right] = \frac{1}{s_t + \alpha_t} \left(-g_t + \frac{-\tau_t' \tilde{F}(\boldsymbol{\lambda}_t)}{\tau_t \tilde{F}(\boldsymbol{\lambda}_t)} \right) \\ g_t \tau_t = -\tau_t' \quad (5)$$

This expression can be solved analytically for α_t , which makes the algorithm as easy to implement as RankBoost, but we avoid using the exact expression in our calculations whenever possible, since it is not easy to work with:

$$\alpha_t = \ln \left[\frac{-g_t d_0 + \sqrt{g_t^2 d_0^2 + (1+g_t)(1-g_t)(1+r_t-d_0)(1-r_t-d_0)}}{(1+g_t)(1-r_t-d_0)} \right]. \quad (6)$$

We now present a theorem stating that the algorithm makes significant progress towards increasing the value of \tilde{G} at every iteration. An analogous statement was an essential tool for proving properties of Approximate Coordinate Ascent Boosting [13], although the proof here (for which we give a sketch) is significantly more difficult since we cannot use the same important trick (i.e., the equivalent of Lemma 1). As usual, the weak learning algorithm always achieves an edge of at least ρ for the calculation to hold.

Theorem 2.

$$g_{t+1} - g_t \geq \frac{1}{2} \frac{\alpha_t (r_t - g_t)}{s_{t+1}}.$$

Sketching the proof, we consider α_t , τ_t , d_+ , and d_- as functions of three basic independent variables $r := r_t$, $g := g_t$ and, d_0 , with ranges $0 < r < 1$, $0 \leq g < r$, and $0 \leq d_0 \leq 1 - r$. Define

$$\Gamma_{r,g,d_0} := \frac{-\ln \tau_t}{\alpha_t}, \quad \text{and} \quad \mathcal{B}_{r,g,d_0} := \frac{\Gamma_{r,g,d_0} - g}{r - g}.$$

Lemma 1. $\mathcal{B}_{r,g,d_0} > 1/2$.

This lemma is a monstrous calculus problem in three variables, for which we omit the proof. Using only this lemma, we can prove the theorem directly.

Proof. (of Theorem 2) Let us unravel the notation a bit. Lemma 1 says:

$$\frac{-\ln \tau_t}{\alpha_t} = \Gamma_{r,g,d_0} > \frac{r_t + g_t}{2} \Rightarrow -\ln \tau_t > \frac{(r_t + g_t)\alpha_t}{2}.$$

Incorporating the recursive equation (4),

$$g_{t+1} - g_t = \frac{1}{s_{t+1}} [-g_t \alpha_t - \ln \tau_t] > \frac{\alpha_t}{s_{t+1}} \left[-g_t + \frac{(r_t + g_t)}{2} \right] = \frac{1}{2} \frac{\alpha_t (r_t - g_t)}{s_{t+1}}.$$

We have finished the proof. Theorem 2 is the main step in proving convergence theorems, for example:

Theorem 3. *The smooth ranking margin ranking algorithm converges to a maximum margin solution, i.e., $\lim_{t \rightarrow \infty} g_t = \rho$.*

Besides Theorem 2, the only other key step in the proof of Theorem 3 is:

Lemma 2.

$$\lim_{t \rightarrow \infty} \frac{\alpha_t}{s_{t+1}} = 0.$$

We omit the proof of Lemma 2, which uses (4), monotonicity and boundedness of the g_t sequence, and then (5).

Proof. (Of Theorem 3) The values of g_t constitute a non-decreasing sequence which is uniformly bounded by 1. Thus, a limit g_∞ exists, $g_\infty := \lim_{t \rightarrow \infty} g_t$. By (3), we know $g_t \leq \rho$ for all t . Thus, $g_\infty \leq \rho$. Suppose $g_\infty < \rho$, i.e., that $\rho - g_\infty \neq 0$. One can use an identical calculation to the one in Rudin et al. [13] to show that this assumption, together with Theorem 2 and Lemma 2 imply that $\lim_{t \rightarrow \infty} \alpha_t = 0$. Using this fact along with (5), we find:

$$g_\infty = \lim_{t \rightarrow \infty} g_t = \liminf_{t \rightarrow \infty} g_t = \liminf_{t \rightarrow \infty} \frac{-\tau_t'}{\tau_t} = \liminf_{t \rightarrow \infty} \frac{-(-d_+ e^{-\alpha_t} + d_- e^{\alpha_t})}{d_+ e^{-\alpha_t} + d_- e^{\alpha_t} + d_0} = \liminf_{t \rightarrow \infty} \frac{r_t}{1} \geq \rho.$$

This is a contradiction with the original assumption that $g_\infty < \rho$. It follows that $g_\infty = \rho$, or $\lim_{t \rightarrow \infty} (\rho - g_t) = 0$. Thus, the smooth ranking algorithm converges to a maximum margin solution.

