

## ARCING CLASSIFIERS

Leo Breiman\*  
 Statistics Department  
 University of California at Berkeley

### ABSTRACT

Recent work has shown that combining multiple versions of unstable classifiers such as trees or neural nets results in reduced test set error. One of the more effective is bagging (Breiman [1996a] ) Here, modified training sets are formed by resampling from the original training set, classifiers constructed using these training sets and then combined by voting. Freund and Schapire [1995,1996] propose an algorithm the basis of which is to adaptively resample and combine (hence the acronym--arc) so that the weights in the resampling are increased for those cases most often misclassified and the combining is done by weighted voting. Arcing is more successful than bagging in test set error reduction. We explore two arcing algorithms, compare them to each other and to bagging, and try to understand how arcing works. We introduce the definitions of bias and variance for a classifier as components of the test set error. Unstable classifiers can have low bias on a large range of data sets. Their problem is high variance. Combining multiple versions either through bagging or arcing reduces variance significantly

\* Partially supported by NSF Grant 1-444063-21445

### 1. Introduction

Some classification and regression methods are unstable in the sense that small perturbations in their training sets or in construction may result in large changes in the constructed predictor. Subset selection methods in regression, decision trees in regression and classification, and neural nets are unstable (Breiman [1996b]).

Unstable methods can have their accuracy improved by perturbing and combining. That is--generate multiple versions of the predictor by perturbing the training set or construction method. Then combine these multiple versions into a single predictor. For instance Ali [1995] generates multiple classification trees by choosing randomly from among the best splits at a node and combines trees using maximum likelihood. Breiman [1996b] adds noise to the response variable in regression to generate multiple subset regressions and then averages these. We use the label P&C (perturb and combine) to designate this group of methods.

One of the more effective P&C methods is bagging (Breiman [1996a]). Bagging perturbs the training set repeatedly to generate multiple predictors and combines these by simple voting (classification) or averaging (regression). Let the training set  $T$  consist of  $N$  cases (instances) labeled by  $n = 1, 2, \dots, N$ . Put equal probabilities  $p(n) = 1/N$  on each case, and using these

probabilities, sample with replacement (bootstrap)  $N$  times from the training set  $T$  forming the resampled training set  $T^{(B)}$ . Some cases in  $T$  may not appear in  $T^{(B)}$ , some may appear more than once. Now use  $T^{(B)}$  to construct the predictor, repeat the procedure and combine. Bagging applied to CART gave dramatic decreases in test set errors.

Freund and Schapire recently [1995], [1996] proposed a P&C algorithm which was designed to drive the training set error rapidly to zero. But if their algorithm is run far past the point at which the training set error is zero, it gives better performance than bagging on a number of real data sets. The crux of their idea is this: start with  $p(n) = 1/N$  and resample from  $T$  to form the first training set  $T^{(1)}$ . As the sequence of classifiers and training sets is being built, increase  $p(n)$  for those cases that have been most frequently misclassified. At termination, combine classifiers by weighted or simple voting. We will refer to algorithms of this type as Adaptive Resampling and Combining, or arcing algorithms. In honor of Freund and Schapire's discovery, we denote their specific algorithm by arc-fs, and discuss their theoretical efforts to relate training set to test set error in the Appendix.

To better understand stability and instability, and what bagging and arcing do, in Section 2 we define the concepts of bias and variance for classifiers. The difference between the test set misclassification error for the classifier and the minimum error achievable is the sum of the bias and variance. Unstable classifiers such as trees characteristically have high variance and low bias. Stable classifiers like linear discriminant analysis have low variance, but can have high bias. This is illustrated on several examples of artificial data. Section 3 looks at the effects of arcing and bagging trees on bias and variance.

The main effect of both bagging and arcing is to reduce variance. Arcing seems to usually do better at this than bagging. Arc-fs does complex things and its behavior is puzzling. But the variance reduction comes from the adaptive resampling and not the specific form of arc-fs. To show this, we define a simpler arc algorithm denoted by arc-x4 whose accuracy is comparable to arc-fs. The two appear to be at opposite poles of the arc spectrum. Arc-x4 was concocted to demonstrate that arcing works not because of the specific form of the arc-fs algorithm, but because of the adaptive resampling.

Freund and Schapire [1996] compare arc-fs to bagging on 27 data sets and conclude that arc-fs has a small edge in test set error rates. We tested arc-fs, arc-x4 and bagging on the 10 real data sets used in our bagging paper and get results more favorable to arcing. These are given in Section 4. arc-fs and arc-x4 finish in a dead heat. On a few data sets one or the other is a little better, but both are almost always significantly better than bagging. We also look at arcing and bagging applied to the US Postal Service digit data base.

