

Experiments with a New Boosting Algorithm

DRAFT — PLEASE DO NOT DISTRIBUTE

Yoav Freund Robert E. Schapire

AT&T Research

600 Mountain Avenue

Rooms {2B-428, 2A-424}

Murray Hill, NJ 07974-0636

{yoav, schapire}@research.att.com

<http://www.research.att.com/orgs/ssr/people/{yoav,schapire}/>

January 22, 1996

Abstract

In an earlier paper [9], we introduced a new “boosting” algorithm called **AdaBoost** which, theoretically, can be used to significantly reduce the error of any learning algorithm that consistently generates classifiers whose performance is a little better than random guessing. We also introduced the related notion of a “pseudo-loss” which is a method for forcing a learning algorithm of multi-label concepts to concentrate on the labels that are hardest to discriminate. In this paper, we describe experiments we carried out to assess how well **AdaBoost** with and without pseudo-loss, performs on real learning problems.

We performed two sets of experiments. The first set compared boosting to Breiman’s [1] “bagging” method when used to aggregate various classifiers (including decision trees and single attribute-value tests). We compared the performance of the two methods on a collection of machine-learning benchmarks. In the second set of experiments, we studied in more detail the performance of boosting using a nearest-neighbor classifier on an OCR problem.

1 Introduction

“Boosting” is a general method for improving the performance of any learning algorithm. In theory, boosting can be used to significantly reduce the error of any “weak” learning algorithm that consistently generates classifiers which need only be a little bit better than random guessing. Despite the potential benefits of boosting promised by the theoretical results, the true practical value of boosting can only be assessed by testing the method on “real” learning problems. In this paper, we present such an experimental assessment of a new boosting algorithm called **AdaBoost**.

Boosting works by repeatedly running a given weak¹ learning algorithm on various distributions over the training data, and then combining the classifiers produced by the weak learner into a single composite classifier. The first provably effective boosting algorithms were presented by Schapire [19] and Freund [8]. More recently, we described and analyzed **AdaBoost**, and we argued that this new boosting algorithm has certain properties which make it more practical and easier to implement than its predecessors [9]. This algorithm, which we used in all our experiments, is described in detail in Section 2.

This paper describes two distinct sets of experiments. In the first set of experiments, described in Section 3, we compared boosting to “bagging,” a method described by Breiman [1] which works in the same general fashion (i.e., by repeatedly rerunning a given weak learning algorithm, and combining the computed classifiers), but which constructs each distribution in a simpler manner. (Details given below.) We compared boosting with bagging because both methods work by combining many classifiers. This comparison allows us to separate out the effect of modifying the distribution on each round (which is done differently by each algorithm) from the effect of voting multiple classifiers (which is done the same by each).

In our experiments, we compared boosting to bagging using a number of different weak learning algorithms of varying levels of sophistication. These include: (1) an algorithm that searches for very simple prediction rules which test on a single attribute (similar to Holte’s very simple classification rules [13]); (2) an algorithm that searches for a single good decision rule that tests on a conjunction of attribute tests (similar in flavor to the rule-formation part of Cohen’s RIPPER algorithm [2] and Fürnkranz and Widmer’s IREP algorithm [10]); and (3) Quinlan’s **C4.5** decision-tree algorithm [17]. We tested these algorithms on a collection of 27 benchmark learning problems taken from the UCI repository.

The main conclusion of our experiments is that boosting performs significantly and uniformly better than bagging when the weak learning algorithm generates fairly simple classifiers (algorithms (1) and (2) above). When combined with **C4.5**, boosting still seems to outperform bagging slightly, but the results are less compelling.

We also found that boosting can be used with very simple rules (algorithm (1)) to construct classifiers that are quite good relative, say, to **C4.5**. Kearns and Mansour [15] argue that **C4.5** can itself be viewed as a kind of boosting algorithm, so a comparison of **AdaBoost** and **C4.5** can be seen as a comparison of two competing boosting algorithms. See Dietterich, Kearns and Mansour’s paper [3] for more detail on this point.

In the second set of experiments, we test the performance of boosting on a nearest neighbor classifier for handwritten digit recognition. In this case the weak learning algorithm is very simple,

¹We use the term “weak” learning algorithm, even though, in practice, boosting might be combined with a quite strong learning algorithm such as **C4.5**.

and this lets us gain some insight to the interaction between the boosting algorithm and the nearest neighbor classifier. We show that the boosting algorithm is an effective way for finding a small subset of prototypes that performs almost as well as the complete set. We also show that it compares favorably to the standard method of Condensed Nearest Neighbor [12] in terms of its test error.

There seem to be two separate reasons for the improvement in performance that is achieved by boosting. The first and better understood effect of boosting is that it generates a hypothesis whose error on the training set is small by combining many hypotheses whose error may be large (but still better than random guessing). It seems that boosting may be helpful on learning problems having either of the following two properties. The first property, which holds for many real-world problems, is that the observed examples tend to have varying degrees of hardness. For such problems, the boosting algorithm tends to generate distributions that concentrate on the harder examples, thus challenging the weak learning algorithm to perform well on these harder parts of the sample space.

The second property is that the learning algorithm be sensitive to changes in the training examples so that significantly different hypotheses are generated for different training sets. In this sense, boosting is similar to Breiman’s bagging [1] which performs best when the weak learner exhibits such “unstable” behavior. However, unlike bagging, boosting tries actively to force the weak learning algorithm to change its hypotheses by constructing a “hard” distribution over the examples based on the performance of previously generated hypotheses.

The second effect of boosting has to do with variance reduction. Intuitively, taking a weighted majority over many hypotheses, all of which were trained on different samples taken out of the same training set, has the effect of reducing the random variability of the combined hypothesis. Thus, like bagging, boosting may have the effect of producing a combined hypothesis whose variance is significantly lower than those produced by the weak learner. However, unlike bagging, boosting may also reduce the bias of the learning algorithm, as discussed above. (See Kong and Dietterich [16] for further discussion of the bias and variance reducing effects of voting multiple hypotheses.) In our first set of experiments, we compare boosting and bagging, and try to use that comparison to separate between the bias and variance reducing effects of boosting.