5 AdaBoost and RankBoost in the Bipartite Problem

In the bipartite ranking problem, every training instance falls into one of two categories, the “positive class” Y_+ and the “negative class” Y_- . The bipartite ranking problem is a special case of the binary ranking problem. Here, $\pi(\mathbf{x}_i, \mathbf{x}_k) = 1$ only when $\mathbf{x}_i \in Y_+$ and $\mathbf{x}_k \in Y_-$. Define $y_i = +1$ when $\mathbf{x}_i \in Y_+$, and $y_i = -1$ otherwise. The function \tilde{F} now becomes an exponentiated version of the AUC, that is, since the step function $\mathbf{1}_{x < 0} \leq e^{-x}$, we have:

$$|Y_+||Y_-|(1 - AUC(\boldsymbol{\lambda})) = \sum_{i \in Y_+} \sum_{k \in Y_-} \mathbf{1}_{(\mathbf{M}\boldsymbol{\lambda})_{ik} < 0} \leq \sum_{i \in Y_+} \sum_{k \in Y_-} e^{-(\mathbf{M}\boldsymbol{\lambda})_{ik}} = \tilde{F}(\boldsymbol{\lambda}).$$

We have written the AUC in the form of the Wilcoxon-Mann-Whitney statistic (see [5]).

We now show that AdaBoost minimizes RankBoost’s loss function \tilde{F} under very natural conditions, namely, whenever the positive and negative instances contribute equally to AdaBoost’s loss function.

We define AdaBoost’s matrix \mathbf{M}^{Ada} element-wise by $M_{ij}^{Ada} = y_i h_j(\mathbf{x}_i)$. Thus, $M_{ikj} = h_j(\mathbf{x}_i) - h_j(\mathbf{x}_k) = y_i h_j(\mathbf{x}_i) + y_k h_j(\mathbf{x}_k) = M_{ij}^{Ada} + M_{kj}^{Ada}$. (To change from $\{-1, 1\}$ hypotheses to $\{0, 1\}$ hypotheses, divide by 2.) Define vector \mathbf{q}_λ element-wise by $q_{\lambda,i} := e^{-(\mathbf{M}^{Ada}\lambda)_i}$. Using this notation, we will write the objective functions for both AdaBoost and RankBoost. First, we define the following:

$$F_+(\lambda) := \sum_{i \in Y_+} q_{\lambda,i} \quad \text{and} \quad F_-(\lambda) := \sum_{k \in Y_-} q_{\lambda,k}.$$

The objective function for AdaBoost is $F(\lambda) := F_+(\lambda) + F_-(\lambda)$. The objective function for RankBoost is: $\tilde{F}(\lambda) = F_+(\lambda)F_-(\lambda)$. Thus, the balance between the positive and negative instances is different between the algorithms.

We now define “F-skew”, which measures the imbalance between positive and negative instances.

$$\text{F-skew}(\lambda) := F_+(\lambda) - F_-(\lambda).$$

Theorem 4. *Assume \mathbf{M}^{Ada} is such that $\inf_\lambda \tilde{F}(\lambda) > 0$ (the non-separable case). For any sequence $\{\lambda_t\}_{t=1}^\infty$ such that*

$$\lim_{t \rightarrow \infty} F(\lambda_t) = \inf_\lambda F(\lambda) \tag{7}$$

and $\lim_{t \rightarrow \infty} \text{F-skew}(\lambda_t) = 0$, then

$$\lim_{t \rightarrow \infty} \tilde{F}(\lambda_t) = \inf_\lambda \tilde{F}(\lambda). \tag{8}$$