The overall results of arcing are exciting--it turns a good but not great classifier (CART) into a procedure that seems to usually get close to the lowest achievable test set error rates. Furthermore, the arc-classifier is off-the-shelf. Its performance does not depend on any tuning or settings for particular problems. Just read in the data and press the start button. It is also, by neural net standards, blazingly fast to construct.

Section 5 gives the results of some experiments aimed at understanding how arc-fs and arc-x4 work. Each algorithm has distinctive and different signatures. Generally, arc-fs uses a smaller number of distinct cases in the resampled training sets and the successive values of  $p(n)$  are highly variable. The successive training sets in arc-fs rock back and forth and there is no convergence to a final set of  $\{p(n)\}$ . The back and forth rocking is more subdued in arc-x4, but there is still no convergence to a final  $\{p(n)\}$ . This variability may be an essential ingredient of successful arcing algorithms.

Instability is an essential ingredient for bagging or arcing to improve accuracy. Nearest neighbors are stable and Breiman[1996a] noted that bagging does not improve nearest neighbor classification. Linear discriminant analysis is also relatively stable (low variance) and in Section 6 our experiments show that neither bagging nor arcing has any effect on linear discriminant error rates.

Sections 7 and 8 contain remarks--mainly aimed at understanding how bagging and arcing work. The reason that bagging reduces error is fairly transparent. But it is not at all clear yet, in other than general terms, how arcing works. Two dissimilar arcing algorithms, arc-fs and arc-x4, give comparable accuracy. It's possible that other arcing algorithms intermediate between arc-fs and arc-x4 will give even better performance. The experiments here, in Freund-Shapire [1995] and in Drucker-Cortes[1995], and in Quinlan[1996] indicate that arcing decision trees may lead to fast and generally accurate classification methods and indicate that additional research aimed at understanding the workings of this class of algorithms will have a high pay-off.

## 2. The Bias and Variance of a Classifier

In classification, the output variable  $y \in \{1, \dots, J\}$  is a class label. The training set  $T$  is of the form  $T = \{(y_n, x_n) \mid n=1, \dots, N\}$  where the  $y_n$  are class labels. Given  $T$ , some method is used to construct a classifier  $C(x, T)$  for predicting future  $y$ -values. Assume that the data in the training set consists of iid selections from the distribution of  $Y, X$ . The misclassification error is defined as:

$$PE(C, T) = P_{X, Y}(C(X, T) \neq Y),$$

and we denote by  $PE(C)$  the expectation of  $PE(C, T)$  over  $T$ . Denote:

$$\begin{aligned} P(j|x) &= P(Y=j \mid X=x) \\ P(dx) &= P(X \in dx) \end{aligned}$$

The minimum misclassification rate is given by the "Bayes classifier  $C^*$ " where

$$C^*(x) = \operatorname{argmax}_j P(j|x)$$

with misclassification rate

$$PE(C^*) = 1 - \int \max_j (P(j|x)) P(dx).$$

Let

$$Q(j|x) = P_T(C(x, T) = j),$$

and define the aggregated classifier as:

$$C_A(x) = \operatorname{argmax}_j Q(j|x).$$

This is aggregation by voting. Consider many independent replicas  $T_1, T_2, \dots$ ; construct the classifiers  $C(x, T_1), C(x, T_2), \dots$ ; and at each  $x$  determine the classification  $C_A(x)$  by having these multiple classifiers vote for the most popular class.

### Definition 2.1

$C(x, \cdot)$  is unbiased at  $x$  if  $C_A(x) = C^*(x)$ .

That is,  $C(x,T)$  is unbiased at  $x$  if, over the replications of  $T$ ,  $C(x,T)$  picks the right class more often than any other class. A classifier that is unbiased at  $x$  is not necessarily an accurate classifier. For instance, suppose that in a two class problem  $P(1|x) = .9$ ,  $P(2|x) = .1$ , and  $Q(1|x) = .6$ ,  $Q(2|x) = .4$ . Then  $C$  is unbiased at  $x$  but the probability of correct classification by  $C$  is  $.6 \times .9 + .4 \times .1 = .58$ . But the Bayes predictor  $C^*$  has probability .9 of correct classification.