Previous work. Drucker, Schapire and Simard [7, 6] performed the first experiments using a boosting algorithm. They used Schapire’s [19] original boosting algorithm combined with a neural net for an OCR problem. Follow-up comparisons to other ensemble methods were done by Drucker et al. [5]. More recently, Drucker and Cortes [4] used **AdaBoost** with a decision-tree algorithm for an OCR task. Jackson and Craven [14] used **AdaBoost** to learn classifiers represented by sparse perceptrons, and tested the algorithm on a set of benchmarks. Finally, Quinlan [18] recently conducted an independent comparison of boosting and bagging combined with **C4.5** on a collection of UCI benchmarks.

2 The boosting algorithm

In this section, we describe our boosting algorithm, called **AdaBoost**. See our earlier paper [9] for more details about the algorithm and its theoretical properties.

We describe two versions of the algorithm which we denote **AdaBoost.M1** and **AdaBoost.M2**. The two versions are equivalent for binary classification problems and differ only in their handling

Algorithm AdaBoost.M1

Input: sequence of m examples $\langle (x_1, y_1), \dots, (x_m, y_m) \rangle$ with labels $y_i \in Y = \{1, \dots, k\}$
weak learning algorithm **WeakLearn**
integer T specifying number of iterations

Initialize $D_1(i) = 1/m$ for all i .

Do for $t = 1, 2, \dots, T$

1. Call **WeakLearn**, providing it with the distribution D_t .
2. Get back a hypothesis $h_t : X \rightarrow Y$.
3. Calculate the error of h_t : $\epsilon_t = \sum_{i:h_t(x_i) \neq y_i} D_t(i)$. If $\epsilon_t > 1/2$, then set $T = t - 1$ and abort loop.
4. Set $\beta_t = \epsilon_t / (1 - \epsilon_t)$.
5. Update distribution D_t : $D_{t+1}(i) = \frac{D_t(i)}{Z_t} \times \begin{cases} \beta_t & \text{if } h_t(x_i) = y_i \\ 1 & \text{otherwise} \end{cases}$
where Z_t is a normalization constant (chosen so that D_{t+1} will be a distribution).

Output the final hypothesis: $h_{\hat{f}_n}(x) = \arg \max_{y \in Y} \sum_{t:h_t(x)=y} \log \frac{1}{\beta_t}$.

Figure 1: The algorithm **AdaBoost.M1**.

of problems with more than two classes.

2.1 AdaBoost.M1

We begin with the simpler version, **AdaBoost.M1**. The boosting algorithm takes as input a training set of m examples $S = \langle (x_1, y_1), \dots, (x_m, y_m) \rangle$ where x_i is an instance drawn from some space X and represented in some manner (typically, a vector of attribute values), and $y_i \in Y$ is the class label associated with x_i . In this paper, we always assume that the set of possible labels Y is of finite cardinality k .

In addition, the boosting algorithm has access to another unspecified learning algorithm, called the weak learning algorithm, which is denoted generically as **WeakLearn**. The boosting algorithm calls **WeakLearn** repeatedly in a series of rounds. On round t , the booster provides **WeakLearn** with a distribution D_t over the training set S . In response, **WeakLearn** computes a classifier or *hypothesis* $h_t : X \rightarrow Y$ which should misclassify a non trivial fraction of the training examples, relative to D_t . That is, the weak learner's goal is to find a hypothesis h_t which minimizes the (training) error $\epsilon_t = \Pr_{i \sim D_t}[h_t(x_i) \neq y_i]$. Note that this error is measured with respect to the distribution D_t that was provided to the weak learner. This process continues for T rounds, and, at last, the booster combines the weak hypotheses h_1, \dots, h_T into a single final hypothesis $h_{\hat{f}_n}$.

Still unspecified are (1) the manner in which D_t is computed on each round, and (2) how $h_{\hat{f}_n}$ is computed. Different boosting schemes answer these two questions in different ways. **AdaBoost.M1** uses the simple rule shown in Figure 1. The initial distribution D_1 is uniform over S so $D_1(i) = 1/m$ for all i . To compute distribution D_{t+1} from D_t and the last weak hypothesis h_t , we multiply the weight of example i by some number $\beta_t \in [0, 1)$ if h_t classifies x_i correctly, and otherwise the weight is left unchanged. The weights are then renormalized by dividing by the normalization constant Z_t . Effectively, "easy" examples that are correctly classified by many

Algorithm AdaBoost.M2

Input: sequence of m examples $\langle (x_1, y_1), \dots, (x_m, y_m) \rangle$ with labels $y_i \in Y = \{1, \dots, k\}$
 weak learning algorithm **WeakLearn**
 integer T specifying number of iterations

Let $B = \{(i, y) : i \in \{1, \dots, m\}, y \neq y_i\}$

Initialize $D_1(i, y) = 1/|B|$ for $(i, y) \in B$.

Do for $t = 1, 2, \dots, T$

1. Call **WeakLearn**, providing it with mislabel distribution D_t .
2. Get back a hypothesis $h_t : X \times Y \rightarrow [0, 1]$.
3. Calculate the pseudo-loss of h_t : $\epsilon_t = \frac{1}{2} \sum_{(i,y) \in B} D_t(i, y) (1 - h_t(x_i, y_i) + h_t(x_i, y))$.
4. Set $\beta_t = \epsilon_t / (1 - \epsilon_t)$.
5. Update D_t : $D_{t+1}(i, y) = \frac{D_t(i, y)}{Z_t} \cdot \beta_t^{(1/2)(1+h_t(x_i, y_i)-h_t(x_i, y))}$
 where Z_t is a normalization constant (chosen so that D_{t+1} will be a distribution).

Output the hypothesis: $h_{fin}(x) = \arg \max_{y \in Y} \sum_{t=1}^T \left(\log \frac{1}{\beta_t} \right) h_t(x, y)$.

Figure 2: The algorithm **AdaBoost.M2**.

of the previous weak hypotheses get lower weight, and “hard” examples which tend often to be misclassified get higher weight. Thus, **AdaBoost** focuses the most weight on the examples which seem to be hardest for **WeakLearn**.

The number β_t is computed as shown in the figure as a function of ϵ_t . The final hypothesis h_{fin} is a weighted vote (i.e. a weighted linear threshold) of the weak hypotheses. That is, for a given instance x , h_{fin} outputs the label y that maximizes the sum of the weights of the weak hypotheses predicting that label. The weight of hypothesis h_t is defined to be $\ln(1/\beta_t)$ so that greater weight is given to hypotheses with lower error.