Proof. It is possible that F or \tilde{F} may have no minimizers. So, to describe (7) and (8), we use the trick from Collins et al. [4], who considered F and \tilde{F} as functions of a variable where the infimum can be achieved. Define, for matrix $\bar{\mathbf{M}}$, the function

$$F_{\bar{\mathbf{M}}}(\lambda) := \sum_i e^{-(\bar{\mathbf{M}}\lambda)_i}.$$

Define $\bar{\mathcal{P}} := \{\mathbf{p} | p_i \geq 0 \forall i, (\mathbf{p}^T \bar{\mathbf{M}})_j = 0 \forall j\}$, and $\bar{\mathcal{Q}} := \{\mathbf{q} | q_i = e^{-(\bar{\mathbf{M}}\lambda)_i} \text{ for some } \lambda\}$. We may thus consider $\bar{F}_{\bar{\mathbf{M}}}$ as a function of $\bar{\mathbf{q}}$, $\bar{F}_{\bar{\mathbf{M}}}(\bar{\mathbf{q}}) = \sum_i \bar{q}_i$, where $\bar{\mathbf{q}} \in \bar{\mathcal{Q}}$. We know that since all \bar{q}_i ’s are positive, the infimum of \bar{F} occurs in a bounded region of $\bar{\mathbf{q}}$ space, which is just what we need.

Theorem 1 of Collins et al., which is taken directly from Lafferty, Della Pietra, and Della Pietra [11] implies that the following are equivalent:

1. $\bar{\mathbf{q}}^* \in \bar{\mathcal{P}} \cap \text{closure}(\bar{\mathcal{Q}})$.
2. $\bar{\mathbf{q}}^* \in \text{argmin}_{\bar{\mathbf{q}} \in \text{closure}(\bar{\mathcal{Q}})} \bar{F}_{\bar{\mathbf{M}}}(\bar{\mathbf{q}})$.

Moreover, either condition is satisfied by exactly one vector $\bar{\mathbf{q}}^*$.

The objective function for AdaBoost is $F = \bar{F}_{\mathbf{M}^{Ada}}$ and the objective for RankBoost is $\tilde{F} = \bar{F}_{\mathbf{M}}$, so the theorem holds for both objectives separately. For

function F , denote $\bar{\mathbf{q}}^*$ as \mathbf{q}^* , also $\bar{\mathcal{P}}$ as \mathcal{P}^{Ada} and $\bar{\mathcal{Q}}$ as \mathcal{Q}^{Ada} . For function \tilde{F} , denote $\bar{\mathbf{q}}^*$ as $\tilde{\mathbf{q}}^*$, also $\bar{\mathcal{P}}$ as $\tilde{\mathcal{P}}$ and $\bar{\mathcal{Q}}$ as $\tilde{\mathcal{Q}}$. Rewriting $\mathbf{q}^* \in \mathcal{P}^{Ada}$:

$$\sum_{i \in \mathcal{I}_+} q_i^* M_{ij}^{Ada} + \sum_{k \in \mathcal{I}_-} q_k^* M_{kj}^{Ada} = 0 \quad \forall j. \quad (9)$$

Define \mathbf{q}_t element-wise by: $q_{t,i} := e^{-(\mathbf{M}^{Ada} \boldsymbol{\lambda}_t)_i}$, where the $\boldsymbol{\lambda}_t$'s are a sequence that obey (7), for example, a sequence produced by AdaBoost. Thus, $\mathbf{q}_t \in \mathcal{Q}^{Ada}$ automatically. Since $F(\mathbf{q}_t)$ converges to the minimum of F , one can show that the sequence of \mathbf{q}_t 's converges to \mathbf{q}^* in ℓ_p . Now define vectors $\tilde{\mathbf{q}}_t$ element-wise by $\tilde{q}_{t,ik} := q_{t,i} q_{t,k} = \exp[-(\mathbf{M}^{Ada} \boldsymbol{\lambda}_t)_i - (\mathbf{M}^{Ada} \boldsymbol{\lambda}_t)_k] = \exp[-(\mathbf{M} \boldsymbol{\lambda}_t)_{ik}]$. Automatically, $\tilde{\mathbf{q}}_t \in \tilde{\mathcal{Q}}$. For any pair i, k the limit of the $\tilde{q}_{t,ik}$'s is $\tilde{q}_{ik}^\infty := q_i^* q_k^*$. Thus, we need only to show $\tilde{\mathbf{q}}^\infty = \tilde{\mathbf{q}}^*$. We will do this by showing $\tilde{\mathbf{q}}^\infty \in \tilde{\mathcal{P}}$; due to the uniqueness of $\tilde{\mathbf{q}}^*$ as $\tilde{\mathcal{P}} \cap \text{closure}(\tilde{\mathcal{Q}})$, this will yield $\tilde{\mathbf{q}}^\infty = \tilde{\mathbf{q}}^*$.