If  $C$  is unbiased at  $x$  then  $C_A(x)$  is optimal. Let  $U$  be the set of all  $x$  at which  $C$  is unbiased, and call  $U$  the unbiased set. The complement of  $U$  is called the bias set and denoted by  $B$ . Define

**Definition 2.2**

The bias of a classifier  $C$  is

$$\text{Bias}(C) = P_{\mathbf{X},Y}(C^*(\mathbf{X}) = Y, \mathbf{X} \in B) - E_T P_{\mathbf{X},Y}(C(\mathbf{X},T) = Y, \mathbf{X} \in B)$$

and its variance is

$$\text{Var}(C) = P_{\mathbf{X},Y}(C^*(\mathbf{X}) = Y, \mathbf{X} \in U) - E_T P_{\mathbf{X},Y}(C(\mathbf{X},T) = Y, \mathbf{X} \in U)$$

This leads to the **Fundamental Decomposition**

$$PE(C) = PE(C^*) + \text{Bias}(C) + \text{Var}(C)$$

Note that aggregating a classifier and replacing  $C$  with  $C_A$  reduces the variance to zero, but there is no guarantee that it will reduce the bias. In fact, it is easy to give examples where the bias will be increased. Thus, if the bias set  $B$  has large probability,  $PE(C_A)$  may be significantly larger than  $PE(C)$ . As defined, bias and variance have these properties:

- a) Bias and variance are always non-negative.
- b) The variance of  $C_A$  is zero.
- c) If  $C$  is deterministic, i.e. does not depend on  $T$ , then its variance is zero.
- d) The bias of  $C^*$  is zero.

The proofs of a)-d) are immediate from the definitions. The variance of  $C$  can be expressed as

$$\text{Var}(C) = \int_U [\max_j P(j|x) - \sum_j Q(j|x) P(j|x)] P(dx).$$

The bias of  $C$  is a similar integral over  $B$ . Clearly, both bias and variance are non-negative. Since  $C_A = C^*$  on  $U$ , its variance is zero. If  $C$  is deterministic, then on  $U$ ,  $C = C^*$ , so  $C$  has zero variance. Finally, it's clear that  $C^*$  has zero bias.

Use of the terms bias and variance in classification is somewhat misleading, since they do not have properties commonly associated with bias and variance in predicting numerical outcomes, i.e. regression (see Geman, Bienenstock, and Doursat[1992]). The most important aspect of our definition is that the variance is the component of the classification error that can be eliminated by aggregation. But the bias may be larger for the aggregated classifier than for the unaggregated.

Friedman[1996] contains a thoughtful analysis of the meaning of bias and variance in two class problems. Using some simplifying assumptions, a definition of "boundary bias" at a point  $x$  is given and it is shown that at points of negative boundary bias, classification error can be reduced by reducing variance in the class probability estimates. If the boundary bias is not negative, decreasing the estimate variance may increase the classification error. The points of negative boundary bias are exactly the points defined above as the unbiased set. Other definitions of bias and variance in classification are given in Dietterich and Kong [1995], Kohavi and Wolpert [1996], and Tibshirani [1996]

### 2.3 Instability, Bias, and Variance

Breiman [1996a] pointed out that some prediction methods were unstable in that small changes in the training set could cause large changes in the resulting predictors. I listed trees and neural nets as unstable, nearest neighbors as stable. Linear discriminant analysis (LDA) is also stable. Unstable classifiers are characterized by high variance. As  $T$  changes, the classifiers  $C(x,T)$  can differ markedly from each other and from the aggregated classifier  $C_A(x)$ . Stable classifiers do not change much over replicates of  $T$ , so  $C(x,T)$  and  $C_A(x)$  will tend to be the same and the variance will be small.

Procedures like CART have high variance, but they are "on average, right", that is, they are largely unbiased-- the optimal class is usually the winner of the popularity vote. Stable methods, like LDA, achieve their stability by having a very limited set of models to fit to the data. The result is low variance. But if the data cannot be adequately represented in the available set of models, large bias can result.

### 2.3 Examples

To illustrate, we compute bias and variance of CART for a few examples. These all consist of artificially generated data, since otherwise  $C^*$  cannot be computed nor  $T$  replicated. In each example, the classes have equal probability and the training sets have 300 cases.

i) *waveform*: This is 21 dimension, 3 class data. It is described in the CART book (Breiman et.al [1984]) and code for generating the data is in the UCI repository (ftp address ftp.ics.uci.edu directory pub/machine-learning-databases).

