

---

# Boosting as Entropy Projection

---

**Jyrki Kivinen\***

Department of Computer Science  
P.O. Box 26 (Teollisuuskatu 23)  
FIN-00014 University of Helsinki  
Finland  
Jyrki.Kivinen@cs.Helsinki.FI

**Manfred K. Warmuth†**

Computer Science Department  
University of California, Santa Cruz  
Santa Cruz, CA 95064  
U.S.A.  
manfred@cse.ucsc.edu

## Abstract

We consider the AdaBoost procedure for boosting weak learners. In AdaBoost, a key step is choosing a new distribution on the training examples based on the old distribution and the mistakes made by the present weak hypothesis. We show how AdaBoost's choice of the new distribution can be seen as an approximate solution to the following problem: Find a new distribution that is closest to the old distribution subject to the constraint that the new distribution is orthogonal to the vector of mistakes of the current weak hypothesis. The distance (or divergence) between distributions is measured by the relative entropy. Alternatively, we could say that AdaBoost approximately projects the distribution vector onto a hyperplane defined by the mistake vector. We show that this new view of AdaBoost as an entropy projection is dual to the usual view of AdaBoost as minimizing the normalization factors of the updated distributions.

## 1 Introduction

Boosting, originally suggested by Schapire [Sch90], is a particular method for improving the performance of a (supervised) learning algorithm by applying it several times on slightly modified training data and then combining the results in a suitable manner. Currently the most popular variants of boosting are based on Freund and Schapire's AdaBoost [FS97b]. The details of the boosting framework of our paper are mainly taken from Schapire and Singer's work on confidence-rated boosting [SS98].

Let us review the basic idea of boosting on a very rough level. We take as our starting point an arbitrary learning algorithm, which in this context is called the *weak learner* (as

opposed to the *master algorithm* that implements the whole boosting procedure). We also have a fixed *training set* of examples. Following Freund [Fre95], we choose some probability distribution over the training set as the initial *training distribution*. We then repeat the following until some termination condition is met. We call the weak learner and draw its training examples from the set of all training examples according to the current training distribution. The weak learner produces a *weak hypothesis*. We use the weak hypothesis to update the old training distribution into a new one. The details of this update will be discussed shortly. We then set the weak hypothesis aside and go to the next iteration with the new training distribution. After the termination condition for the iterations is met, the master algorithm outputs as its hypothesis a suitable weighted combination of all the weak hypotheses produced during this process.

The basic question we consider in this paper is how to update the training distribution between the calls to the weak learner. Assuming we have  $m$  examples in the training set, we represent the training distribution for the  $t$ th call to the weak learner as a distribution vector  $\mathbf{d}_t \in [0, 1]^m$  such that  $\sum_i d_{t,i} = 1$ . Lacking a reason to do otherwise we would typically choose the initial distribution to be uniform, with  $d_{1,i} = 1/m$  for all  $i$ . Based on the training distribution  $\mathbf{d}_t$ , the weak learner produces a weak hypothesis  $h_t$ . We describe the performance of the weak hypothesis  $h_t$  on the training set by the vector  $\mathbf{u}_t \in [-1, 1]^m$ , where  $u_{t,i}$  indicated the goodness of the algorithm on example number  $i$ . In the most basic case we would choose  $u_{t,i} = 1$  if  $h_t$  predicts correctly on the  $i$ th example and  $u_{t,i} = -1$  otherwise, but more fine-grained measures are also possible. The fundamental idea of boosting is now to concentrate the new training distribution on those examples on which the current weak hypothesis performs badly. In particular, in AdaBoost [FS97b] the updated distribution has the following *exponential form*

$$d_{t+1,i} = \frac{1}{Z_t} d_{t,i} \exp(-\alpha_t u_{t,i}), \quad (1.1)$$

where  $\alpha_t > 0$  regulates the amount of change and  $Z_t$  is a normalization factor that gives  $\sum_i d_{t+1,i} = 1$ . The final hypothesis  $H$  of the master algorithm is given by  $H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right)$ . Thus  $\alpha_t$  also acts as the weight given to  $h_t$ . Schapire and Singer [SS98] show in an elegant proof that the training error of the final hypothesis after  $T$  boosting iterations is bounded by the prod-

---

\*Supported by ESPRIT Working Group NeuroCOLT2

†Supported by NSF grant CCR 9700201

uct  $\prod_{t=1}^T Z_t$  of the normalization factors. They show that the choice of  $\alpha_t$  in the original AdaBoost [FS97b], namely  $\alpha_t = \ln((1 + \mathbf{d}_t \cdot \mathbf{u}_t)/(1 - \mathbf{d}_t \cdot \mathbf{u}_t))/2$ , minimizes  $Z_t$  in the discrete case (when  $u_{t,i} \in \{-1, 1\}$ ). In the continuous-valued case (when  $u_{t,i} \in [-1, 1]$ ) the same choice of  $\alpha_t$  only minimizes a certain upper bound for  $Z_t$ . As an alternative, they suggest choosing  $\alpha_t$  in the continuous valued case so that  $Z_t$  is exactly minimized. Further, they show that minimizing  $Z_t$  exactly occurs at a unique value  $\alpha_t$  such that the dot product  $\mathbf{d}_{t+1} \cdot \mathbf{u}_t = \sum_i d_{t+1,i} u_{t,i}$  is zero. We call the update (1.1), with  $\alpha_t$  such that  $\mathbf{d}_{t+1} \cdot \mathbf{u}_t = 0$ , the *corrective update*.

Our central idea is to view the corrective update as a solution to a relative entropy minimization problem and the choice of  $\alpha_t$  used by AdaBoost as an approximate solution to the same problem. Define the relative entropy between distribution vectors  $\mathbf{d}$  and  $\tilde{\mathbf{d}}$  by

$$\Delta(\mathbf{d}, \tilde{\mathbf{d}}) = \sum_{i=1}^m d_i \ln \frac{d_i}{\tilde{d}_i} . \quad (1.2)$$

and consider the minimization problem

$$\min_{\mathbf{d} \in P_m} \Delta(\mathbf{d}, \mathbf{d}_t) \quad \text{subject to } \mathbf{d} \cdot \mathbf{u}_t = 0 . \quad (1.3)$$

Here  $P_m = \{ \mathbf{d} \in \mathbf{R}^m \mid \sum_{i=1}^m d_i = 1, d_i \geq 0 \}$  is the set of  $m$ -dimensional distribution vectors, or the  $m$ -dimensional *probability simplex*. We show that the corrective update is the solution to the above constrained minimization problem:

$$\mathbf{d}_{t+1} = \operatorname{argmin}_{\mathbf{d} \in P_m} \Delta(\mathbf{d}, \mathbf{d}_t) \quad \text{subject to } \mathbf{d} \cdot \mathbf{u}_t = 0 . \quad (1.4)$$

More specifically, employing the standard concepts and tools from constrained convex optimization [Lue84, BSS93, HUL91], we show that the constrained minimization problem (1.3) is the dual of the unconstrained problem of maximizing  $-\ln Z_t$  as a function of  $\alpha_t$ . The variable  $\alpha_t$  of the unconstrained problem is effectively the Lagrange multiplier used to enforce the constraint  $\mathbf{d} \cdot \mathbf{u}_t = 0$  in (1.4). Also the value of the minimization and maximization problems are the same, i.e.,

$$\min_{\substack{\mathbf{d} \in P_m \\ \mathbf{d} \cdot \mathbf{u}_t = 0}} \Delta(\mathbf{d}, \mathbf{d}_t) = \max_{\alpha \in \mathbf{R}} (-\ln Z_t(\alpha)) . \quad (1.5)$$

The relative entropy is of course a very commonly used tool in statistics and in computational learning theory. We wish here to consider two different aspects of applying the relative entropy in the context of on-line learning. First, we can use the relative entropy to analyse the convergence and other properties of existing on-line algorithms. Second, we can use relative entropy to motivate new algorithms. The purpose of the present paper is to bring out explicitly in the context of boosting the connection between these two aspects by means of the duality property (1.5). Thus, we see more clearly the relationship between boosting as a minimizer for  $Z_t$  [SS98] and the analyses of boosting-style algorithms in terms of the relative entropy [FS97a, FS97b].

Of course, since we are basically considering two mathematically equivalent derivations for the single update rule (1.1), most if not all of this is already implicit in earlier work. In particular, a procedure more or less equivalent with the

corrective boosting algorithm, but in a context somewhat different from boosting weak learners, was analysed using a duality relation similar to (1.5) by Della Pietra et al. [DDL97]; see Lafferty [Laf99] for connecting this work to boosting as it is understood in computational learning theory.

Considering problems other than boosting, one should notice work on on-line prediction algorithms using experts [FSSW97, KW99] and linear regression [KW97]. In this context, the relative entropy has been used in the same kind of double role as here, both in deriving updates and then proving (worst-case) performance bounds for them.

Outside the context of on-line learning theory, and its worst-case bounds, relative entropy minimization with linear constraints is of course an important method of statistics [KK92, Jum90]. Even more generally, the relative entropy is a special case of a *Bregman divergence* [Bre67, Csi91]. Iterative projection algorithms with respect to arbitrary Bregman divergences in the more general case of inequality constraints have been studied extensively in convex optimization [Bre67, CL81, JB90]. In Appendix B we give some notes on generalizing the boosting update (1.4) and the duality connection (1.5) to arbitrary Bregman divergences (but at this point we are unable to show that such updates actually boost weak learners). Similar generalizations have been done in parallel work by Lafferty [Laf99].