The important theoretical property about **AdaBoost.M1** is stated in the following theorem. This theorem shows that if the weak hypotheses consistently have error only slightly better than $1/2$, then the error of the final hypothesis h_{fin} drops to zero exponentially fast. For binary classification problems, this means that the weak hypotheses need be only slightly better than random.

Theorem 1 ([9]) *Suppose the weak learning algorithm **WeakLearn**, when called by **AdaBoost.M1**, generates hypotheses with errors $\epsilon_1, \dots, \epsilon_T$, where ϵ_t is as defined in Figure 1. Assume each $\epsilon_t \leq 1/2$, and let $\gamma_t = 1/2 - \epsilon_t$. Then the following upper bound holds on the error of the final hypothesis h_{fin} :*

$$\frac{1}{m} |\{i : h_{fin}(x_i) \neq y_i\}| \leq \prod_{t=1}^T \sqrt{1 - 4\gamma_t^2} \leq \exp \left(-2 \sum_{t=1}^T \gamma_t^2 \right).$$

The main disadvantage of **AdaBoost.M1** is that it is unable to handle weak hypotheses with error greater than $1/2$. The expected error of a hypothesis which randomly guesses the label is $1 - 1/k$, where k is the number of possible labels. Thus **AdaBoost.M1** requirement for $k = 2$ is that the prediction is just slightly better than random guessing. However, when $k > 2$, the requirement of **AdaBoost.M1** is much stronger than that, and might be hard to meet.

2.2 AdaBoost.M2

The second version of **AdaBoost** attempts to overcome this difficulty by extending the communication between the boosting algorithm and the weak learner. First, we allow the weak learner to generate more expressive hypotheses whose output is a vector in $[0, 1]^k$, rather than a single label in Y . Intuitively, the y th component of this vector represents a “degree of belief” that the correct label is y . The components with values close to 1 or 0 correspond to those labels considered to be plausible or implausible, respectively.

While we give the weak learning algorithm more expressive power, we also place a more complex requirement on the performance of the weak hypotheses. Rather than using the usual prediction error, we ask that the weak hypotheses do well with respect to a more sophisticated error measure that we call the pseudo-loss. Unlike ordinary error which is computed with respect to a distribution over examples, pseudo-loss is computed with respect to a distribution over the set of all pairs of examples and incorrect labels. By manipulating this distribution, the boosting algorithm can focus the weak learner not only on hard-to-classify examples, but more specifically, on the incorrect labels that are hardest to discriminate. We will see that the boosting algorithm **AdaBoost.M2**, which is based on these ideas, achieves boosting if each weak hypothesis has pseudo-loss slightly better than random guessing.

More formally, a *mislabel* is a pair (i, y) where i is the index of a training example and y is an *incorrect* label associated with example i . Let B be the set of all mislabels:

$$B = \{(i, y) : i \in \{1, \dots, m\}, y \neq y_i\}.$$

A *mislabel distribution* is a distribution defined over the set B of all mislabels.

On each round t of boosting, **AdaBoost.M2** (Figure 2) supplies the weak learner with a mislabel distribution D_t . In response, the weak learner computes a hypothesis h_t of the form $h_t : X \times Y \rightarrow [0, 1]$.

Intuitively, we interpret each mislabel (i, y) as representing a binary question of the form: “Do you predict that the label associated with example x_i is y_i (the correct label) or y (one of the incorrect labels)?” With this interpretation, the weight $D_t(i, y)$ assigned to this mislabel represents the importance of distinguishing incorrect label y on example x_i .

A weak hypothesis h_t is then interpreted in the following manner. If $h_t(x_i, y_i) = 1$ and $h_t(x_i, y) = 0$, then h_t has (correctly) predicted that x_i ’s label is y_i , not y (since h_t deems y_i to be “plausible” and y “implausible”). Similarly, if $h_t(x_i, y_i) = 0$ and $h_t(x_i, y) = 1$, then h_t has (incorrectly) made the opposite prediction. If $h_t(x_i, y_i) = h_t(x_i, y)$, then h_t ’s prediction is taken to be a random guess. (Values for h_t in $(0, 1)$ are interpreted probabilistically.)

This interpretation leads us to define the *pseudo-loss* of hypothesis h_t with respect to mislabel distribution D_t by the formula

$$\epsilon_t = \frac{1}{2} \sum_{(i,y) \in B} D_t(i, y) (1 - h_t(x_i, y_i) + h_t(x_i, y)).$$

Space limitations prevent us from giving a complete derivation of this formula which is explained in detail in our earlier paper [9]. It should be clear, however, that the pseudo-loss is minimized when correct labels y_i are given values near 1 and incorrect labels $y \neq y_i$ values near 0. Further, note that pseudo-loss $1/2$ is trivially achieved by any constant-valued hypothesis h_t , and moreover

that a hypothesis h_t with pseudo-loss greater than $1/2$ can be replaced by the hypothesis $1 - h_t$ whose pseudo-loss is less than $1/2$.

The weak learner’s goal is to find a weak hypothesis h_t with small pseudo-loss. Thus, standard “off-the-shelf” learning algorithms may need some modification to be used in this manner, although this modification is often straightforward. After receiving h_t , the mislabel distribution is updated using a rule similar to the one used in **AdaBoost.M1**. The final hypothesis h_{fin} outputs, for a given instance x , the label y that maximizes a weighted average of the weak hypothesis values $h_t(x, y)$.

The following theorem gives a bound on the training error of the final hypothesis. Note that this theorem requires only that the weak hypotheses have pseudo-loss less than $1/2$, i.e., only slightly better than a trivial (constant-valued) hypothesis, regardless of the number of classes. Also, although the weak hypotheses h_t are evaluated with respect to the pseudo-loss, we of course evaluate the final hypothesis h_{fin} using the ordinary error measure.

Theorem 2 ([9]) *Suppose the weak learning algorithm **WeakLearn**, when called by **AdaBoost.M2** generates hypotheses with pseudo-losses $\epsilon_1, \dots, \epsilon_T$, where ϵ_t is as defined in Figure 2. Let $\gamma_t = 1/2 - \epsilon_t$. Then the following upper bound holds on the error of the final hypothesis h_{fin} :*

$$\frac{1}{m} |\{i : h_{fin}(x_i) \neq y_i\}| \leq (k - 1) \prod_{t=1}^T \sqrt{1 - 4\gamma_t^2} \leq (k - 1) \exp\left(-2 \sum_{t=1}^T \gamma_t^2\right)$$

where k is the number of classes.