ii) *twonorm*: This is 20 dimension, 2 class data. Each class is drawn from a multivariate normal distribution with unit covariance matrix. Class #1 has mean  $(a,a, \dots, a)$  and class #2 has mean  $(-a,-a, \dots, -a)$ .  $a=2/(20)^{1/2}$

iii) *threenorm*: This is 20 dimension, 2 class data. Class #1 is drawn with equal probability from a unit multivariate normal with mean  $(a,a, \dots, a)$  and from a unit multivariate normal with mean  $(-a,-a, \dots, -a)$ . Class #2 is drawn from a unit multivariate normal with mean at  $(a,-a,a,-a, \dots, a)$ .  $a=2/(20)^{1/2}$

iv) *ringnorm*: This is 20 dimension, 2 class data. Class #1 is multivariate normal with mean zero and covariance matrix 4 times the identity. Class #2 has unit covariance matrix and mean  $(a,a, \dots, a)$ .  $a=1/(20)^{1/2}$

Monte Carlo techniques were used to compute bias and variance for each of the above distributions. One hundred training sets of size 300 were generated and a CART tree grown on each one. An additional set of size 18,000 was generated from the same distribution and the aggregated classifier computed for each point in this latter set. The Bayes predictor can be

analytically computed, so each point in the 18,000 can be assigned to either the bias set or its complement. The results of averaging over this set are in Table 1.

Table 1 Bias, Variance and Error of CART (%)

Data Set	PE(C*)	Bias	Variance	Error
waveform	13.2	1.7	14.1	29.0
twonorm	2.3	.1	19.6	22.1
threenorm	10.5	1.4	20.9	32.8
ringnorm	1.3	1.5	18.5	21.4

These problems are difficult for CART. For instance, in twonorm the optimal separating surface is an oblique plane. This is hard to approximate by the multidimensional rectangles used in CART. In ringnorm, the separating surface is a sphere, again difficult for a rectangular approximation. Threenorm is the most difficult, with the separating surface formed by the continuous join of two oblique hyperplanes. Yet in all examples CART has low bias. The problem is its variance.

We will explore, in the following sections, methods for reducing variance by combining CART classifiers trained on perturbed versions of the training set. In all of the trees that are grown, only the default options in CART are used, i.e. splitting continues as far as possible. No special parameters are set, nor is anything done to optimize the performance of CART on these data sets.

### 3. Bias and Variance for Arcing and Bagging

Given the ubiquitous low bias of tree classifiers, if their variances can be reduced accurate classifiers may result. The general direction toward reducing variance is indicated by the classifier  $C_A(x)$ . This classifier has zero variance and low bias. Specifically, on the four problems above its bias is 2.9, .4, 2.6, 3.4. Thus, it is nearly optimal. Recall that it is based on generating independent replicates of  $T$ , constructing a CART classifier on each replicate training sets, and then letting these classifiers vote for the most popular class. It is not possible, given real data, to generate independent replicates of the training set. But imitations are possible and do work.

#### 3.1 Bagging

The simplest implementation of the idea of generating quasi-replicate training sets is bagging (Breiman[1996a]). Define the probability of the  $n$ th case in the training set to be  $p(n)=1/N$ . Now sample  $N$  times from the distribution  $\{p(n)\}$ . Equivalently, sample from  $T$  *with replacement*. This forms a resampled training set  $T'$ . Cases in  $T$  may not appear in  $T'$  or may appear more than once.  $T'$  is more familiarly called a bootstrap sample from  $T$ .

Denote the distribution on  $T$  given by  $\{p(n)\}$  as  $P^{(B)}$ .  $T'$  is iid from  $P^{(B)}$ . Repeat this sampling procedure, getting a sequence of independent bootstrap training sets. Construct a corresponding sequence of classifiers by using the same classification algorithm applied to each one of the bootstrap training sets. Then let these multiple classifiers vote for each class. For any point  $x$ ,  $C_A(x)$  really depends on the underlying probability  $P$  that the training sets are drawn from i.e.  $C_A(x) = C_A(x, P)$ . The bagged classifier is  $C_A(x, P^{(B)})$ . The hope is that this is a good enough approximation to  $C_A(x, P)$  that considerable variance reduction will result.

### 3.2 Arcing

Arcing is a more complex procedure. Again, multiple classifiers are constructed and vote for classes. But the construction is sequential, with the construction of the  $(k+1)$ st classifier depending on the performance of the  $k$  previously constructed classifiers. We give a brief description of the Freund-Schapire arc-fs algorithm. Details are contained in Section 4.

At the start of each construction, there is a probability distribution  $\{p(n)\}$  on the cases in the training set. A training set  $T$  is constructed by sampling  $N$  times from this distribution. Then the probabilities are updated depending on how the cases in  $T$  are classified by  $C(x,T)$ . A factor  $\beta > 1$  is defined which depends on the misclassification rate--the smaller it is, the larger  $\beta$  is. If the  $n$ th case in  $T$  is misclassified by  $C(x,T)$ , then put weight  $\beta p(n)$  on that case. Otherwise define the weight to be  $p(n)$ . Now divide each weight by the sum of the weights to get the updated probabilities for the next round of sampling. After a fixed number of classifiers have been constructed, they do a weighted voting for the class.