Our assumption that the F-skew vanishes can be rewritten as:

$$\begin{aligned} \lim_{t \rightarrow \infty} \left[\sum_{i \in \mathcal{I}_+} q_{t,i} - \sum_{k \in \mathcal{I}_-} q_{t,k} \right] &= 0, \quad \text{i.e.,} \\ \sum_{i \in \mathcal{I}_+} q_i^* &= \sum_{k \in \mathcal{I}_-} q_k^*. \end{aligned} \quad (10)$$

Consider the quantities $(\tilde{\mathbf{q}}^{\infty T} \mathbf{M})_j$. Remember, if these quantities are zero for every j , then $\tilde{\mathbf{q}}^\infty \in \tilde{\mathcal{P}}$ and we have proved the theorem.

$$(\tilde{\mathbf{q}}^{\infty T} \mathbf{M})_j = \left(\sum_{k \in Y_-} q_k^* \right) \left(\sum_{i \in Y_+} q_i^* M_{ij}^{Ada} \right) + \left(\sum_{i \in Y_+} q_i^* \right) \left(\sum_{k \in Y_-} q_k^* M_{kj}^{Ada} \right).$$

Incorporating (10), which is the condition that $\text{F-skew}(\mathbf{q}^*) = 0$, (11) becomes:

$$(\tilde{\mathbf{q}}^{\infty T} \mathbf{M})_j = \left(\sum_{i \in Y_+} q_i^* \right) \left[\sum_{i \in Y_+} q_i^* M_{ij}^{Ada} + \sum_{k \in Y_-} q_k^* M_{kj}^{Ada} \right].$$

In fact, according to (9), the bracket in this expression is zero for all j . Thus, $\tilde{\mathbf{q}}^\infty \in \tilde{\mathcal{P}}$. We have proved the theorem.

Corollary 1. *If the constant weak hypothesis $h_j(\mathbf{x}) = 1 \quad \forall \mathbf{x}$ is one of the weak classifiers used to construct \mathbf{M}^{Ada} , and the $\{\boldsymbol{\lambda}_t\}_t$ sequence obeys (7), then $\lim_{t \rightarrow \infty} \text{F-skew}(\boldsymbol{\lambda}_t) = 0$, and $\{\boldsymbol{\lambda}_t\}_t$ thus obeys (8).*

This result and Theorem 4 together imply that if the constant weak hypothesis is included in the set of weak classifiers, then the F-skew vanishes. That is, any algorithm which minimizes F (such as AdaBoost) solves the ranking problem whenever the constant hypothesis is included.

Proof. Recall that $\mathbf{q}^* \in \mathcal{P}^{Ada}$. Specifically writing this condition just for the constant weak classifier yields:

$$0 = \sum_{i \in Y_+} q_i^* y_i 1 + \sum_{k \in Y_-} q_k^* y_k 1 = \sum_{i \in Y_+} q_i^* - \sum_{k \in Y_-} q_k^* = \lim_{t \rightarrow \infty} \text{F-skew}(\boldsymbol{\lambda}_t).$$

Thus, AdaBoost and RankBoost are closely related indeed, since under very weak conditions (e.g., when the constant weak classifier is included), AdaBoost minimizes RankBoost's objective function. Given these results, it is no longer surprising that AdaBoost performs so well for the ranking problem.