For solving the minimization problem (1.3), one can use standard methods of constrained convex optimization; see [HUL91] for an overview. Here we want to point out two earlier papers that use the actual boosting update (1.1) to solve (ostensibly) a different numerical problem. Littlestone, Long and Warmuth [LLW92] suggest this update for solving iteratively a system of linear equations with a sparse solution. Cesa-Bianchi, Krogh, and Warmuth [CBKW94] developed the same algorithm in the context of finding a maximum likelihood model from an exponential family. Both papers actually give a more general algorithm that corresponds to (1.4) generalized to allow multiple linear constraints, but the algorithms can naturally be specialized to the one-constraint case (1.4). It turns out that in both cases a single iteration step of the one-constraint algorithm is exactly the same as the update step of the original AdaBoost. In particular, the choice of  $\alpha_t$  is the same.

In the exponential form update (1.1), there is no obvious reason why we need to have only a single real parameter  $\alpha_t$  to adjust at update  $t$ . The update (1.1) can naturally be generalized to

$$d_{t+1,i} = \frac{1}{Z'_t} d_{t,i} \exp \left( - \sum_{q=1}^t \alpha_{t,q} u_{q,i} \right) \quad (1.6)$$

where again  $Z'_t$  is the normalization factor and now a parameter  $\alpha_{t,q}$  is chosen for *each* of the past  $t$  weak hypotheses. (The update (1.1) uses only a parameter  $\alpha_t$  for the most recent hypothesis.) Again, the product of the normalization factors  $Z'_t$  bounds the training error of the final hypothesis, and it is natural to choose a parameter vector  $\boldsymbol{\alpha}_t$  such that  $Z'_t$  is minimized. Analogously with the corrective case, the unconstrained problem of maximizing  $-\ln Z'_t$  as a function of  $\boldsymbol{\alpha}_t \in \mathbf{R}^t$  is dual to the following constrained problem: Minimize the relative entropy  $\Delta(\mathbf{d}_{t+1}, \mathbf{d}_t)$  subject to  $\mathbf{d}_{t+1} \cdot \mathbf{u}_q = 0$  for all  $q \leq t$ . Again the original variables

$\alpha_{t,q}$  become Lagrange multipliers in the dual problem. Note that the update (1.6) may be seen as an extended exponential form. We call this update, when  $\alpha_t$  is chosen such that  $\mathbf{d}_{t+1} \cdot \mathbf{u}_q = 0$  holds for  $1 \leq q \leq t$ , the *totally corrective update*.

As Schapire and Singer [SS98] observe, the property  $\mathbf{d}_{t+1} \cdot \mathbf{u}_t = 0$  of the corrective update has the intuitive meaning that the new distribution should be uncorrelated with the mistakes made by the current weak hypothesis. Then it seems that the new weak hypothesis, trained on the new distribution, should be more likely to give us information not present in the current weak hypothesis. Given this intuitive motivation, it would seem perhaps even better to have the new distribution uncorrelated with all the previous weak hypotheses, leading us to the totally corrective algorithm.

First consider briefly implementing the corrective and totally corrective algorithms. For the corrective algorithm, there is only one parameter  $\alpha_t$  that can be determined by a simple line search [SS98]. For the totally corrective algorithm there are situations in which all the  $t$  constraints  $\mathbf{d}_{t+1} \cdot \mathbf{u}_q = 0$  cannot be simultaneously satisfied. Also, even if a solution exists, finding it is a  $t$ -dimensional numerical problem that seems to be nontrivial. A simple method is to repeatedly cycle over the past hypotheses updating one parameter at a time with an approximated corrective algorithm, such as AdaBoost. We discuss these numerical issues briefly in Appendix A. One could of course also use the minimization procedures of [LLW92, CBKW94], or any other general convex optimization algorithm for implementing the totally corrective algorithm. However, the convergence bounds in [LLW92, CBKW94] are given in terms of quantities that do not have a natural interpretation in the boosting context.

Instead of getting too involved with the implementation problems of the totally corrective algorithm, we prefer to ignore them and keep the totally corrective algorithm mostly as a conceptual tool for comparison with the corrective algorithm. Actually, we feel that the totally corrective algorithm may not always be the right approach. This is partly because of the problems just mentioned, but we also expect that the totally corrective update might lead boosting to overfit in certain circumstances. However, Della Pietra et al. [DDL97] have successfully used a method analogous to totally corrective boosting also in practice.

As an alternative to the corrective update, equivalent with (1.4), we suggest the update based on

$$\mathbf{d}_{t+1} = \underset{\mathbf{d} \in P_m}{\operatorname{argmin}} (\Delta(\mathbf{d}, \mathbf{d}_t) + \eta L(\mathbf{d} \cdot \mathbf{u}_t)), \quad (1.7)$$

where  $\eta$  is a positive parameter and  $L$  is some loss function. We assume  $L(z) \geq 0$  with equality holding iff  $z = 0$ . In the limit of  $\eta$  approaching infinity, (1.7) reduces to minimizing the relative entropy subject to  $L(\mathbf{d} \cdot \mathbf{u}_t) = 0$  (or equivalently  $\mathbf{d} \cdot \mathbf{u}_t = 0$ ) (1.4). However, by choosing different values of  $\eta$  we can control the trade-off between the tendency to be corrective and the tendency to be conservative, i.e., not move too much in a single update. Thus  $\eta$  can be considered a learning rate parameter. The Exponentiated Gradient algorithm for on-line linear regression has been derived by Kivinen and Warmuth [KW97] as an approximate solution to (1.7). A particularly intriguing connection is that (1.7) with  $\eta = 1$  and a certain entropic loss function  $L$  gives ex-

actly the AdaBoost update. It would be very interesting to see what kind of boosting results could be proved for algorithms based on other loss functions  $L$  and values of  $\eta$ .

We continue by giving in Section 2 a brief review of the boosting algorithms and the error bound of Schapire and Singer [SS98]. Section 3 shows the details of the minimum relative entropy interpretations of the corrective and totally corrective boosting algorithms. We use the minimum relative entropy interpretation in Section 4 for developing some geometric intuitions for the corrective update. The connection to on-line regression algorithms through (1.7) is pursued further in Section 5. Appendix A considers briefly some iterative methods for approximately solving the corrective and totally corrective updates. Appendix B discusses generalizations from relative entropy to other Bregman divergences.

## 2 The boosting algorithms

We take our framework for boosting from Schapire and Singer's work on confidence-rated boosting [SS98]. We consider classifying elements of an arbitrary set  $X$  into two classes, which we denote by  $-1$  and  $+1$ . Our training set consists of a set of  $m$  examples  $(x_i, y_i) \in X \times \{-1, 1\}$ , for  $i = 1, \dots, m$ . The interpretation of this input data is that for each *instance*  $x_i$ , the *label*  $y_i$  gives the correct classification of  $x_i$  according to some unknown target classifier. We allow our weak hypotheses to be arbitrary *confidence-rated* classifiers, i.e., mappings from  $X$  to  $[-1, 1]$ . Such a mapping  $h$  can be interpreted as predicting classifications for the elements of  $X$ , with  $\operatorname{sign}(h(x))$  the predicted classification and  $|h(x)|$  a confidence rating. The output of the master hypothesis will still be a strict classifier, i.e., a mapping from  $X$  to  $\{-1, 1\}$ . We assume that our weak learner receives a distribution over the training set as its input. We represent these distributions as vectors  $\mathbf{d}$  from the simplex  $P_m$ .

In boosting, we run the weak learner with  $T$  different distributions  $\mathbf{d}_t$ , for some suitable number  $T$  of rounds, and then combine the resulting weak hypotheses  $h_t$  by a weighted majority vote [Sch90, Fre95]. Figure 1 shows the details of AdaBoost and the corrective boosting algorithm, which differ only in the choice of the parameter  $\alpha_t$  regulating the amount of change at update  $t$ . Otherwise they both share the exponential form (1.1) of the updated distribution, and also the weighted majority form (2.3) of the master algorithm's hypothesis. As a practical point, a value  $\alpha_t$  such that  $\mathbf{d}_{t+1} \cdot \mathbf{u}_t = 0$  can be found by a line search except for the degenerate case in which all the components  $u_{t,i}$  have the same sign [SS98]. We give in Appendix A some very simple bounds for this line search.

Schapire and Singer noticed that for both AdaBoost and the corrective boosting algorithm the training error of the master hypothesis can be bounded by the product of the normalization factors as

$$\frac{1}{m} |\{i \mid H(x_i) \neq y_i\}| \leq \prod_{t=1}^T Z_t(\alpha_t). \quad (2.4)$$

Since  $Z_t(\alpha)$  is minimized when  $\alpha$  is such that  $\mathbf{d}_{t+1}(\alpha) \cdot \mathbf{u}_t = 0$ , the corrective update is motivated as a minimizer of this upper bound. For AdaBoost,  $Z_t(\alpha)$  is replaced by an upper bound (which is exact in the discrete case  $u_{t,i} \in \{-1, 1\}$ ) and then  $\alpha_t$  is chosen by minimizing this upper bound.