The intuitive idea of arcing is that the points most likely to be selected for the replicate data sets are those most likely to be misclassified. Since these are the troublesome points, focusing on them using the adaptive resampling scheme of arc-fs may do better than the neutral bagging approach.

### 3.3 Results

Bagging and arc-fs were run on the artificial data set described above. For each one, the procedure consisted of generating 100 replicate training sets of size 300. On each training set bagging and arc-fs were run 50 times using CART as the classifier. For each training set this gave an arced classifier and a bagged classifier. An additional 18,000 member Monte Carlo set was generated and the results of aggregating the 100 bagged and arced classifiers computed at each member of this latter set and compared with the Bayes classification. This enabled the bias and variance to be computed. The results are given in Table 2 and compared with the CART results.

Table 2. Bias and Variance (%)

Data Set		CART	Bagging	Arcing
waveform	bias	1.7	1.4	1.0
	var	14.1	5.3	3.6
twonorm	bias	0.1	0.1	1.2
	var	19.6	5.0	1.3
threenorm	bias	1.4	1.3	1.4
	var	20.9	8.6	6.9
ringnorm	bias	1.5	1.4	1.1
	var	18.5	8.3	4.5

Although both bagging and arcing reduce bias a bit, their major contribution to accuracy is in the large reduction of variance. Arcing does better than bagging because it does better at variance reduction.

### 3.4 The effect of combining more classifiers.

The experiments with bagging and arcing above used combinations of 50 tree classifiers. A natural question is what happens if more classifiers are combined. To explore this, we ran arc-fs and bagging on the waveform and twonorm data using combinations of 50, 100, 250 and 500 trees. Each run consisted of 100 repetitions. In each run, a training set of 300 and a test set of 1500 were generated, the prescribed number of trees constructed and combined and the test set error computed. These errors were averaged over 100 repetitions to give the results shown in Table 4. Standard errors were less than 0.16%

Table 3 Test Set Error(%) for 50, 100, 250, 500 Combinations

<u>Data Set</u>	50	100	250	500
waveform				
arc-fs	17.8	17.3	16.6	16.8
bagging	19.8	19.5	19.2	19.2
twonorm				
arc-fs	4.9	4.1	3.8	3.7
bagging	7.3	6.8	6.5	6.5

Arc-fs error rates decrease significantly out to 250 combination, reaching rates close to the Bayes minimums (13.2% for waveform and 2.3% for twonorm). One standard of comparison is linear discriminant analysis, which should be almost optimal for twonorm. It has an error rate of 2.8%, averaged over 100 repetitions. Bagging error rates also decrease out to 250 combinations, but the decreases are smaller than for arc-fs.

## 4. Arcing Algorithms

This section specifies the two arc algorithms and looks at their performance over a number of data sets.

### 4.1. Definitions of the arc algorithms.

Both algorithms proceed in sequential steps with a user defined limit of how many steps until termination. Initialize probabilities  $\{p(n)\}$  to be equal. At each step, the new training set is selected by sampling from the original training set using probabilities  $\{p(n)\}$ . After the classifier based on this resampled training set is constructed, the  $\{p(n)\}$  are updated depending on the misclassifications up to the present step. On termination the classifiers are combined using weighted (arc-fs) or unweighted (arc-x4) voting. The arc-fs algorithm is based on a boosting theorem given in Freund and Schapire [1995]. Arc-x4 is an ad hoc invention.

#### arc-fs specifics:

- i) At the  $k$ th step, using the current probabilities  $\{p^{(k)}(n)\}$ , sample with replacement from  $T$  to get the training set  $T^{(k)}$  and construct classifier  $C_k$  using  $T^{(k)}$ .
- ii) Run  $T$  down the classifier  $C_k$  and let  $d(n)=1$  if the  $n$ th case is classified incorrectly, otherwise zero.
- iii) Define

$$\epsilon_k = \sum_n p^{(k)}(n)d(n), \quad \beta_k = (1 - \epsilon_k)/\epsilon_k$$



and the updated (k+1)st step probabilities by

$$p^{(k+1)}(n) = p^{(k)}(n)\beta_k^{d(n)} / \sum p^{(k)}(n)\beta_k^{d(n)}$$

After K steps, the  $C_1, \dots, C_K$  are combined using weighted voting with  $C_k$  having weight  $\log(\beta_k)$ . Two revisions to this algorithm are necessary. If  $\epsilon_k$  becomes equal to or greater than 1/2, then the original Freund and Schapire algorithm exits from the construction loop. We found that better results were gotten by setting all  $\{p(n)\}$  equal and restarting. If  $\epsilon_k$  equals zero, making the subsequent step undefined, we again set the probabilities equal and restart.