3 Boosting and bagging

In this section, we describe our experiments comparing boosting and bagging on the UCI benchmarks.

We first mention briefly a small implementation issue: Many learning algorithms can be modified to handle examples that are weighted by a distribution such as the one created by the boosting algorithm. When this is possible, the booster’s distribution D_t is supplied directly to the weak learning algorithm, a method we call boosting by *reweighting*. However, some learning algorithms require an unweighted set of examples. For such a weak learning algorithm, we instead choose a set of examples from S independently at random according to the distribution D_t with replacement. The number of examples to be chosen on each round is a matter of discretion; in our experiments, we chose m examples on each round, where m is the size of the original training set S . We refer to this method as boosting by *resampling*.

Boosting by resampling is also possible when using the pseudo-loss. In this case, a set of mislabels are chosen from the set B of all mislabels with replacement according to the given distribution D_t . Such a procedure is consistent with the interpretation of mislabels discussed in Section 2.2. In our experiments, we chose a sample of size $|B| = m(k - 1)$ on each round when using the resampling method.

3.1 The weak learning algorithms

As mentioned in the introduction, we used three weak learning algorithms in these experiments. In all cases, the examples are described by a vector of values which corresponds to a fixed set of

features or attributes. These values may be discrete or continuous. Some of the examples may have missing values. All three of the weak learners build hypotheses which classify examples by repeatedly testing the values of chosen attributes.

The first and simplest weak learner, which we call **FindAttrTest**, searches for the single attribute-value test with minimum error (or pseudo-loss) on the training set. More precisely, **FindAttrTest** computes a classifier which is defined by an attribute a , a value v and three predictions p_0 , p_1 and $p_?$. This classifier classifies a new example x as follows: if the value of attribute a is missing on x , then predict $p_?$; if attribute a is discrete and its value on example x is equal to v , or if attribute a is continuous and its value on x is at most v , then predict p_0 ; otherwise predict p_1 . If using ordinary error (**AdaBoost.M1**), these “predictions” $p_0, p_1, p_?$ would be simple classifications; for pseudo-loss, the “predictions” would be vectors in $[0, 1]^k$ (where k is the number of classes).

The algorithm **FindAttrTest** searches exhaustively for the classifier of the form given above with minimum error or pseudo-loss with respect to the distribution provided by the booster. In other words, all possible values of a , v , p_0 , p_1 and $p_?$ are considered. This search can be carried out in time linear in the size of the training set (per round of boosting).

For this algorithm, we used boosting with reweighting.

The second weak learner does a somewhat more sophisticated search for a decision rule that tests on a conjunction of attribute-value tests. We sketch the main ideas of this algorithm, which we call **FindDecRule**, but omit some of the finer details for lack of space. These details will be provided in the full paper.

First, the algorithm requires an unweighted training set, so we use the resampling version of boosting. The given training set is randomly divided into a growing set using 70% of the data, and a pruning set with the remaining 30%. In the first phase, the growing set is used to grow a list of attribute-value tests. Each test compares an attribute a to a value v , similar to the tests used by **FindAttrTest**. We use an entropy-based potential function to guide the growth of the list of tests. The list is initially empty, and one test is added at a time, each time choosing the test that will cause the greatest drop in potential. After the test is chosen, only one branch is expanded, namely, the branch with the highest remaining potential. The list continues to be grown in this fashion until no test remains which will further reduce the potential.

In the second phase, the list is pruned by selecting the prefix of the list with minimum error (or pseudo-loss) on the pruning set.

The third weak learner is Quinlan’s **C4.5** decision-tree algorithm [17]. We used all the default options with pruning turned on. Since **C4.5** expects an unweighted training sample, we used resampling. Also, we did not attempt to use **AdaBoost.M2** since **C4.5** is designed to minimize error, not pseudo-loss.

3.2 Bagging

We compared boosting to Breiman’s [1] “bootstrap aggregating” or “bagging” method for training and combining multiple copies of a learning algorithm. Briefly, the method works by training each copy of the algorithm on a bootstrap sample, i.e., a sample of size m chosen uniformly at random with replacement from the original training set S (of size m). The multiple hypotheses that are computed are then combined using simple voting; that is, the final composite hypothesis classifies an example x to the class most often assigned by the underlying “weak” hypotheses. See

his paper for more details. The method can be quite effective, especially, according to Breiman, for “unstable” learning algorithms for which a small change in the data effects a large change in the computed hypothesis.

In order to compare **AdaBoost.M2**, which uses pseudo-loss, to bagging, we also extended bagging in a natural way for use with a weak learning algorithm that minimizes pseudo-loss rather than ordinary error. As described in Section 2.2, such a weak learning algorithm expects to be provided with a distribution over the set B of all mislabels. On each round of bagging, we construct this distribution using the bootstrap method; that is, we select $|B|$ mislabels from B (chosen uniformly at random with replacement), and assign each mislabel weight $1/|B|$ times the number of times it was chosen. The hypotheses h_t computed in this manner are then combined using voting in a natural manner; namely, given x , the combined hypothesis outputs the label y which maximizes $\sum_t h_t(x, y)$.

For either error or pseudo-loss, the differences between bagging and boosting can be summarized as follows: (1) bagging always uses resampling rather than reweighting; (2) bagging does not modify the distribution over examples or mislabels, but instead always uses the uniform distribution; and (3) in forming the final hypothesis, bagging gives equal weight to each of the weak hypotheses.

3.3 The experiments

We conducted our experiments on a collection of machine learning datasets available from the repository at University of California at Irvine.² A summary of some of the properties of these datasets is given in Table 3 in the appendix. Some datasets are provided with a test set. For these, we reran each algorithm 20 times (since some of the algorithms are randomized), and averaged the results. For datasets with no provided test set, we used 10-fold cross validation, and averaged the results over 10 runs (for a total of 100 runs of each algorithm on each dataset).

In all our experiments, we set the number of rounds of boosting or bagging to be $T = 100$.