**Input** a set of  $m$  examples  $(x_i, y_i) \in X \times \{-1, 1\}$ ,  $i = 1, \dots, m$ .

**Initialize**  $d_{1,i} = 1/m$  for  $i = 1, \dots, m$ .

**Repeat** for  $t = 1, \dots, T$ :

- Call the weak learner with the distribution  $\mathbf{d}_t$  over the examples  $(x_i, y_i)$ ; let the resulting hypothesis be  $h_t$ .
- Choose a parameter  $\alpha_t \in \mathbf{R}$  as follows. Define  $\mathbf{u}_t \in [-1, 1]^m$  by  $u_{t,i} = y_i h_t(x_i)$ . Depending on the algorithm, use the following values.

**AdaBoost:** Choose

$$\alpha_t = \frac{1}{2} \ln \frac{1 + \mathbf{d}_t \cdot \mathbf{u}_t}{1 - \mathbf{d}_t \cdot \mathbf{u}_t} .$$

**Corrective:** Choose  $\alpha_t$  such that

$$\sum_{i=1}^m u_{t,i} d_{t,i} \exp(-\alpha_t u_{t,i}) = 0 .$$

- Update the distribution by  $\mathbf{d}_{t+1} = \mathbf{d}_{t+1}(\alpha_t)$  with

$$d_{t+1,i}(\alpha) = \frac{1}{Z_t(\alpha)} d_{t,i} \exp(-\alpha u_{t,i}) \quad (2.1)$$

and

$$Z_t(\alpha) = \sum_{i=1}^m d_{t,i} \exp(-\alpha u_{t,i}) . \quad (2.2)$$

**Output** the master hypothesis  $H$  defined by

$$H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right) . \quad (2.3)$$

Figure 1: AdaBoost and the corrective and totally corrective boosting algorithm

The requirement  $\mathbf{d}_{t+1} \cdot \mathbf{u}_t = 0$  can also be interpreted as requiring that the new distribution makes the current weak hypothesis totally uncorrelated with the training data. Intuitively, the weak learner is then forced to learn something new for  $h_{t+1}$ . Given this motivation, it could seem natural to consider a more general algorithm that enforces the constraint  $\mathbf{d}_{t+1} \cdot \mathbf{u}_q$  not only for  $q = t$  but also for all  $q < t$ . Satisfying all these constraints naturally requires more than one free variable.

Hence, instead of having at update  $t$  just one coefficient  $\alpha_t$  for the  $t$ th weak hypothesis, we take  $t$  coefficients  $\alpha_{t,q}$ , one for each past weak hypothesis  $h_q$ ,  $1 \leq q \leq t$ . Using this  $t$ -dimensional parameter vector  $\boldsymbol{\alpha}_t$  we now write the update as  $\mathbf{d}_{t+1} = \mathbf{d}_{t+1}(\boldsymbol{\alpha}_t)$  where

$$d_{t+1,i}(\boldsymbol{\alpha}) = \frac{1}{Z_t(\boldsymbol{\alpha})} d_{t,i} \exp \left( - \sum_{q=1}^t \alpha_q u_{q,i} \right) \quad (2.5)$$

and

$$Z_t(\boldsymbol{\alpha}) = \sum_{i=1}^m d_{t,i} \exp \left( - \sum_{q=1}^t \alpha_q u_{q,i} \right) . \quad (2.6)$$

We call the algorithm with the extended exponential form (2.5), with  $\boldsymbol{\alpha}_t$  chosen such that  $\mathbf{d}_{t+1} \cdot \mathbf{u}_q$  holds for all  $1 \leq q \leq t$ , the *totally corrective* algorithm. In the definition (2.3) of the master hypothesis we use the weights  $\alpha_t = \sum_{q=1}^T \alpha_{q,t}$ .

Schapire and Singer's proof of the error bound (2.4) generalizes easily to give

$$\frac{1}{m} |\{i \mid H(x_i) \neq y_i\}| \leq \prod_{t=1}^T Z_t(\boldsymbol{\alpha}_t)$$

for the totally corrective algorithm. Analogously with the corrective algorithm, the totally corrective algorithm chooses at update  $t$  the parameter vector  $\boldsymbol{\alpha}_t$  such that  $Z_t(\boldsymbol{\alpha}_t)$  is minimized. It is also easy to see that assuming the same set of weak hypotheses for the corrective and totally corrective algorithm, the bound  $\prod_t Z_t(\boldsymbol{\alpha}_t)$  for the totally corrective algorithm is no larger than the bound  $\prod_t Z_t(\alpha_t)$  for the corrective algorithm. (But of course we would not expect to get the same weak hypotheses with different distributions.)

It should be noted that there may not exist any  $\boldsymbol{\alpha}_t$  such that  $\mathbf{d}_{t+1} \cdot \mathbf{u}_q$  holds for all  $1 \leq q \leq t$ , and in any case finding such a vector  $\boldsymbol{\alpha}_t$  would be an  $t$ -dimensional optimization problem. We consider this issue briefly in Appendix A. However, we are not claiming that the totally corrective algorithm would necessarily be a practical learning algorithm. We introduce it here mainly as a theoretical comparison point for the corrective algorithm.

### 3 Boosting as relative entropy minimization

We saw in Section 2 various boosting algorithms whose updated distribution has exponential form can be seen as minimizing the factors  $Z_t(\boldsymbol{\alpha}_t)$  that appear in the bound (2.4). This minimization problem explains the choice of the values  $\alpha_t$ . The use of the exponential form (2.1) is essential. In particular, the proof of (2.4) uses the property that  $e^{-\alpha y h(x)} \leq 1$  if and only if  $\text{sign}(h(x)) \neq \text{sign}(y)$ . Thus the exponential gives a nice approximation to the discrete loss [SS98]. (See [FHT98] for more discussion).

We now suggest an alternative view, in which the corrective and totally corrective updates appear as solutions to constrained relative entropy minimization problems. The exponential form of the update is an immediate consequence of using the relative entropy as a measure of divergence between distributions. If the relative entropy was replaced by a different Bregman divergence, then the update would have another form. The values of the parameters  $\alpha_t$  are the Lagrange multiplier for enforcing the constraints. For simplicity, we show the details only for the corrective algorithm and then explain briefly how the results generalize to the totally corrective algorithm.

We view the corrective update as pursuing two conflicting goals. First, the updated distribution should be uncorrelated with the mistakes made by the previous weak hypothesis, i.e.  $\mathbf{d}_{t+1} \cdot \mathbf{u}_t = 0$ . Otherwise the updated distribution should stay closest to the last distribution so as to retain

changes made in previous updates and also resist overreacting to noise. The distance is measured by the relative entropy  $\Delta(\mathbf{d}_{t+1}, \mathbf{d}_t)$  defined in (1.2).

The following theorem is basically an application of standard duality techniques from convex optimization [Lue84]. Intuitively, a constrained minimization problem for  $\Delta(\mathbf{d}, \mathbf{d}_t)$  turns out to be equivalent to an unconstrained maximization problem for  $-\ln Z_t(\alpha)$ . A similar result, but with more emphasis on the special properties of the relative entropy, is given by Della Pietra et al. [DDL97]. The free variable  $\alpha$  of the unconstrained problem becomes a Lagrange multiplier in the constrained problem. Although the theorem is rather basic, we give the proof in complete detail for clarity.

**Theorem 1** Define  $\mathbf{d}_{t+1}(\alpha)$  and  $Z_t(\alpha)$  as in (2.1) and (2.2), and assume that  $\mathbf{d}_{t+1}(\alpha) \cdot \mathbf{u}_t = 0$  for some  $\alpha \in \mathbf{R}$ . Then

$$\min_{\substack{\mathbf{d} \in P_m \\ \mathbf{d} \cdot \mathbf{u}_t = 0}} \Delta(\mathbf{d}, \mathbf{d}_t) = \max_{\alpha \in \mathbf{R}} (-\ln Z_t(\alpha)) . \quad (3.1)$$

Further,

$$\operatorname{argmin}_{\substack{\mathbf{d} \in P_m \\ \mathbf{d} \cdot \mathbf{u}_t = 0}} \Delta(\mathbf{d}, \mathbf{d}_t) = \mathbf{d}_{t+1}(\alpha_t) \quad (3.2)$$

where

$$\alpha_t = \operatorname{argmax}_{\alpha \in \mathbf{R}} (-\ln Z_t(\alpha)) .$$

Recall that the corrective update is given by

$$\mathbf{d}_{t+1} = \mathbf{d}_{t+1}(\alpha_t) \quad \text{where } \alpha_t = \operatorname{argmax}_{\alpha \in \mathbf{R}} (-\ln Z_t(\alpha)) .$$

Equivalently  $\alpha_t$  is such that  $\mathbf{d}_{t+1}(\alpha_t) \cdot \mathbf{u}_t = 0$ . Hence, Theorem 1 gives a relative entropy interpretation both for the corrective update and the normalization factor  $Z_t(\alpha_t)$  in the error bound (2.4).

**Proof of Theorem 1** As discussed earlier, the assumption that  $\mathbf{d}_{t+1}(\alpha) \cdot \mathbf{u}_t = 0$  for some  $\alpha$  implies that this  $\alpha$  is actually the unique minimum point for  $Z_t(\alpha)$ . Hence, in particular, the maximum on the right-hand side of (3.1) is attained at  $\alpha = \alpha_t$  where  $\mathbf{d}_{t+1}(\alpha_t) \cdot \mathbf{u}_t = 0$ .