References

- [1] Shivani Agarwal, Thore Graepel, Ralf Herbich, Sarel Har-Peled, and Dan Roth. Generalization bounds for the area under the ROC curve. *Journal of Machine Learning Research*, 2005. Accepted pending revision.
- [2] Olivier Bousquet. New approaches to statistical learning theory. *Annals of the Institute of Statistical Mathematics*, 55(2):371–389, 2003.
- [3] Rich Caruana and Alexandru Niculescu-Mizil. An empirical comparison of supervised learning algorithms using difference performance metrics. Technical report, Cornell University, 2005.
- [4] Michael Collins, Robert E. Schapire, and Yoram Singer. Logistic regression, AdaBoost and Bregman distances. *Machine Learning*, 48(1/2/3), 2002.
- [5] Corinna Cortes and Mehryar Mohri. AUC optimization vs. error rate minimization. In *Advances in Neural Information Processing Systems 16*, 2004.
- [6] Felipe Cucker and Steve Smale. On the mathematical foundations of learning. *Bull. Amer. Math. Soc.*, (39):1–49, 2002.
- [7] Yoav Freund, Raj Iyer, Robert E. Schapire, and Yoram Singer. An efficient boosting algorithm for combining preferences. In *Machine Learning: Proceedings of the Fifteenth International Conference*, 1998.
- [8] Yoav Freund and Robert E. Schapire. Adaptive game playing using multiplicative weights. *Games and Economic Behavior*, 29:79–103, 1999.
- [9] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Additive logistic regression: A statistical view of boosting. *The Annals of Statistics*, 38(2):337–374, April 2000.
- [10] Vladimir Koltchinskii and Dmitry Panchenko. Empirical margin distributions and bounding the generalization error of combined classifiers. *The Annals of Statistics*, 30(1), February 2002.
- [11] John D. Lafferty, Stephen Della Pietra, and Vincent Della Pietra. Statistical learning algorithms based on Bregman distances. In *Proceedings of the Canadian Workshop on Information Theory*, 1997.
- [12] Cynthia Rudin, Ingrid Daubechies, and Robert E. Schapire. The dynamics of AdaBoost: Cyclic behavior and convergence of margins. *Journal of Machine Learning Research*, 5:1557–1595, December 2004.
- [13] Cynthia Rudin, Robert E. Schapire, and Ingrid Daubechies. Analysis of boosting algorithms using the smooth margin function: A study of three algorithms. Submitted, 2004.
- [14] Cynthia Rudin, Robert E. Schapire, and Ingrid Daubechies. Boosting based on a smooth margin. In *Proceedings of the Sixteenth Annual Conference on Computational Learning Theory*, pages 502–517, 2004.
- [15] Robert E. Schapire, Yoav Freund, Peter Bartlett, and Wee Sun Lee. Boosting the margin: A new explanation for the effectiveness of voting methods. *The Annals of Statistics*, 26(5):1651–1686, October 1998.

A Proof of Theorem 1

We owe inspiration for this proof to the works of Cucker and Smale [6], Koltchinskii and Panchenko [10], and Bousquet [2].

We define a Lipschitz function $\phi : \mathbb{R} \rightarrow \mathbb{R}$ (with Lipschitz constant $\text{Lip}(\phi)$) which acts as our loss function, and give us the margin. We will eventually use the same piecewise linear definition of ϕ as Koltchinskii and Panchenko [10], but for now, we require $\phi(z) \geq 0 \forall z$ and $\phi(z) = 1$ for $z < 0$. Since $\phi(z) \geq \mathbf{1}_{[z \leq 0]}$, we can define an upper bound for the misranking probability defined above. We have $\mathbb{P}_D\{\text{misrank}_f\} \leq \mathbb{P}\phi_f$, where:

$$\mathbb{P}\phi_f := \frac{\mathbb{E}_{\bar{\mathbf{x}}, \tilde{\mathbf{x}} \sim D}[\phi(f(\bar{\mathbf{x}}) - f(\tilde{\mathbf{x}}))\pi(\bar{\mathbf{x}}, \tilde{\mathbf{x}})]}{\mathbb{E}_{\bar{\mathbf{x}}, \tilde{\mathbf{x}} \sim D}\pi(\bar{\mathbf{x}}, \tilde{\mathbf{x}})}.$$

The empirical error associated with $\mathbb{P}\phi_f$ is:

$$\mathbb{P}_S\phi_f := \frac{\sum_{i=1}^m \sum_{k=1}^m \phi(f(\mathbf{x}_i) - f(\mathbf{x}_k))\pi(\mathbf{x}_i, \mathbf{x}_k)}{\sum_{i=1}^m \sum_{k=1}^m \pi(\mathbf{x}_i, \mathbf{x}_k)}$$

First, we upper bound the misranking probability by two terms: the empirical error term $\mathbb{P}_S\phi_f$, and a term characterizing the deviation of $\mathbb{P}_S\phi_f$ from $\mathbb{P}\phi_f$ uniformly for all $f \in \mathcal{F}$.

$$\mathbb{P}_D\{\text{misrank}_f\} \leq \mathbb{P}\phi_f = \mathbb{P}\phi_f - \mathbb{P}_S\phi_f + \mathbb{P}_S\phi_f \leq \sup_{\bar{f} \in \mathcal{F}} (\mathbb{P}\phi_{\bar{f}} - \mathbb{P}_S\phi_{\bar{f}}) + \mathbb{P}_S\phi_f$$