A referee pointed out that the updating definition in arc-fs leads to the interesting result that  $\sum_n p^{(k+1)}(n)d(n) = 1/2$ . That is, the probability weight at the (k+1)st step is equally divided between the points misclassified on the kth step and those correctly classified.

arc-x4 specifics:

- i) Same as for arc-fs
- ii) Run T down the classifier  $C_k$  and let  $m(n)$  be the number of misclassifications of the nth case by  $C_1, \dots, C_k$ .
- iii) The updated k+1 step probabilities are defined by

$$p^{(k+1)}(n) = (1 + m(n)^4) / \sum (1 + m(n)^4)$$

After K steps the  $C_1, \dots, C_K$  are combined by unweighted voting.

After a training set T' is selected by sampling from T with probabilities  $\{p(n)\}$ , another set T'' is generated the same way. T' is used for tree construction, T'' is used as a test set for pruning. By eliminating the need for cross-validation pruning, 50 classification trees can be grown and pruned in about the same cpu time as it takes for 5 trees grown and pruned using 10-fold cross-validation. This is also true for bagging. Thus, both arcing and bagging, applied to decision trees, grow classifiers relatively fast. Parallel bagging can be easily implemented but arc is essentially sequential.

Here is how arc-x4 was devised. After testing arc-fs I suspected that its success lay not in its specific form but in its adaptive resampling property, where increasing weight was placed on those cases more frequently misclassified. To check on this, I tried three simple update schemes for the probabilities  $\{p(n)\}$ . In each, the update was of the form  $1 + m(n)^h$ , and  $h=1,2,4$  was tested on the waveform data. The last one did the best and became arc-x4. Higher values of h were not tested so further improvement is possible.

#### 4.2 Experiments on data sets.

Our experiments used the 6 moderate sized data sets and 4 larger ones used in the bagging paper (Breiman [1996a] plus a handwritten digit data set. The data sets are summarized in Table 4.

Table 4 Data Set Summary

<u>Data Set</u>	#Training	#Test	#Variables	#Classes
heart	1395	140	16	2
breast cancer	699	70	9	2
ionosphere	351	35	34	2
diabetes	768	77	8	2
glass	214	21	9	6
soybean	683	68	35	19
-----				
letters	15,000	5000	16	26
satellite	4,435	2000	36	6
shuttle	43,500	14,500	9	7
DNA	2,000	1,186	60	3
digit	7,291	2,007	256	10

Of the first six data sets, all but the heart data are in the UCI repository. Brief descriptions are in Breiman[1996a]. The procedure used on these data sets consisted of 100 iterations of the following steps:

- i) Select at random 10% of the training set and set it aside as a test set.
- ii) Run 50 steps each of arc-fs and arc-x4 on the remaining 90% of the data.
- iii) Get error rates on the 10% test set.

The error rates computed in iii) are averaged over the 100 iterations to get the final numbers shown in Table 5. We note that using arc-fs in the soybean data set, frequently  $\epsilon_k \geq .5$  causing restarting.

The five larger data sets came with separate test and training sets. Again, each of the arcing algorithms was used to generate 50 classifiers (100 in the digit data) which were then combined into the final classifier. The test set errors are also shown in Table 5.

Table 5 Test Set Error (%)

<u>Data Set</u>	arc-fs	arc-x4	bagging	CART
heart	1.1	1.0	2.8	4.9
breast cancer	3.2	3.3	3.7	5.9
ionosphere	6.4	6.3	7.9	11.2
diabetes	26.6	25.0	23.9	25.3
glass	22.0	21.6	23.2	30.4
soybean	5.8	5.7	6.8	8.6
-----				
letters	3.4	4.0	6.4	12.4
satellite	8.8	9.0	10.3	14.8
shuttle	.007	.021	.014	.062
DNA	4.2	4.8	5.0	6.2
digit	6.2	7.5	10.5	27.1

The first four of the larger data sets were used in the Statlog Project (Michie et.al. 1994) which compared 22 classification methods. Based on their results arc-fs ranks best on three of the four and is barely edged out of first place on DNA. Arc-x4 is close behind.