Consider now the constrained minimization on the left-hand side of (3.1). Define the Lagrangian

$$F_t(\mathbf{d}, \alpha) = \Delta(\mathbf{d}, \mathbf{d}_t) + \alpha \mathbf{d} \cdot \mathbf{u}_t .$$

The key step of the proof of (3.1) is the minimax equation

$$\min_{\mathbf{d} \in P_m} \max_{\alpha \in \mathbf{R}} F_t(\mathbf{d}, \alpha) = \max_{\alpha \in \mathbf{R}} \min_{\mathbf{d} \in P_m} F_t(\mathbf{d}, \alpha) . \quad (3.3)$$

Before going into the proof of (3.3), let us see how it gives the main equality (3.1).

First, it is clear that  $\max_{\alpha \in \mathbf{R}} F_t(\mathbf{d}, \alpha)$  is  $\Delta(\mathbf{d}, \mathbf{d}_t)$  if  $\mathbf{d} \cdot \mathbf{u}_t = 0$  and  $\infty$  otherwise. Therefore, we have

$$\min_{\substack{\mathbf{d} \in P_m \\ \mathbf{d} \cdot \mathbf{u}_t = 0}} \Delta(\mathbf{d}, \mathbf{d}_t) = \min_{\mathbf{d} \in P_m} \max_{\alpha \in \mathbf{R}} F_t(\mathbf{d}, \alpha) . \quad (3.4)$$

Now (3.3) gives us

$$\min_{\substack{\mathbf{d} \in P_m \\ \mathbf{d} \cdot \mathbf{u}_t = 0}} \Delta(\mathbf{d}, \mathbf{d}_t) = \max_{\alpha \in \mathbf{R}} \min_{\mathbf{d} \in P_m} F_t(\mathbf{d}, \alpha) .$$

To finish the proof of (3.1), it is now sufficient to show

$$\min_{\mathbf{d} \in P_m} F_t(\mathbf{d}, \alpha) = -\ln Z_t(\alpha) \quad (3.5)$$

for all  $\alpha \in \mathbf{R}$ . This last claim follows by direct substitution from the more specific result that the value  $F_t(\mathbf{d}, \alpha)$  for fixed  $\alpha$  is minimized when  $\mathbf{d} = \mathbf{d}_{t+1}(\alpha)$ :

$$\operatorname{argmin}_{\mathbf{d} \in P_m} F_t(\mathbf{d}, \alpha) = \mathbf{d}_{t+1}(\alpha) . \quad (3.6)$$

Equation (3.6) is obtained by a straightforward differentiation with an additional Lagrange coefficient to enforce  $\sum_i d_i = 1$ . Since  $\Delta(\mathbf{w}, \mathbf{w}_t)$  is convex in  $\mathbf{w}$ , the zero of the derivatives is the minimum point.

Thus, the proof of (3.3) remains. First, we have

$$\begin{aligned} \max_{\alpha \in \mathbf{R}} \min_{\mathbf{d} \in P_m} F_t(\mathbf{d}, \alpha) &= \min_{\mathbf{d} \in P_m} F_t(\mathbf{d}, \alpha_t) \\ &\leq \min_{\mathbf{d} \in P_m} \max_{\alpha \in \mathbf{R}} F_t(\mathbf{d}, \alpha) , \end{aligned}$$

where we used the fact shown above that the maximum of  $\min_{\mathbf{d} \in P_m} F_t(\mathbf{d}, \alpha)$  occurs at  $\alpha = \alpha_t$ . On the other hand, by applying (3.5) and the fact  $\mathbf{d}_{t+1}(\alpha_t) \cdot \mathbf{u}_t = 0$  we get

$$\begin{aligned} \max_{\alpha \in \mathbf{R}} \min_{\mathbf{d} \in P_m} F_t(\mathbf{d}, \alpha) &= F_t(\mathbf{d}_{t+1}(\alpha_t), \alpha_t) \\ &\geq \min_{\mathbf{d} \in P_m} F_t(\mathbf{d}, \alpha_t) . \end{aligned}$$

From (3.4) we now see that

$$\max_{\alpha \in \mathbf{R}} \min_{\mathbf{d} \in P_m} F_t(\mathbf{d}, \alpha) \geq \min_{\mathbf{d} \in P_m} \max_{\alpha \in \mathbf{R}} F_t(\mathbf{d}, \alpha) ,$$

and (3.3) follows.

The proof of (3.3) implies that actually

$$F_t(\mathbf{d}_{t+1}(\alpha_t), \alpha_t) = \min_{\mathbf{d} \in P_m} F_t(\mathbf{d}, \alpha_t) ,$$

so  $\Delta(\mathbf{d}_{t+1}(\alpha_t), \mathbf{d}_t) = \min_{\mathbf{d} \in P_m} \Delta(\mathbf{d}, \mathbf{d}_t)$ . Since  $\Delta(\mathbf{d}, \mathbf{d}_t)$  is convex in  $\mathbf{d}$  and thus has a unique minimum, (3.2) follows.  $\square$

Consider now the totally corrective algorithm. The only difference to the corrective one is that now for  $\mathbf{d}_{t+1}$  we have  $t$  constraints instead of just one. Thus, let

$$C = \{ \mathbf{d} \in P_m \mid \mathbf{d} \cdot \mathbf{u}_q = 0 \text{ for } 1 \leq q \leq t \} .$$

An argument similar to the proof of Theorem 1 shows that

$$\min_{\mathbf{d} \in C} (\Delta(\mathbf{d}_{t+1}, \mathbf{d}_t)) = \max_{\alpha \in \mathbf{R}^t} (-\ln Z_t(\alpha)) ,$$

and further  $\alpha_t = \operatorname{argmax}_{\alpha} (-\ln Z_t(\alpha))$  and

$$\mathbf{d}_{t+1}(\alpha_t) = \operatorname{argmin}_{\mathbf{d} \in C} \Delta(\mathbf{d}, \mathbf{d}_t) .$$

## 4 Geometric interpretations for boosting

We next show some properties of the corrective update that follow naturally when we interpret (3.2) in a geometric fashion: the new distribution  $\mathbf{d}_{t+1}$  is obtained by projecting the old distribution  $\mathbf{d}_t$  onto the hyperplane  $U_t = \{ \mathbf{d} \mid \mathbf{d} \cdot \mathbf{u}_t = 0 \}$ . Here the projection of the point  $\mathbf{d}_t$  to the plane  $U_t$  is defined as the point  $\mathbf{d}$  on the plane that is closest to the starting point  $\mathbf{d}_t$ . The relative entropy  $\Delta(\mathbf{d}, \mathbf{d}_t)$  is used as our measure of distance. Although geometric

metaphors are very illuminating here, it must be remembered that the relative entropy is not a metric and we have to be careful when we use our intuitions about distance. The ideas sketched here can be applied directly to the totally corrective update by replacing the hyperplane  $U_t$  by the intersection of  $t$  hyperplanes  $U_q, 1 \leq q \leq t$ .

Assume now that there is at least one  $\alpha \in \mathbf{R}$  such that the distribution  $\mathbf{d}_{t+1}(\alpha)$  in the exponential form (2.1) satisfies  $\mathbf{d}_{t+1}(\alpha) \cdot \mathbf{u}_t = 0$ . As discussed earlier, this is a reasonable assumption in the boosting setting, and then actually there is a unique value  $\alpha_t$  such that  $\mathbf{d}_{t+1}(\alpha) \cdot \mathbf{u}_t = 0$  holds if and only if  $\alpha = \alpha_t$ . Now it is easy to show that all distributions  $\mathbf{d}_{t+1}(\alpha)$  in exponential form (2.1) project to the same point  $\mathbf{d}_{t+1}$  on the hyperplane  $U_t$ , i.e., for any  $\alpha \in \mathbf{R}$  we have

$$\operatorname{argmin}_{\mathbf{d} \in P_m \cap U_t} \Delta(\mathbf{d}, \mathbf{d}_t) = \operatorname{argmin}_{\mathbf{d} \in P_m \cap U_t} \Delta(\mathbf{d}, \mathbf{d}_{t+1}(\alpha)) = \mathbf{d}_{t+1} .$$

Thus, for any distribution  $\mathbf{d}$  on the curve  $\{\mathbf{d}_{t+1}(\alpha) \mid \alpha \in \mathbf{R}\}$ , the projection of  $\mathbf{d}$  on the hyperplane  $U_t$  is the unique point where the curve intersects the hyperplane.

Another interesting property is analogous to the Pythagorean Theorem. Consider  $\mathbf{d}_t$  and its projection  $\mathbf{d}_{t+1}$  onto the hyperplane  $U_t$ . Then for any  $\mathbf{d}_*$  on the plane  $U_t$ ,

$$\Delta(\mathbf{d}_*, \mathbf{d}_t) = \Delta(\mathbf{d}_*, \mathbf{d}_{t+1}) + \Delta(\mathbf{d}_{t+1}, \mathbf{d}_t) . \quad (4.1)$$

It is easy to check this using the properties noted in Section 3. If we replace in (4.1) the relative entropy  $\Delta(\mathbf{d}, \mathbf{d}')$  by the squared Euclidean distance  $\|\mathbf{d} - \mathbf{d}'\|_2^2$  and  $\mathbf{d}_{t+1}$  by the usual Euclidean projection of  $\mathbf{d}_t$  onto  $U_t$ , then (4.1) becomes the familiar Pythagorean Theorem. This property of minimum distance projections onto sets defined by linear constraints holds even more generally for all Bregman divergences [Bre67]. For applications in on-line learning theory, see [HW98].