The proof of the theorem involves an upper bound on the first term. First, define $L(f)$ as follows: $L(f) := \mathbb{P}\phi_f - \mathbb{P}_S\phi_f$. The following lemma (for which the proof is omitted) is true for every training set S :

Lemma 3. For any two functions $f_1, f_2 \in L_\infty(\mathcal{X})$,

$$L(f_1) - L(f_2) \leq 4\text{Lip}(\phi)\|f_1 - f_2\|_\infty.$$

The following step is due to Cucker and Smale [6]. Let $\ell_\epsilon := \mathcal{N}\left(\mathcal{F}, \frac{\epsilon}{8\text{Lip}(\phi)}\right)$, the covering number of \mathcal{F} by L_∞ disks of radius $\frac{\epsilon}{8\text{Lip}(\phi)}$. Define $f_1, f_2, \dots, f_{\ell_\epsilon}$ to be the centers of such a cover, i.e., the collection of L_∞ disks D_p centered at f_p and with radius $\frac{\epsilon}{8\text{Lip}(\phi)}$ is a cover for \mathcal{F} . In the proof of the theorem, we will use the center of each disk to act as a representative for the whole disk. The following lemma (proof omitted) shows we do not lose too much by using f_p as a representative for disk D_p .

Lemma 4.

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left\{ \sup_{f \in D_p} |L(f)| \geq \epsilon \right\} \leq \mathbb{P}_{S \sim \mathcal{D}^m} \left\{ |L(f_p)| \geq \frac{\epsilon}{2} \right\}$$

Now we incorporate the fact that the training set is chosen randomly.

Lemma 5.

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left\{ |L(f)| \geq \epsilon/2 \right\} \leq 4 \exp \left[-2m \left(\frac{\epsilon E}{4 + \epsilon} \right)^2 \right]$$

Proof. To make notation easier for this lemma, we introduce some shorthand notation:

$$\begin{aligned} \text{top}_D &:= \mathbb{E}_{\bar{\mathbf{x}}, \tilde{\mathbf{x}} \sim D} \phi(f(\bar{\mathbf{x}}) - f(\tilde{\mathbf{x}}))\pi(\bar{\mathbf{x}}, \tilde{\mathbf{x}}), & \text{bot}_D &:= E := \mathbb{E}_{\bar{\mathbf{x}}, \tilde{\mathbf{x}} \sim D} \pi(\bar{\mathbf{x}}, \tilde{\mathbf{x}}), \\ \text{top}_S &:= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{k=1}^m \phi(f(\mathbf{x}_i) - f(\mathbf{x}_k))\pi(\mathbf{x}_i, \mathbf{x}_k), & \text{bot}_S &:= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{k=1}^m \pi(\mathbf{x}_i, \mathbf{x}_k). \end{aligned}$$

Since diagonal terms are $\pi(\mathbf{x}_i, \mathbf{x}_i) = 0$, $\text{top}_D = \mathbb{E}_{S \sim \mathcal{D}^m} \text{top}_S$ and $\text{bot}_D = \mathbb{E}_{S \sim \mathcal{D}^m} \text{bot}_S$. Thus, we can bound the difference between top_S and top_D using large deviation bounds; same for bot_S and bot_D . One can show that the replacement of one instance changes top_S (or bot_S) by at most $1/m$. Thus, McDiarmid's inequality implies, for every $\epsilon_1 > 0$:

$$\mathbb{P}\{|\text{top}_S - \text{top}_D| \geq \epsilon_1\} \leq 2 \exp[-2\epsilon_1^2 m] \quad \text{and} \quad \mathbb{P}\{|\text{bot}_S - \text{bot}_D| \geq \epsilon_1\} \leq 2 \exp[-2\epsilon_1^2 m].$$

We will specify ϵ_1 in terms of ϵ later. Consider the following event:

$$|\text{top}_S - \text{top}_D| < \epsilon_1 \quad \text{and} \quad |\text{bot}_S - \text{bot}_D| < \epsilon_1.$$