3.4 Results and discussion

The results of our experiments are shown in Table 1. The figures indicate test error rate averaged over multiple runs of each algorithm. Columns indicate which weak learning algorithm was used, and whether pseudo-loss (**AdaBoost.M2**) or error (**AdaBoost.M1**) was used. Columns labeled “–” indicate that the weak learning algorithm was used by itself (with no boosting or bagging). Columns using boosting or bagging are marked “boost” and “bag,” respectively.

One of our goals in carrying out these experiments was to determine if boosting using pseudo-loss (rather than error) is worthwhile. These experiments indicate that pseudo-loss is definitely worth the effort. Using pseudo-loss did dramatically better than error on every non-binary problem (except it did slightly worse on “iris” with three classes). Because **AdaBoost.M2** did so much better than **AdaBoost.M1**, we will only discuss **AdaBoost.M2** in the remaining discussion.

Using pseudo-loss with bagging gave mixed results in comparison to ordinary error. Overall, pseudo-loss gave better results, but occasionally, using pseudo-loss hurt considerably.

For the simpler weak learning algorithms (**FindAttrTest** and **FindDecRule**), boosting did significantly and uniformly better than bagging. The boosting error rate was worse than the bagging

²URL “<http://www.ics.uci.edu/~mllearn/MLRepository.html>”

name	FindAttrTest					FindDecRule					C4.5		
	error		pseudo-loss			error		pseudo-loss			error		
	-	boost	bag	boost	bag	-	boost	bag	boost	bag	-	boost	bag
soybean-small	57.6	56.4	48.7	0.2	20.5	51.8	56.0	45.7	0.4	2.9	2.2	3.4	2.2
labor	25.1	8.8	19.1	9.0	18.9	24.0	7.3	14.6	7.2	15.7	15.8	13.1	11.3
promoters	29.7	8.9	16.6	9.1	17.2	25.9	8.3	13.7	8.5	14.2	22.0	5.0	12.7
iris	35.2	4.7	28.4	4.8	7.1	38.3	4.3	18.8	4.8	5.5	5.9	5.0	5.0
hepatitis	19.7	18.6	16.8	18.3	17.4	21.6	18.0	20.1	18.4	20.6	21.2	16.3	17.5
sonar	25.9	16.5	25.9	16.8	25.9	31.4	16.2	26.1	15.4	25.7	28.9	19.0	24.3
glass	51.5	51.1	50.9	29.4	54.2	49.7	48.5	47.2	25.0	52.0	31.7	22.7	25.7
audiology.stand	57.7	57.7	57.7	26.9	77.2	57.7	57.7	57.7	18.5	63.5	15.4	15.4	10.2
cleve	27.8	18.8	22.4	18.8	21.9	27.4	19.7	20.3	20.6	19.9	26.6	21.7	20.9
soybean-large	64.8	64.5	59.0	9.8	74.2	73.6	73.6	73.6	7.2	66.0	13.3	6.8	12.2
ionosphere	17.8	8.5	17.3	8.5	17.2	10.3	6.6	9.3	6.4	9.4	8.9	5.8	6.2
house-votes-84	4.4	3.7	4.4	3.7	4.4	5.0	4.4	4.4	4.3	4.5	3.5	5.1	3.6
votes1	12.7	8.9	12.7	8.9	12.7	13.2	9.4	11.2	9.4	10.7	10.3	10.4	9.2
crx	14.5	14.4	14.5	14.4	14.5	14.5	13.5	14.5	13.6	14.5	15.8	13.8	13.6
breast-cancer-w	8.4	4.4	6.7	4.4	6.6	8.1	4.1	5.3	4.0	5.2	5.0	3.3	3.2
pima-indians-di	26.1	24.4	26.1	24.5	26.0	27.8	25.3	26.4	25.4	26.6	28.4	25.7	24.4
vehicle	64.3	64.4	57.6	26.1	56.1	61.3	61.2	61.0	25.0	54.3	29.9	22.6	26.1
vowel	81.8	81.8	76.8	18.2	74.7	82.0	72.7	71.6	6.5	63.2	2.2	0.0	0.0
german	30.0	24.9	30.4	24.9	30.3	30.0	25.4	29.6	25.6	29.7	29.4	25.0	24.6
segmentation	75.8	75.8	54.5	4.2	72.5	73.7	53.3	54.3	2.4	58.0	3.6	1.4	2.7
hypothyroid	2.2	1.0	2.2	1.0	2.2	0.8	1.0	0.7	1.0	0.7	0.8	1.0	0.8
sick-euthyroid	5.6	3.0	5.6	3.0	5.6	2.4	2.4	2.2	2.4	2.1	2.2	2.1	2.1
splice	37.0	9.2	35.6	4.4	33.4	29.5	8.0	29.5	4.0	29.5	5.8	4.9	5.2
kr-vs-kp	32.8	4.4	30.7	4.4	31.3	24.6	0.7	20.8	0.6	21.7	0.5	0.3	0.6
satimage	58.3	58.3	58.3	14.9	41.6	57.6	56.5	56.7	13.1	30.0	14.8	8.9	10.6
agaricus-lepiot	11.3	0.0	11.3	0.0	11.3	8.2	0.0	8.2	0.0	8.2	0.0	0.0	0.0
letter-recognit	92.9	92.9	91.9	34.1	93.7	92.3	91.8	91.8	30.4	93.7	13.8	3.3	6.8

Table 1: Test error rates of various algorithms on benchmark problems.

error rate (using either pseudo-loss or error) on a very small number of benchmark problems, and on these, the difference in performance was quite small. On average, for **FindAttrTest**, boosting improved the error rate over using **FindAttrTest** alone by 55.1%, compared to bagging which gave an improvement of only 10.6% using pseudo-loss or 8.4% using error. For **FindDecRule**, boosting improved the error rate by 53.1%, bagging by only 19.0% using pseudo-loss, 13.1% using error.

When using **C4.5** as the weak learning algorithm, boosting and bagging seem more evenly matched, although boosting still seems to have a slight advantage. On average, boosting improved the error rate by 23.6%, bagging by 20.7%. Boosting beat bagging by more than 2% on 6 of the benchmarks; while bagging beat boosting by this amount on only 1 benchmark. For the remaining 20 benchmarks, the difference in performance was less than 2%.

Using boosting with **FindAttrTest** does quite well as a learning algorithm in its own right, in comparison, say, to **C4.5**. This algorithm beat **C4.5** on 10 of the benchmarks (by at least 2%), tied on 13, and lost on 4. As mentioned above, its average performance relative to using **FindAttrTest** by itself was 55.1%. In comparison, **C4.5**'s improvement in performance over **FindAttrTest** was 49.9%.