The digit data set is the famous US Postal Service data set as preprocessed by Le Cun et. al [1990] to result in 16x16 grey-scale images. This data set has been used as a test bed for many adventures in classification at AT&T Bell Laboratories. The best two layer neural net gets 5.9% error rate. A five layer network gets down to 5.1%. Hastie and Tibshirani used deformable prototypes [1994] and get to 5.5% error. Using a very smart metric and nearest neighbors gives the lowest error rate to date--2.7% (P. Simard et. al [1993]). All of these classifiers were specifically tailored for this data.

The interesting Support Vector machines described by Vapnik [1995] are off-the-shelf, but require specification of some parameters and functions. Their lowest error rates are slightly over 4%. Use of the arcing algorithms and CART requires nothing other than reading in the training set, yet arc-fs gives accuracy competitive with the hand-crafted classifiers. It is also relatively fast. The 100 trees constructed in arc-fs on the digit data took about 4 hours of CPU time on a Sparc 20. Some uncomplicated reprogramming would get this down to about one hour of CPU time.

Looking over the test set error results, there is little to choose between arc-fs and arc-x4. Arc-x4 has a slight edge on the smaller data sets, while arc-fs does a little better on the larger ones.

## 5. Properties of the arc algorithms

Experiments were carried out on the six smaller sized data sets listed in table 1 plus the artificial waveform data. Arc-fs and arc-x4 were each given lengthy 1000-step runs on each data set. In each run, information on various characteristics was gathered. We used this information to better understand the algorithms, their similarities and differences. Arc-fs and arc-x4 probably stand at opposite extremes of effective arcing algorithms. In arc-fs the constructed trees change considerably from one step to the next. In arc-x4 the changes are more gradual.

### 5.1 Preliminary Results

Resampling with equal probabilities from a training set, about 37% of the cases do not appear in the resampled data set--put another way, only about 63% of the data is used. With adaptive resampling, more weight is given to some of the cases and less of the data is used. Table 6 gives the average percent of the data used by the arc algorithms in constructing each classifier in the sequence of 1000 steps. The third column is the average value of beta used by the arc-fs algorithm in constructing its sequence.

Table 6 Percent of data Used

<u>Data Set</u>	arc-x4	arc-fs	av. beta
waveform	60	51	5
heart	49	30	52
breast cancer	35	13	103
ionosphere	43	25	34
diabetes	53	36	13
glass	53	38	11
soybean	38	39	17

Arc-x4 data use ranges from 35% to 60%. Arc-fs uses considerably smaller fractions of the data--ranging down to 13% on the breast cancer data set--about 90 cases per tree. The average values of beta are surprisingly large. For instance, for the breast cancer data set, a misclassification of a training set case lead to amplification of its (unnormalized) weight by a factor of 103. The shuttle data (unlisted) leads to more extreme results. On average, only 3.4% of the data is used in constructing each arc-fs tree in the sequence of 50 and the average value of beta is 145,000.

## 5.2 A variability signature

Variability is a characteristic that differed significantly between the algorithms. One signature was derived as follows: In each run, we kept track of the average value of  $N \cdot p(n)$  over the 1000-step runs for each  $n$ . If the  $\{p(n)\}$  were equal, as in bagging, these average values would be about 1.0. The standard deviation of  $N \cdot p(n)$  for each  $n$  was also computed. Figure 1 gives plots of the standard deviations vs. the averages for the six data sets and for each algorithm. The upper point cloud in each graph corresponds to the arc-fs values; the lower to the arc-x4 values. The graph for the soybean data set is not shown because the frequent restarting causes the arc-fs values to be anomalous.

Figure 1

For arc-fs the standard deviations of  $p(n)$  is generally larger than its average, and increase linearly with the average. The larger  $p(n)$ , the more volatile it is. In contrast, the standard deviations for arc-x4 are quite small and only increase slowly with average  $p(n)$ . Further, the range of  $p(n)$  for arc-fs is 2-3 times larger than for arc-x4. Note that, modulo scaling, the shapes of the point sets are similar between data sets.

## 5.3 A mysterious signature

In the 1000-step runs, we also kept track of the number of times the  $n$ th case appeared in a training set and the number of times it was misclassified (whether or not it was in the training set). For both algorithms, the more frequently a point is misclassified, the more its probability increases, and the more frequently it will be used in a training set. This seems intuitively obvious, so we were mystified by the graphs of figure 2.

Figure 2

For each data set, number of times misclassified was plotted vs. number of times in a training set. The plots for arc-x4 behave as expected. Not so for arc-fs. Their plots rise sharply to a plateau. On this plateau, there is almost no change in misclassification rate vs. rate in training set. Fortunately, this mysterious behavior has a rational explanation in terms of the structure of the arc-fs algorithm.