## 5 Boosting in a regression framework

It has been pointed out [ROM98] that on highly noisy training sets, AdaBoost may tend to overfit. Considering this, the strict constraint  $\mathbf{d}_{t+1} \cdot \mathbf{u}_t = 0$  of the corrective algorithm, and the even tighter constraint of the totally corrective algorithm, seems a little uncautious. As a possible means of avoiding this problem, we suggest replacing the constrained minimization problem (1.4) by

$$\mathbf{d}_{t+1} = \operatorname{argmin}_{\mathbf{d} \in P_m} (\Delta(\mathbf{d}, \mathbf{d}_t) + \eta_t L(\mathbf{d} \cdot \mathbf{u}_t)) , \quad (5.1)$$

where  $L$  is some loss function such that  $L(z)$  is minimized for  $z = 0$  and  $\eta_t > 0$  is a parameter controlling the trade-off between minimizing the loss and staying close to the previous distribution. The corrective update is obtained from (5.1) in the limit when  $\eta_t$  approaches infinity.

This framework for deriving parameter updates was introduced by Kivinen and Warmuth [KW97] in the context of on-line linear regression. The solution to (5.1) satisfies

$$d_{t+1,i} = \frac{1}{N_t} d_{t,i} \exp(-\eta_t L'(\mathbf{d}_{t+1} \cdot \mathbf{u}_t) u_{t,i}) , \quad (5.2)$$

where  $L'$  is the derivative of  $L$  and  $N_t$  a normalization factor. For general  $L$ , (5.2) cannot be solved in closed form because of how  $\mathbf{d}_{t+1}$  appears on the right-hand side. Therefore,

Kivinen and Warmuth suggest approximating the derivative  $L'(\mathbf{d}_{t+1} \cdot \mathbf{u}_t)$  on the right-hand side by its old value  $L'(\mathbf{d}_t \cdot \mathbf{u}_t)$ , which is a reasonable approximation at least for small values of  $\eta_t$ . This results in the Exponentiated Gradient update [KW97]

$$d_{t+1,i} = \frac{1}{M_t} d_{t,i} \exp(-\eta_t L'(\mathbf{d}_t \cdot \mathbf{u}_t) u_{t,i}) , \quad (5.3)$$

where again  $M_t$  is a normalization factor. The trade-off parameter  $\eta_t$  can be interpreted as a learning rate. This approach can be generalized by replacing the relative entropy in (5.1) by any Bregman divergence.

It is here particularly interesting to apply (5.1) with  $L(z) = L_{\text{ent}}(z, 0)$ , where  $L_{\text{ent}}(z, \hat{z})$  for  $z, \hat{z} \in [-1, 1]$  is the usual entropic loss

$$L_{\text{ent}}(z, \hat{z}) = \frac{1-z}{2} \ln \frac{1-z}{1-\hat{z}} + \frac{1+z}{2} \ln \frac{1+z}{1+\hat{z}} .$$

Then  $L'(z) = \frac{1}{2} \ln((1+z)/(1-z))$ . By comparing (5.3) with (1.1) and recalling the value  $\alpha_t = \ln((1 + \mathbf{d}_t \cdot \mathbf{u}_t)/(1 - \mathbf{d}_t \cdot \mathbf{u}_t))/2$  used by AdaBoost, we see that AdaBoost can be interpreted as Exponentiated Gradient with the entropic loss function and learning rate  $\eta_t = 1$ .

## 6 Conclusions

We have considered the update step of the standard boosting algorithms as a constrained relative entropy minimizer, or alternatively as projection with respect to the relative entropy distance measure. We hope that our simple observations will be useful in designing better boosting algorithms. Many of the basic properties of the relative entropy are shared more generally by all Bregman divergences [Bre67]. It would be interesting to see whether some other divergences might lead to useful boosting procedures. Some updates motivated by different Bregman divergences are briefly discussed in Appendix B, but without any results on the training error of the resulting boosting procedure. Perhaps the GeoLev procedure [DH99] could be related to projections with respect to the squared Euclidean distances. Note that the relative entropy is a special divergence in that it is defined on the simplex  $P_m$  and this is the natural domain for boosting. For other divergence, an additional projection onto  $P_m$  would be needed, but this is not necessarily a problem. See [HW98] for examples of using projections onto arbitrary convex sets in a regression setting.

Another interesting subject for further study is noncorrective updates motivated analogously to regression algorithms as in Section 5. Hopefully, they will provide a means for making the boosting algorithms less prone to overfitting.

### Acknowledgments

We wish to thank Nigel Duffy, David Helmbold, Mark Herbster, John Lafferty, and Sandra Panizza for inspiring discussions.

## A Finding the corrective parameter values

### A.1 The corrective algorithm

As Schapire and Singer [SS98] have observed, there is a unique value  $\alpha_t$  such that  $\mathbf{d}_{t+1}(\alpha_t) \cdot \mathbf{u}_t = 0$  for the corrective

algorithm, unless all the components  $u_{t,i}$  have the same sign. Such value can be found by a simple line search. However, if the absolute values  $|u_{t,i}|$  can be arbitrarily small, then the value  $\alpha_t$  can be arbitrarily large. We now give a crude estimate of the range we need to search.

Thus, consider an updated distribution  $\mathbf{d}_{t+1}(\alpha)$  as in (2.1). We are looking for some upper and lower bounds for the value  $\alpha_t$  such that  $\mathbf{d}_{t+1}(\alpha_t) \cdot \mathbf{u}_t = 0$ . Assume that  $\mathbf{d}_t \cdot \mathbf{u}_t > 0$ . (If  $\mathbf{d}_t \cdot \mathbf{u}_t < 0$ , replace  $\mathbf{u}_t$  with  $-\mathbf{u}_t$ . If  $\mathbf{d}_t \cdot \mathbf{u}_t = 0$ , then the solution is  $\alpha_t = 0$ .) We try to find a value  $\beta > 0$  such that  $\mathbf{d}_{t+1}(\beta) \cdot \mathbf{u}_t < 0$ , so  $0 < \alpha_t < \beta$ . Let  $k$  be such that  $u_{t,k} = \min\{u_{t,i} \mid d_{t,i} \neq 0\}$ . Define  $P = \{i \mid u_{t,i} > 0, d_{t,i} \neq 0\}$ , and let  $j$  be such that  $u_{t,j} = \min_{i \in P} u_{t,i}$ . We are assuming  $P \neq \emptyset$  and  $u_{t,i} < 0$ . Consider any  $\alpha$  such that  $\alpha \geq 1/u_{t,j}$ . The function  $f_\alpha$  given by  $f_\alpha(x) = xe^{-\alpha x}$  has  $f'_\alpha(x) < 0$  when  $x > 1/\alpha$ . Hence, in particular, we have  $f_\alpha(u_{t,j}) \geq f_\alpha(u_{t,i})$  for all  $i \in P$ . We can therefore write

$$\begin{aligned} Z_t(\alpha) \mathbf{d}_{t+1}(\alpha) \cdot \mathbf{u}_t &= d_{t,k} u_{t,k} e^{-\alpha u_{t,k}} + \sum_{i \neq k} d_{t,i} f_\alpha(u_{t,i}) \\ &\leq d_{t,k} u_{t,k} e^{-\alpha u_{t,k}} + \sum_{i \in P} d_{t,i} f_\alpha(u_{t,i}) \\ &\leq d_{t,k} u_{t,k} e^{-\alpha u_{t,k}} + f_\alpha(u_{t,j}) \sum_{i \in P} d_{t,i} \\ &= e^{-\alpha u_{t,k}} \cdot \left( d_{t,k} u_{t,k} + u_{t,j} e^{-\alpha(u_{t,j} - u_{t,k})} \sum_{i \in P} d_{t,i} \right), \end{aligned}$$

from which we see that  $\mathbf{d}_{t+1}(\alpha) \cdot \mathbf{u}_t \leq 0$  holds for  $\alpha \geq \alpha^*$  where

$$\alpha^* = \frac{1}{u_{t,j} - u_{t,k}} \ln \frac{u_{t,j} \sum_{i \in P} d_{t,i}}{-d_{t,k} u_{t,k}}.$$

Hence, we can take  $\beta = \max\{1/u_{t,j}, \alpha^*\}$ .

We now can do a binary search for  $\alpha_t$  in the interval  $(0, \beta)$ . Alternatively we can start a search for a point  $\gamma$  s.t.  $\mathbf{d}_{t+1}(\gamma) \cdot \mathbf{u}_t < 0$  by starting with  $\gamma = 1$  and then doubling  $\gamma$  iteratively. This iterative procedure will terminate quickly because  $\gamma$  can never be much larger than  $\beta$ . Once we found a  $\gamma$  with the property we want we can start the binary search for  $\alpha_t$  in a small region.

## A.2 The totally corrective algorithm

As we have seen, the problem of finding the values for the parameters  $\alpha$  for the corrective and totally corrective algorithm is a relative entropy minimization problem. Obviously, there are a large number of general optimization algorithms than could be used to solve such problems. We present here an analysis of an iterative algorithm that is specifically tailored to the boosting case. The analysis is very close the bound given in [LLW92] on the number of iterations required for finding an approximate solution to a system of equations.

Of course, we need to assume that there is at least one distribution  $\mathbf{d}$  that satisfies all the constraints  $\mathbf{d} \cdot \mathbf{u}_t = 0$ . This may not be the case. Consider for example  $\mathbf{u}_1 = (-1/3, 1/2, 0, 0)$ ,  $\mathbf{u}_2 = (0, 0, 1/2, -1/3)$ , and  $\mathbf{u}_3 =$

$(0, 1/2, 0, 1/3)$ . Then it is easy to see that there are vectors  $\mathbf{w}$  that satisfy the three constraints  $\mathbf{w} \cdot \mathbf{u}_t = 0$ , but any such vector has both positive and negative components and cannot therefore be normalized into a distribution.