Using boosting with **FindDecRule** did somewhat better. The win-tie-lose numbers for this

algorithm (compared to **C4.5**) were 12-12-3, and its average improvement over **FindAttrTest** was 58.2%.

4 Boosting a nearest-neighbor classifier

In this section we study the performance of a learning algorithm which combines **AdaBoost** and a variant of the nearest-neighbor classifier. We test the combined algorithm on the problem of recognizing handwritten digits. Our goal is not to improve on the accuracy of the nearest neighbor classifier, but rather, to speed it up. Speedup is achieved by reducing the number of prototypes in the hypothesis and the number of required distance calculations without increasing the error rate. It is a similar approach to that of nearest-neighbor editing [11, 12] in which one tries to find the minimal set of prototypes that is sufficient to label all the training set correctly.

The dataset comes from the US Postal Service (USPS) and consists of 9709 training examples and 2007 test examples. The training and test examples are evidently drawn from rather different distributions as there is a very significant improvement in the performance if the partition of the data into training and testing is done at random (rather than using the given partition). We report results both on the original partitioning and on a training set and a test set of the same sizes that were generated by randomly partitioning the union of the original training and test sets.

Each image is represented by a 16×16 -matrix of 8-bit pixels. The metric that we use is the standard Euclidean distance between the images (viewed as vectors in \mathbb{R}^{256}). This is a very naive metric, but it gives reasonably good performance. A nearest-neighbor classifier which uses all the training examples as prototypes achieves a test error of 5.7% (2.3% on randomly partitioned data). Using tangent distance [20] is in our future plans.

Our weak learning algorithm is simply to use a random set of examples as prototypes, chosen according to the distribution provided by the boosting algorithm. A standard nearest-neighbor classifier would predict the label of a test example according to the identity of its closest prototype. However, we found that significant improvement was achieved by using a variant of this algorithm that is optimized for use with **AdaBoost.M2**. To do this we consider all of the training set examples that are closest to each selected prototype, and choose the fixed prediction that minimizes the pseudo-loss for this (weighted) set of examples. This prediction choice is described in detail in Example 5 in [9]. In addition, we make the following modifications to the basic scheme:

- When using **AdaBoost.M2**, it is possible that hypothesis h_t has pseudo-loss less than $1/2$ on distribution D_{t+1} . When this happens, we reuse in hypothesis h_{t+1} the same set of prototypes that were used for h_t . Without increasing the size of the final hypothesis, this method reduces the theoretical error bound, and, as we observe in experiments, also the actual error on the training set.
- Instead of just selecting the prototypes at random, the algorithm repeatedly selects ten random prototypes and adds the one which causes the largest reduction in the pseudo-loss.

We ran 30 iterations of the boosting algorithm, and the number of prototypes we used were 10 for the first weak hypothesis, 20 for the second, 40 for the third, 80 for the next five, and 100 for the remaining weak hypotheses. These sizes were chosen so that the errors of all of the weak hypotheses are approximately equal.

round	# Proto- types	random partition					given partition				
		AdaBoost			Strawman		AdaBoost			Strawman	
		theory	train	test	train	test	theory	train	test	train	test
1	10	524.6	42.8	43.5	44.8	43.5	536.3	49.1	47.6	39.2	41.9
5	230	86.4	5.8	8.7	5.5	7.8	83.0	5.4	13.7	4.3	9.6
10	670	16.0	0.2	5.5	2.5	5.3	10.9	0.1	9.1	1.4	7.7
13	970	4.5	0.0	4.6	1.7	4.7	3.3	0.0	7.8	0.9	7.3
15	1170	2.4	0.0	3.9	1.4	4.7	1.5	0.0	7.8	0.7	7.4
20	1670	0.4	0.0	3.4	1.4	4.9	0.2	0.0	7.2	0.5	6.7
25	2170	0.1	0.0	3.3	1.2	4.2	0.0	0.0	6.8	0.3	7.1
30	2670	0.0	0.0	3.1	1.1	3.9	0.0	0.0	6.6	0.2	6.5

Table 2: Errors on randomly selected training and test sets, in percent. For columns labeled “random partition,” a random partition of the union of the training and test sets was used; “given partition” means the provided partition into training and test sets was used. Columns labeled “theory” give theoretical upper bound on training error (see Theorem 2).

We compared the performance of our algorithm to a strawman algorithm which uses a single set of prototypes. Similar to our algorithm, the prototype set is generated incrementally, comparing ten prototype candidates at each step. We compared the performance of the boosting algorithm to that of the strawman hypothesis that uses the same number of prototypes. We also compared our performance to that of the condensed nearest neighbor rule [12] (CNN), a greedy method for finding a small set of prototypes which correctly classify the entire training set.

4.1 Results and discussion

The results of our experiments are summarized in Table 2 and Figure 4.

In the left part of Table 2 we have the results of an experiment in which the test set was selected randomly from the union of the test and training data of the USPS dataset. We see that the boosting algorithm outperforms the strawman algorithm; the difference is dramatic on the training set but is also significant on the test set. Comparing these results to those of CNN, we find that the boosting algorithm reached zero error on the training set after 970 prototypes, which is about the same as CNN which reduced the set of prototypes to 964. However, the *test* error of this CNN was 5.7%, while the test error achieved by the 970 prototypes found by boosting was 4.5% and was further reduced to 3.1% when 2670 prototypes were used. Better editing algorithms might reduce the set of prototypes further, but it seems unlikely that this will reduce the test error. The error achieved by using the full training set (9709 prototypes) was 2.3%.

These results are also described by the graphs in Figure 4. The uppermost jagged line is a concatenation of the errors of the weak hypotheses with respect to the corresponding weights on the training set. Each peak followed by a valley corresponds to the beginning and end errors of a weak hypothesis. As we see the weighted error always started around 50% on the beginning of a boosting iteration and reached 20% – 30% by its end. The heaviest line describes the upper bound on the training error that is guaranteed by Theorem 2, and the two bottom lines describe the training and test error of the final combined hypothesis.