By the union bound, this event is true with probability at least $1 - 4 \exp[-2\epsilon_1^2 m]$. When the event is true, we can rearrange the equations to be a bound on $L(f) = \frac{\text{top}_D}{\text{bot}_D} - \frac{\text{top}_S}{\text{bot}_S}$. Namely, as long as ϵ_1 is smaller than top_S and bot_S ,

$$\left| \frac{\text{top}_S}{\text{bot}_S} - \frac{\text{top}_D}{\text{bot}_D} \right| < \frac{\text{top}_D + \epsilon_1}{\text{bot}_D - \epsilon_1} - \frac{\text{top}_D}{\text{bot}_D} =: \epsilon/2.$$

Above, we have just specified the value for ϵ_1 in terms of ϵ . Let us solve for ϵ_1 :

$$\epsilon_1 = \frac{\epsilon \text{bot}_D}{2 + \epsilon + 2 \frac{\text{top}_D}{\text{bot}_D}} \geq \frac{\epsilon E}{4 + \epsilon}$$

Here, we have used $E := \text{bot}_D$, and by definition, $\text{top}_D \leq \text{bot}_D$. We directly have:

$$1 - 4 \exp[-2\epsilon_1^2 m] \geq 1 - 4 \exp\left(-2m \left[\frac{\epsilon E}{4 + \epsilon}\right]^2\right)$$

Therefore, from our earlier application of McDiarmid, with probability at least $1 - 4 \exp\left(-2m \left[\frac{\epsilon E}{4 + \epsilon}\right]^2\right)$ the following holds: $|L(f)| = \left| \frac{\text{top}_S}{\text{bot}_S} - \frac{\text{top}_D}{\text{bot}_D} \right| < \epsilon/2$.

Proof. (of Theorem 1) Since the D_p are a cover of \mathcal{F} , it is true that

$$\sup_{f \in \mathcal{F}} |L(f)| \geq \epsilon \iff \exists p \leq \ell_\epsilon \text{ such that } \sup_{f \in D_p} |L(f)| \geq \epsilon.$$

First applying the union bound, then applying Lemma 4, and then Lemma 5, we find:

$$\begin{aligned} \mathbb{P}_S \left\{ \sup_{f \in \mathcal{F}} |L(f)| \geq \epsilon \right\} &\leq \sum_{p=1}^{\ell_\epsilon} \mathbb{P} \left\{ \sup_{f \in D_p} |L(f)| \geq \epsilon \right\} \leq \sum_{p=1}^{\ell_\epsilon} \mathbb{P} \{ |L(f_p)| \geq \epsilon/2 \} \\ &\leq \sum_{p=1}^{\ell_\epsilon} 4 \exp \left[-2m \left(\frac{\epsilon E}{4 + \epsilon} \right)^2 \right] = \mathcal{N} \left(\mathcal{F}, \frac{\epsilon}{4 \text{Lip}(\phi)} \right) 4 \exp \left[-2m \left(\frac{\epsilon E}{4 + \epsilon} \right)^2 \right]. \end{aligned}$$

Now, with probability at least $1 - \mathcal{N} \left(\mathcal{F}, \frac{\epsilon}{4 \text{Lip}(\phi)} \right) 4 \exp \left[-2m \left(\frac{\epsilon E}{4 + \epsilon} \right)^2 \right]$, we have

$\mathbb{P}_D \{\text{misrank}_f\} \leq \mathbb{P}_S \phi_f + \epsilon$. Now let $\phi(z) = 1$ for $z \leq 0$, $\phi(z) = 0$ for $z \geq \theta$, and linear in between, with slope $1/\theta$. Thus, $\text{Lip}(\phi) = 1/\theta$. Since $\phi(z) \leq 1$ for $z \leq \theta$, we have:

$$\mathbb{P}_S \phi_f = \frac{\sum_{i=1}^m \sum_{k=1}^m \phi(f(\mathbf{x}_i) - f(\mathbf{x}_k)) \pi(\mathbf{x}_i, \mathbf{x}_k)}{\sum_{i=1}^m \sum_{k=1}^m \pi(\mathbf{x}_i, \mathbf{x}_k)} \leq \mathbb{P}_S \{\text{margin}_f \leq \theta\}.$$

Thus, with probability at least $1 - \mathcal{N} \left(\mathcal{F}, \frac{\epsilon \theta}{4} \right) 4 \exp \left[-2m \left(\frac{\epsilon E}{4 + \epsilon} \right)^2 \right]$, we have $\mathbb{P}_D \{\text{misrank}_f\} \leq \mathbb{P}_S \{\text{margin}_f \leq \theta\} + \epsilon$. Thus, the theorem has been proved.