Our approach is quite similar to so-called row-action optimization methods [Bre67, CL81]. We wish to find a distribution  $\mathbf{d}_{t+1} \in P_m$  that minimizes the relative entropy  $\Delta(\mathbf{d}_{t+1}, \mathbf{d}_t)$  subject to  $t$  constraints  $\mathbf{d}_{t+1} \cdot \mathbf{u}_q = 0$  for  $1 \leq q \leq t$ . We define a sequence  $\hat{\mathbf{d}}_j$ ,  $j = 1, 2, \dots$ , as follows. We start with  $\hat{\mathbf{d}}_1 = \mathbf{d}_t$ . Then at step  $j$  pick one constraint that is not satisfied with the current distribution  $\hat{\mathbf{d}}_j$ , and define  $\hat{\mathbf{d}}_{j+1}$  be the projection of  $\hat{\mathbf{d}}_j$  onto the hyperplane defined by that constraint. The limit point to which the sequence  $\hat{\mathbf{d}}_j$  converges must then satisfy all the constraints. Further, because of our choice of starting point and because of properties of projections such as discussed in Section 4, this limit point also minimized the relative entropy from  $\mathbf{d}_t$ .

As a minor change to the above outline, we do not here do the projections exactly but are satisfied with the AdaBoost update step that does an approximate projection, as discussed earlier. Also notice that in the totally corrective algorithm the correct distribution  $\mathbf{d}_t$  has been obtained from the initial distribution  $\mathbf{d}_1$  by (perhaps approximate) entropy projections with respect to constraints  $\mathbf{d}_t \cdot \mathbf{u}_q = 0$ ,  $q < t$ . Therefore, by the properties of projections in Section 4, the constrained relative entropy minimization problem does not really change if we take  $\mathbf{d}_1$  instead of  $\mathbf{d}_t$  as the starting point. Thus, we have an iterative procedure that starts with  $\hat{\mathbf{d}}_1 = \mathbf{d}_1$  and then for  $j = 1, 2, \dots$  repeats the following:

1. Let  $q_j$  be such that  $|\hat{\mathbf{d}}_j \cdot \mathbf{u}_{q_j}|$  is maximized.
2. Let  $\hat{\alpha}_j = \ln((1 + \hat{\mathbf{d}}_j \cdot \mathbf{u}_{q_j}) / (1 - \hat{\mathbf{d}}_j \cdot \mathbf{u}_{q_j})) / 2$ .
3. Define  $\hat{\mathbf{d}}_{j+1}$  by

$$\hat{\mathbf{d}}_{j+1,i} = \frac{1}{\hat{Z}_j} \hat{\mathbf{d}}_{j,i} \exp(-\hat{\alpha}_j \mathbf{u}_{q_j,i})$$

where  $\hat{Z}_j$  is the normalization factor.

Then  $\hat{Z}_j \leq (1 - (\hat{\mathbf{d}}_j \cdot \mathbf{u}_{q_j})^2)^{1/2}$ , as shown by the usual analysis of AdaBoost [FS97b]. Fix a distribution vector  $\mathbf{d}_* \in P_m$ . It follows straight from the definitions that

$$\Delta(\mathbf{d}_*, \hat{\mathbf{d}}_j) - \Delta(\mathbf{d}_*, \hat{\mathbf{d}}_{j+1}) = -\hat{\alpha}_j \mathbf{d}_* \cdot \mathbf{u}_{q_j} - \ln \hat{Z}_j.$$

Fix now some parameter  $0 < \gamma < 1$ , and let  $J$  be the largest index such that  $|\hat{\mathbf{d}}_j \cdot \mathbf{u}_{q_j}| \geq \gamma$  holds for all  $j \leq J$ . Then for  $j \leq J$  we have

$$\begin{aligned} -\ln \hat{Z}_j &\geq -(1/2) \ln(1 - (\hat{\mathbf{d}}_j \cdot \mathbf{u}_{q_j})^2) \\ &\geq -(1/2) \ln(1 - \gamma^2), \end{aligned}$$

so in particular  $\sum_{j=1}^J (-\ln \hat{Z}_j) \geq -(J/2) \ln(1 - \gamma^2)$ . On the other hand, assume now that  $\mathbf{d}_*$  can be chosen such that  $\mathbf{d}_* \cdot \mathbf{u}_q = 0$  holds for all  $1 \leq q \leq t$ . (That is, there is a

solution  $\mathbf{d}_*$  that satisfies the constraints.) Then

$$\begin{aligned} \sum_{j=1}^J (-\ln \widehat{Z}_j) &= \sum_{j=1}^J \left( \Delta(\mathbf{d}_*, \widehat{\mathbf{d}}_j) - \Delta(\mathbf{d}_*, \widehat{\mathbf{d}}_{j+1}) \right) \\ &= \Delta(\mathbf{d}_*, \widehat{\mathbf{d}}_1) - \Delta(\mathbf{d}_*, \widehat{\mathbf{d}}_{J+1}) \\ &\leq \ln m \end{aligned}$$

if we in addition choose  $\mathbf{d}_1 = \widehat{\mathbf{d}}_1$  to be the uniform start vector. Hence, we have

$$J \leq \frac{2 \ln m}{-\ln(1 - \gamma^2)} \leq \frac{2 \ln m}{\gamma^2} .$$

This is an upper bound for the number of iteration rounds before all the dot products  $\widehat{\mathbf{d}}_j \cdot \mathbf{u}_q$  get smaller than  $\gamma$  in absolute value.

The above iterative method can also be used in the case where there is just one constraint (i.e., the case of the corrective update). Thus iterating AdaBoost on the last weak hypothesis can be used to find the corrective update. This is a simple alternate to the binary search method discussed in the previous section.

## B Other distance measures

As we have mentioned above, relative entropy is a special case of *Bregman divergences* [Bre67], and much of the discussion about the boosting update and its motivation applies directly to the case of general Bregman divergences. In this section we explain the connection of the minimax results of Theorem 1 to the general duality properties of constrained convex optimization problems. Similar duality properties have been analysed in parallel work by Lafferty [Laf99]. Notice that these ideas only generalize the motivation of boosting as relative entropy minimization. The training error bounds for boosting [FS97b, SS98] are based on the specific way the exponential function appears in the update that minimizes relative entropy, and we know of no way of generalizing this to updates minimizing other Bregman divergences.

The discussion here is on a general level, and we omit regularity conditions such as having the optimal solution lie in the interior of the feasible region. To obtain a rigorous proof, such details would need to be considered, but it seems easier to do this individually for each Bregman divergence we wish to consider (like we did in Theorem 1 for the relative entropy) rather than to try to obtain general necessary and sufficient conditions. For a more complete treatment of duality in convex optimization, see standard textbooks such as Luenberger [Lue84, pp. 396–401] or Bazaraa et al. [BSS93, pp. 199–210].

Consider now a continuously differentiable strictly convex function  $F$  from some convex set  $X \subseteq \mathbf{R}^m$  to  $\mathbf{R}$ . Hence, the gradient  $\nabla F = f$  is a one-to-one mapping from  $\mathbf{R}^m$  to  $\mathbf{R}^m$ . We define the Bregman divergence  $\Delta_F$  for vectors  $\mathbf{w}$  and  $\tilde{\mathbf{w}}$  by

$$\Delta_F(\tilde{\mathbf{w}}, \mathbf{w}) = F(\tilde{\mathbf{w}}) - F(\mathbf{w}) - (\tilde{\mathbf{w}} - \mathbf{w}) \cdot f(\mathbf{w}) . \quad (\text{B.1})$$

Thus  $\Delta_F(\tilde{\mathbf{w}}, \mathbf{w})$  is the difference between  $F(\tilde{\mathbf{w}})$  and its approximation based on the first order Taylor polynomial of  $F$  around  $\mathbf{w}$ . Since  $F$  is strictly convex, this difference is strictly positive for  $\tilde{\mathbf{w}} \neq \mathbf{w}$ .

Let us now consider a generalized boosting update, which we obtain by replacing the relative entropy by an arbitrary Bregman divergence in (1.4):

$$\mathbf{w}_{t+1} = \underset{\mathbf{w}}{\operatorname{argmin}} \Delta_F(\mathbf{w}, \mathbf{w}_t) \quad \text{subject to } \mathbf{w} \cdot \mathbf{u}_t = 0 . \quad (\text{B.2})$$

Notice that we have omitted the constraint that the weights must be in the probability simplex and, with this in mind, use the symbol  $\mathbf{w}$  instead of  $\mathbf{d}$ . If one wants to keep the weights as probability vectors, which of course is needed in the standard boosting scenario, one can enforce the constraint  $\mathbf{w} \in P_m$  by the usual method of Lagrange multipliers. (This involves one multiplier for the constraint  $\sum_i d_i = 1$  and  $m$  multipliers for the constraints  $d_i \geq 0$ .) Another possibility would be to first obtain a solution  $\tilde{\mathbf{w}}$  to the minimization problem ignoring the constraint  $\mathbf{w} \in P_m$  and then obtain the final solution  $\mathbf{w}_{t+1} \in P_m$  as the projection of  $\tilde{\mathbf{w}}$  into  $P_m$  with respect to  $\Delta_F$ . This method was used by Herbert and Warmuth [HW98] in a regression setting. It is also possible to entirely ignore the interpretation of the weights as probabilities and apply the boosting framework simply as a certain kind of a parameter fitting procedure. This is the approach taken by Friedman et al. [FHT98], and also the method suggested by Della Pietra et al. in an earlier work [DDL97] (see Lafferty [Laf99] for later developments). Of course, if the weights do not represent a probability distribution, the weak learner must somehow be able to use the weights directly instead of via sampling. Because of this variety of possibilities, we concentrate here on properties of the update (B.2) on a level that leaves our position with respect to the constraint  $\mathbf{w} \in P_m$  open.