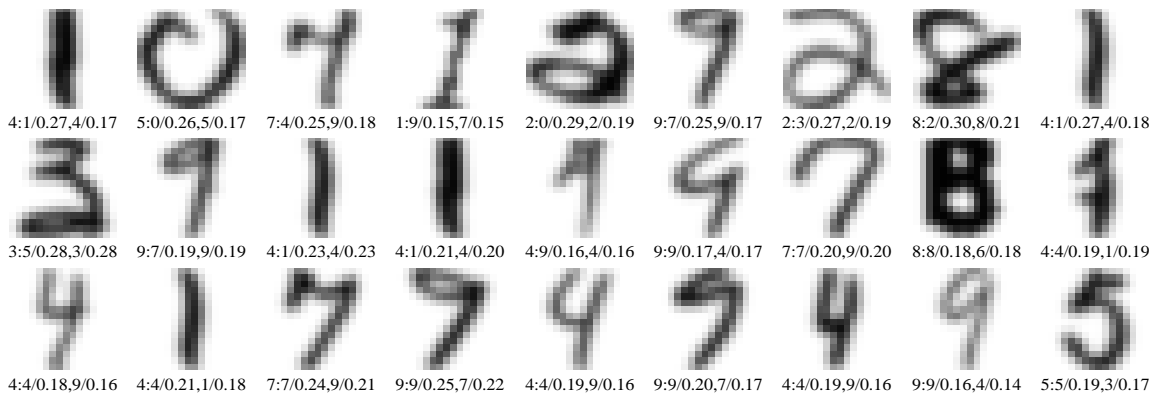


Figure 3: A sample of the examples that have the largest weight after 3 out of the 30 boosting iterations. The first line is after iteration 4, the second after iteration 12 and the third after iteration 25. Underneath each image we have a line of the form $d:\ell_1/w_1,\ell_2/w_2$, where d is the label of the example, ℓ_1 and ℓ_2 are the labels that get the highest and second highest vote from the combined hypothesis at that point in the run of the algorithm, and w_1, w_2 are the corresponding votes.

It is interesting that the performance of the boosting algorithm on the test set improved significantly after the error on the training set has already become zero. This is surprising because an “Occam’s razor” argument would predict that increasing the complexity of the hypothesis after the error has been reduced to zero is likely to degrade the performance on the test set.

The right hand side of Table 2 summarizes the results of running our algorithm on the training and test set as they were defined. Here the performance of boosting was similar to that of the strawman algorithm. However, it was still significantly better than that of CNN, which achieved a test error of 7.8% using 835 prototypes. It seems that the difference between the distribution of the test set and the training set removed the advantage that boosting had over the strawman algorithm.

We observed that when calculating the prediction of the combined hypothesis as we add the weighted vote of one weak hypothesis at a time, we can sometimes conclude what the final vote will be before calculating all of the hypotheses. This is possible when the difference between the current largest vote and the current second largest vote is larger than the total weight of the remaining hypotheses. In our experiments we found that the average number of weak hypotheses that had to be considered is 24.0 for the randomly chosen training set and 23.6 for the original training set. We can thus, on average, reduce the number of distance calculations that are required for evaluating the hypothesis from 2670 to 2070 without changing the predictions.

It is instructive to observe the examples that are given large weights by the boosting algorithm. A sample of these is given in Figure 3. There seem to be two types of “hard” examples. First are examples which are very atypical or wrongly labeled (such as example 2 on the first line and examples 3, 4 and 9 on the second line). The second type, which tends to dominate on later iterations, consists of examples that are very similar to each other but have different labels (such as examples 3 versus 4 and 1 versus 8 on the third line). Although the algorithm at this point was correct on all training examples, it is clear from the votes it assigned to different labels for these example pairs that it was still trying to improve the discrimination between these hard to discriminate pairs. This agrees with our intuition that the pseudo-loss is a mechanism that causes

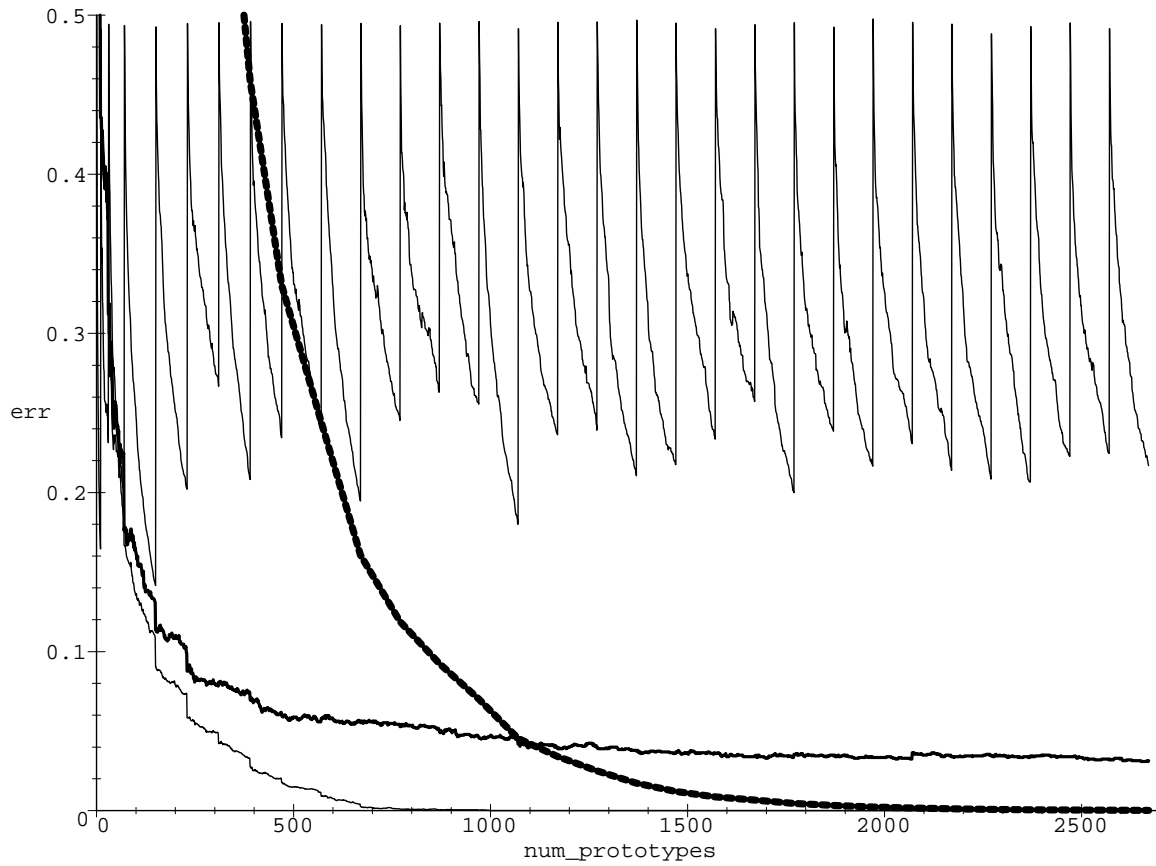


Figure 4: Graphs of the performance of the boosting algorithm on the randomly partitioned USPS dataset. The horizontal axis indicates the total number of prototypes that were added to the combined hypothesis, and the vertical axis indicates error. The topmost jagged line indicates the error of the weak hypothesis that is trained at this point on the weighted training set. The bold curve is the bound on the training error that is calculated based on the performance of the weak learner. The lowest thin curve is the performance of the combined hypothesis on the training set. The medium-bold curve is the performance of the combined hypothesis on the test set.

the boosting algorithm to concentrate on the hard to discriminate labels of hard examples.

Acknowledgements

Thanks to Jason Catlett and William Cohen for extensive advice on the design of our experiments. Thanks to Ross Quinlan for first suggesting a comparison of boosting and bagging. Thanks also to Leo Breiman, Corinna Cortes, Harris Drucker, Jeff Jackson, Michael Kearns, Ofer Matan, Partha Niyogi, Warren Smith, and David Wolpert for helpful comments, suggestions and criticisms.

References

- [1] Leo Breiman. Bagging predictors. Technical Report 421, Department of Statistics, University of California at Berkeley, 1994.

- [2] William Cohen. Fast effective rule induction. In *Proceedings of the Twelfth International Conference on Machine Learning*, pages 115–123, 1995.
- [3] Tom Dietterich, Michael Kearns, and Yishay Mansour. Applying the weak learning framework to understand and improve C4.5. Unpublished manuscript, 1996.
- [4] Harris Drucker and Corinna Cortes. Boosting decision trees. In *Advances in Neural Information Processing Systems 8*, 1996.
- [5] Harris Drucker, Corinna Cortes, L. D. Jackel, Yann LeCun, and Vladimir Vapnik. Boosting and other ensemble methods. *Neural Computation*, 6(6):1289–1301, 1994.
- [6] Harris Drucker, Robert Schapire, and Patrice Simard. Boosting performance in neural networks. *International Journal of Pattern Recognition and Artificial Intelligence*, 7(4):705–719, 1993.
- [7] Harris Drucker, Robert Schapire, and Patrice Simard. Improving performance in neural networks using a boosting algorithm. In *Advances in Neural Information Processing Systems 5*, pages 42–49, 1993.
- [8] Yoav Freund. Boosting a weak learning algorithm by majority. *Information and Computation*, 121(2):256–285, 1995.
- [9] Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. Unpublished manuscript available electronically (on our web pages, or by email request). An extended abstract appeared in *Computational Learning Theory: Second European Conference, EuroCOLT '95*, pages 23–37, Springer-Verlag, 1995.
- [10] Johannes Fürnkranz and Gerhard Widmer. Incremental reduced error pruning. In *Machine Learning: Proceedings of the Eleventh International Conference*, pages 70–77, 1994.
- [11] Geoffrey W. Gates. The reduced nearest neighbor rule. *IEEE Transactions on Information Theory*, pages 431–433, 1972.
- [12] Peter E. Hart. The condensed nearest neighbor rule. *IEEE Transactions on Information Theory*, IT-14:515–516, May 1968.
- [13] Robert C. Holte. Very simple classification rules perform well on most commonly used datasets. *Machine Learning*, 11(1):63–91, 1993.
- [14] Jeffrey C. Jackson and Mark W. Craven. Learning sparse perceptrons. In *Advances in Neural Information Processing Systems 8*, 1996.
- [15] Michael Kearns and Yishay Mansour. On the boosting ability of top-down decision tree learning algorithms. In *Proceedings of the Twenty-Eighth Annual ACM Symposium on the Theory of Computing*, 1996.
- [16] Eun Bae Kong and Thomas G. Dietterich. Error-correcting output coding corrects bias and variance. In *Proceedings of the Twelfth International Conference on Machine Learning*, pages 313–321, 1995.

- [17] J. Ross Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, 1993.
- [18] J. Ross Quinlan. Bagging, boosting, and C4.5. unpublished manuscript, 1996.
- [19] Robert E. Schapire. The strength of weak learnability. *Machine Learning*, 5(2):197–227, 1990.
- [20] Patrice Simard, Yann Le Cun, and John Denker. Efficient pattern recognition using a new transformation distance. In *Advances in Neural Information Processing Systems*, volume 5, pages 50–58, 1993.

A The Irvine datasets

name	# examples		# classes	# attributes		missing values
	train	test		discrete	cont.	
soybean-small	47	-	4	35	-	-
labor	57	-	2	8	8	✓
promoters	106	-	2	57	-	-
iris	150	-	3	-	4	-
hepatitis	155	-	2	13	6	✓
sonar	208	-	2	-	60	-
glass	214	-	7	-	9	-
audiology.stand	226	-	24	69	-	✓
cleve	303	-	2	7	6	✓
soybean-large	307	376	19	35	-	✓
ionosphere	351	-	2	-	34	-
house-votes-84	435	-	2	16	-	✓
votes1	435	-	2	15	-	✓
crx	690	-	2	9	6	✓
breast-cancer-w	699	-	2	-	9	✓
pima-indians-di	768	-	2	-	8	-
vehicle	846	-	4	-	18	-
vowel	528	462	11	-	10	-
german	1000	-	2	13	7	-
segmentation	2310	-	7	-	19	-
hypothyroid	3163	-	2	18	7	✓
sick-euthyroid	3163	-	2	18	7	✓
splice	3190	-	3	60	-	-
kr-vs-kp	3196	-	2	36	-	-
satimage	4435	2000	6	-	36	-
agaricus-lepiot	8124	-	2	22	-	-
letter-recognit	16000	4000	26	-	16	-

Table 3: The benchmark machine learning problems used in the experiments.