To solve (B.2), introduce now the Lagrangian by

$$U_t(\mathbf{w}, \alpha) = \Delta_F(\mathbf{w}, \mathbf{w}_t) + \alpha \mathbf{w} \cdot \mathbf{u}_t . \quad (\text{B.3})$$

First we wish to solve for a fixed  $\alpha$  the value

$$\mathbf{w}_{t+1}(\alpha) = \underset{\mathbf{w}}{\operatorname{argmin}} U_t(\mathbf{w}, \alpha) . \quad (\text{B.4})$$

By substituting the definition (B.1) into (B.3) and differentiating with respect to  $\mathbf{w}$  we see that

$$f(\mathbf{w}_{t+1}(\alpha)) = f(\mathbf{w}_t) - \alpha \mathbf{u}_t . \quad (\text{B.5})$$

Notice that since  $f$  is one-to-one, this determines  $\mathbf{w}_{t+1}(\alpha)$  uniquely. Now, in particular, the solution to (B.2) is given by  $\mathbf{w}_{t+1} = \mathbf{w}_{t+1}(\alpha_t)$  where  $\alpha_t$  is such that the constraint  $\mathbf{w}_{t+1}(\alpha) \cdot \mathbf{u}_t = 0$  is satisfied.

As we noticed,  $f$  is one-to-one; let  $g$  be the inverse of  $f$ , so we can write  $\mathbf{w} = g(\boldsymbol{\theta})$  for any  $\boldsymbol{\theta} = f(\mathbf{w})$ . We can then alternatively give the solution as an additive update in the  $\boldsymbol{\theta}$  parameters: we have  $\mathbf{w}_t = g(\boldsymbol{\theta}_t)$  where  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_{t+1}(\alpha_t)$  for

$$\boldsymbol{\theta}_{t+1}(\alpha) = \boldsymbol{\theta}_t - \alpha \mathbf{u}_t . \quad (\text{B.6})$$

To see the connection to the original boosting update, consider  $F(\mathbf{w}) = \sum_{i=1}^m (w_i \ln w_i - w_i)$ , for which the gradient is given by  $f_i(\mathbf{w}) = \ln w_i$  and the Bregman divergence is the *unnormalized relative entropy*

$$\Delta_F(\mathbf{w}, \mathbf{w}_t) = \sum_{i=1}^m \left( w_i \ln \frac{w_i}{w_{t,i}} - w_i + w_{t,i} \right) . \quad (\text{B.7})$$



The update corresponding to this divergence satisfies

$$\ln w_{t+1,i} = \ln w_{t,i} - \alpha_t u_{t,i} ,$$

which differs from the boosting update

$$\ln w_{t+1,i} = \ln w_{t,i} - \alpha_t u_{t,i} - \ln Z_t$$

only by the normalization. As we mentioned, the normalization constraint  $\sum_i w_i = 1$  can be included into the problem (B.2) via an additional Lagrange multiplier in (B.3). The positivity constraints  $w_i \geq 0$  turn out to hold automatically for the special case of unnormalized relative entropy.

Consider now the dual problem. For the Lagrangian  $U_t$  we define the usual dual function  $Q_t$  by

$$Q_t(\alpha) = \min_{\mathbf{w}} U_t(\mathbf{w}, \alpha) . \quad (\text{B.8})$$

Basic duality results (see, e.g., [BSS93, Theorem 6.2.4]) now imply that the constrained problem of minimizing  $\Delta_F(\mathbf{w}, \mathbf{w}_t)$  subject to  $\mathbf{w} \cdot \mathbf{u}_t = 0$  is equivalent to the unconstrained problem of maximizing  $Q_t(\alpha)$ . To be more precise, notice first that  $\Delta_F(\mathbf{w}, \mathbf{w}_t)$  is convex in  $\mathbf{w}$ , and  $Q_t(\alpha)$  can easily be shown to be concave in  $\alpha$ , so they have a unique minimum and maximum point, respectively. These are now known to give the same value, i.e.,

$$\min_{\mathbf{w} \cdot \mathbf{u}_t = 0} \Delta_F(\mathbf{w}, \mathbf{w}_t) = \max_{\alpha} Q_t(\alpha) . \quad (\text{B.9})$$

Further, the minimum and maximum points correspond to each other [BSS93, Theorem 6.2.5] in the sense that

$$\operatorname{argmin}_{\mathbf{w} \cdot \mathbf{u}_t = 0} \Delta_F(\mathbf{w}, \mathbf{w}_t) = \mathbf{w}_{t+1}(\alpha_t)$$

where  $\mathbf{w}_{t+1}(\alpha)$  is as in (B.4) and

$$\alpha_t = \operatorname{argmax}_{\alpha} Q_t(\alpha) . \quad (\text{B.10})$$

To make more use of this, let us write  $Q_t(\alpha)$  out in a more explicit form. For this, it is useful to introduce the *convex conjugate* of  $F$  [Roc70]. This is the function  $G$  that satisfies

$$G(\boldsymbol{\theta}) + F(\mathbf{w}) = \boldsymbol{\theta} \cdot \mathbf{w} \quad (\text{B.11})$$

for  $\boldsymbol{\theta} = f(\mathbf{w})$ . From the convexity of  $F$ , we know that  $G$  is well-defined and convex. Further, by differentiating the definition (B.11) we see directly the gradient  $g = \nabla G$  is the inverse of the gradient of  $F$ , i.e.,  $g(\boldsymbol{\theta}) = \mathbf{w}$  when  $f(\mathbf{w}) = \boldsymbol{\theta}$ . By substituting (B.11) into the definition (B.1) we see the connection [Ama85, AW99]

$$\Delta_G(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) = \Delta_F(\tilde{\mathbf{w}}, \mathbf{w})$$

(notice the change in the order of variables). In particular, from (B.6) we get

$$\begin{aligned} Q_t(\alpha) &= \Delta_F(\mathbf{w}_{t+1}(\alpha), \mathbf{w}_t) + \alpha \mathbf{w}_{t+1}(\alpha) \cdot \mathbf{u}_t \\ &= \Delta_G(\boldsymbol{\theta}_t, \boldsymbol{\theta}_{t+1}(\alpha)) + \alpha \mathbf{w}_{t+1}(\alpha) \cdot \mathbf{u}_t \\ &= G(\boldsymbol{\theta}_t) - G(\boldsymbol{\theta}_{t+1}(\alpha)) \\ &\quad - (\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1}(\alpha)) \cdot \mathbf{w}_{t+1}(\alpha) \\ &\quad + \alpha \mathbf{w}_{t+1}(\alpha) \cdot \mathbf{u}_t \\ &= G(\boldsymbol{\theta}_t) - G(\boldsymbol{\theta}_t - \alpha \mathbf{u}_t) . \end{aligned}$$

Hence, (B.9) becomes

$$\min_{\mathbf{w} \cdot \mathbf{u}_t = 0} \Delta_F(\mathbf{w}, \mathbf{w}_t) = \max_{\alpha} (G(\boldsymbol{\theta}_t) - G(\boldsymbol{\theta}_t - \alpha \mathbf{u}_t)) ,$$

and (B.10) gives

$$\alpha_t = \operatorname{argmax}_{\alpha} (-G(\boldsymbol{\theta}_t - \alpha \mathbf{u}_t)) .$$

As a simple example, consider the unnormalized relative entropy (B.7), which is the Bregman divergence for the convex function  $F = \sum_{i=1}^m (w_i \ln w_i - w_i)$  defined in  $X = \mathbf{R}_+^m$ . The gradient  $f$  now is given by  $f_i(\mathbf{w}) = \ln w_i$ , so for its inverse  $g$  we get  $g(\boldsymbol{\theta}) = e^{\boldsymbol{\theta}}$ . Clearly  $g = \nabla G$  for  $G(\boldsymbol{\theta}) = \sum_{i=1}^m e^{\theta_i}$ , and indeed this combination of  $F$  and  $G$  satisfies the condition (B.11). The update we get from this divergence is then  $w_{t+1,i} = w_{t,i} \exp(-\alpha_t u_{t,i})$ , and

$$\begin{aligned} \min_{\mathbf{w} \cdot \mathbf{u}_t = 0} \Delta_F(\mathbf{w}, \mathbf{w}_t) &= \Delta_F(\mathbf{w}_{t+1}, \mathbf{w}_t) \\ &= \max_{\alpha} (G(\boldsymbol{\theta}_t) - G(\boldsymbol{\theta}_t - \alpha \mathbf{u}_t)) \\ &= \max_{\alpha} \left( \sum_{i=1}^m w_{t,i} (1 - \exp(-\alpha u_{t,i})) \right) . \end{aligned}$$

If we constrain the weights to satisfy  $\sum_i w_i = 1$ , then of course the unnormalized relative entropy becomes the usual relative entropy, but the above derivation for the value  $\Delta_F(\mathbf{w}_{t+1}, \mathbf{w}_t)$  in terms of  $G$  becomes invalid. To see the algorithm for the relative entropy, i.e., the corrective boosting algorithm, in this light, first notice that the boosting update (1.1) can be written as  $\mathbf{w}_t = g(\boldsymbol{\theta}_t)$  with  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_{t+1}(\alpha)$  when  $g$  is the *softmax function*

$$g_i(\boldsymbol{\theta}) = \frac{e^{\theta_i}}{\sum_{j=1}^m e^{\theta_j}} .$$

Then  $g = \nabla G$  for

$$G(\boldsymbol{\theta}) = \ln \left( \sum_{i=1}^m e^{\theta_i} \right) .$$

We then get for the relative entropy  $\Delta(\mathbf{w}_{t+1}, \mathbf{w}_t)$  the value  $G(\boldsymbol{\theta}_t) - G(\boldsymbol{\theta}_{t+1})$  as expected. Unfortunately,  $G$  is not *strictly* convex, and accordingly  $g$  is not one-to-one, so the derivation given above is not valid without modifications. The simplest way to resolve this is to represent the weights  $\mathbf{w} \in P_m$  by lower-dimensional weights  $\mathbf{w}' \in [0, 1]^{m-1}$  with  $w_i = w'_i$  for  $1 \leq i \leq m-1$  and  $w_m = 1 - \sum_{i=1}^{m-1} w'_i$ . We omit the details of this reduction, but the result is that we also get Theorem 1 for the usual relative entropy as a special case of the derivation given here.

Another related divergence is the sum of binary relative entropies used by Bylander [Byl97] to analyse on-line linear regression. This divergence is defined for vectors in  $[0, 1]^m$ , with

$$\Delta_F(\tilde{\mathbf{w}}, \mathbf{w}) = \sum_{i=1}^m \left( \tilde{w}_i \ln \frac{\tilde{w}_i}{w_i} + (1 - \tilde{w}_i) \ln \frac{1 - \tilde{w}_i}{1 - w_i} \right)$$

for  $F(\mathbf{w}) = \sum_{i=1}^m (w_i \ln w_i + (1 - w_i) \ln(1 - w_i))$ . The gradient is now given by  $f_i(\mathbf{w}) = \ln(w_i/(1 - w_i))$ , from which we get  $g_i(\boldsymbol{\theta}) = e^{\theta_i}/(1 + e^{\theta_i})$ . The update then becomes

$$w_{t+1,i} = \frac{w_{t,i} e^{-\alpha_t u_{t,i}}}{1 - w_{t,i} + w_{t,i} e^{-\alpha_t u_{t,i}}} .$$

We also get  $G(\boldsymbol{\theta}) = \sum_{i=1}^m \ln(1 + e^{\theta_i})$ , so

$$\begin{aligned} \min_{\mathbf{w} \cdot \mathbf{u}_t = 0} \Delta_F(\mathbf{w}, \mathbf{w}_t) &= \Delta_F(\mathbf{w}_{t+1}, \mathbf{w}_t) \\ &= \max_{\alpha} \left( \sum_{i=1}^m (\ln(1 + e^{\theta_{t,i}}) - \ln(1 + e^{\theta_{t,i} - \alpha u_{t,i}})) \right) \\ &= \max_{\alpha} \left( \sum_{i=1}^m -\ln(1 - w_{t,i} + w_{t,i} e^{-\alpha u_{t,i}}) \right). \end{aligned}$$

Finally, if we take  $F(\mathbf{w}) = \|\mathbf{w}\|^2/2 = \sum_{i=1}^m w_i^2/2$ , we get the squared Euclidean distance  $\Delta_F(\tilde{\mathbf{w}}, \mathbf{w}) = \|\tilde{\mathbf{w}} - \mathbf{w}\|^2/2$ . The gradient is the identity function, so  $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{u}_t$ , and  $G = F$ , so

$$\begin{aligned} \min_{\mathbf{w} \cdot \mathbf{u}_t = 0} \Delta_F(\mathbf{w}, \mathbf{w}_t) &= \Delta_F(\mathbf{w}_{t+1}, \mathbf{w}_t) \\ &= \frac{1}{2} \max_{\alpha} (\|\mathbf{w}_t\|^2 - \|\mathbf{w}_t - \alpha \mathbf{u}_t\|^2) \\ &= \frac{(\mathbf{w}_t \cdot \mathbf{u}_t)^2}{\|\mathbf{u}_t\|^2}. \end{aligned}$$

(In this special case we were thus able to solve the maximization in closed form.) Geometrically,  $\mathbf{w}_{t+1}$  is the point closest to  $\mathbf{w}_t$  on the hyperplane  $\mathbf{w} \cdot \mathbf{u}_t = 0$ .

## References

- [Ama85] S. Amari. *Differential Geometrical Methods in Statistics*. Springer, Berlin, 1985.
- [AW99] K. Azoury and M. K. Warmuth. Relative loss bounds for on-line density estimation with the exponential family of distributions. In *Proc. 15th Conf. on Uncertainty in Artificial Intelligence*. Morgan Kaufmann, San Francisco, CA, 1999.
- [Bre67] L. M. Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Physics*, 7:200–217, 1967.
- [BSS93] M. S. Bazaraa, H. D. Sherali, and C. M. Shetty. *Nonlinear Programming: Theory and Algorithms*. Wiley, New York, NY, 1993. Second edition.
- [Byl97] T. Bylander. The binary exponentiated gradient algorithm for learning linear functions. In *Proc. 10th Annu. Conf. on Comput. Learning Theory*, pages 184–192. ACM, New York, 1997.
- [CBKW94] N. Cesa-Bianchi, A. Krogh, and M. K. Warmuth. Bounds on approximate steepest descent for likelihood maximization in exponential families. *IEEE Transaction on Information Theory*, 40(4):1215–1220, 1994.
- [CL81] Y. Censor and A. Lent. An iterative row-action method for interval convex programming. *Journal of Optimization Theory and Applications*, 34(3):321–353, 1981.
- [Csi91] I. Csiszar. Why least squares and maximum entropy? An axiomatic approach for linear inverse problems. *The Annals of Statistics*, 19(4):2032–2066, 1991.
- [DDL97] S. Della Pietra, V. Della Pietra, and J. Lafferty. Inducing features of random fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(4):380–393, 1997.
- [DH99] N. Duffy and D. P. Helmbold. A geometric approach to leveraging weak learners. In *Computational Learning Theory: 4th European Conference (EuroCOLT '99)*, pages 18–33. Springer, Berlin, 1999.
- [FHT98] J. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: a statistical view of boosting. Technical report, Stanford University, 1998.
- [Fre95] Y. Freund. Boosting a weak learning algorithm by majority. *Information and Computation*, 121(2):256–285, September 1995.
- [FS97a] Y. Freund and R. E. Schapire. Adaptive game playing using multiplicative weights. To appear in *Games and Economic Behavior*.
- [FS97b] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.
- [FSSW97] Y. Freund, R. E. Schapire, Y. Singer, and M. K. Warmuth. Using and combining predictors that specialize. In *Proc. 29th Annual ACM Symposium on Theory of Computing*, pages 334–343. ACM, New York, 1997.
- [HUL91] J.-B. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms II*. Springer, Berlin, 1991.
- [HW98] M. Herbster and M. K. Warmuth. Tracking the best regressor. In *Proc. 11th Annu. Conf. on Comput. Learning Theory*, pages 24–31. ACM, New York, 1998.
- [JB90] L. Jones and C. Byrne. General entropy criteria for inverse problems, with applications to data compression, pattern classification and cluster analysis. *IEEE Transactions on Information Theory*, 36(1):23–30, 1990.
- [Jum90] G. Jumarie. *Relative information*. Springer, Berlin, 1990.
- [KK92] J. N. Kapur and H. K. Kesavan. *Entropy Optimization Principles with Applications*. Academic Press, 1992.
- [KW97] J. Kivinen and M. K. Warmuth. Additive versus exponentiated gradient updates for linear prediction. *Information and Computation*, 132(1):1–64, 1997.
- [KW99] J. Kivinen and M. K. Warmuth. Averaging expert predictions. In *Computational Learning Theory: 4th European Conference (EuroCOLT '99)*, pages 153–167. Springer, Berlin, 1999.
- [Laf99] J. Lafferty. Additive models, boosting, and inference for generalized divergences. In *Proc. 12th Annu. Conf. on Comput. Learning Theory*.

- ACM, New York, 1999.
- [LLW92] N. Littlestone, P. M. Long, and M. K. Warmuth. On-line learning of linear functions. Technical Report UCSC-CRL-91-29, University of California, Santa Cruz, 1992. Short version appeared in *Journal of Computational Complexity*, 5:1–23, 1995.
- [Lue84] D. G. Luenberger. *Linear and Nonlinear Programming*. Addison-Wesley, Reading, MA, 1984. Second edition.
- [Roc70] R. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- [ROM98] G. Rätsch, T. Onoda, and K. Müller. Soft margins for AdaBoost. Technical Report NC-TR-1998-021, NeuroCOLT2 Technical Report Series, 1998.
- [Sch90] R. E. Schapire. The strength of weak learnability. *Machine Learning*, 5(2):197–227, 1990.
- [SS98] R. E. Schapire and Y. Singer. Improved boosting algorithms using confidence-rated predictions. In *Proc. 11th Annu. Conf. on Comput. Learning Theory*, pages 80–91. ACM, New York, 1998.