Assume that there are  $K$  iterations and that  $\beta_k$  is constant equal to  $\beta$  (in our experiments, the values of  $\beta_k$  had moderate sd/mean values for  $K$  large). For each  $n$ , let  $r(n)$  be the proportion of times that the  $n$ th case was misclassified. Then

$$p(n) \cong \beta \sum K r(n) / \sum \beta K r(n)$$

Let  $r^* = \max_n r(n)$ ,  $L$  the set of indices such that  $r(n) > r^* - \epsilon$  for a small positive  $\epsilon$ , and  $|L|$  the cardinality of  $L$ . If  $|L|$  is too small, then there will be an increasing number of misclassifications for those cases not in  $L$  that are not accurately classified by training sets drawn from  $L$ . Thus, their misclassification rates will increase until they get close to  $r^*$ . To

illustrate this, Figure 3 shows the misclassification rates as a function of the number of iterations for two cases in the twonorm data discussed in the next subsection. The top curve is for a case with consistently large  $p(n)$ . The lower curve is for a case with  $p(n)$  almost vanishingly small.

Figure 3

There are also a number of cases that are more accurately classified by training sets drawn from  $L$ . These are characterized by lower values of the misclassification rate, and by small  $p(n)$ . That is, they are the cases that cluster on the y-axes of figure 2. More insight is provided by Figure 4. This is a percentile plot of the proportion of the training sets that the 300 cases of the twonorm data are in (10,000 iterations). About 40% of the cases are in a very small number of the training sets. The rest have a uniform distribution across the proportion of training sets.

Figure 4

#### 5.4 Do hard-to classify points get more weight?

To explore this question, we used the twonorm data. The ratio of the probability densities of the two classes at the point  $x$  depends only on the value of  $|\langle x, \mathbf{1} \rangle|$  where  $\mathbf{1}$  is the vector whose coordinates are all one. The smaller  $|\langle x, \mathbf{1} \rangle|$  is, the closer the ratio of the two densities to one, and the more difficult the point  $x$  is to classify. If the idea underlying the arc algorithms is valid, then the probabilities of inclusion in the resampled training sets should increase as  $|\langle x, \mathbf{1} \rangle|$  decreases. Figure 5 plots the average of  $p(n)$  over 1000 iterations vs.  $|\langle x(n), \mathbf{1} \rangle|$  for both arc algorithms.

Figure 5

While  $\text{av}(p(n))$  generally increases with decreasing  $|\langle x(n), \mathbf{1} \rangle|$  the relation is noisy. It is confounded by other factors that I have not yet been able to pinpoint.

### 6. Linear Discriminant Analysis Isn't Improved by Bagging or Arcing.

Linear discriminant analysis (LDA) is fairly stable with low variance and it should come as no surprise that its test set error is not significantly reduced by use of bagging or arcing. Here our test bed was four of the first six data sets of Table 1. Ionosphere and soybean were eliminated because the within class covariance matrix was singular, either for the full training set (ionosphere) or for some of the bagging or arc-fs training sets (soybean).

The experimental set-up was similar to that used in Section 2. Using a leave-out-10% as a test set, 100 repetitions were run using linear discriminant analysis alone and the test set errors averaged. Then this was repeated, but in every repetition, 25 combinations of linear discriminants were built using bagging or arc-fs. The test set errors of these combined classifiers were also averaged. The results are listed in Table 7.

Table 7 Linear Discriminant Test Set Error(%).

<u>Data Set</u>	LDA	LDA: bag	LDA: arc	Restart Freq.
heart	25.8	25.8	26.6	1/9
breast cancer	3.9	3.9	3.8	1/8
diabetes	23.6	23.5	23.9	1/9
glass	42.2	41.5	40.6	1/5

Recall that for arc-fs, if  $\epsilon_k \geq .5$ , then the construction was restarted with equal  $\{p(n)\}$ . The last column of Table 7 indicates how often restarting occurred. For instance, in the heart data, on the average, it occurred about once every 9 times. In contrast, in the runs combining trees restarting was encountered only on the soybean data. The frequency of restarting was also a consequence of the stability of linear discriminant analysis. If the procedure is stable, the same cases tend to be misclassified even with the changing training sets. Then their weights increase and so does the weighted training set error.

These results indicate that linear discriminant analysis is generally a low variance procedure. It fits a simple parametric normal model that does not change much with replicate training sets. To illustrate, we did a Monte Carlo estimation of bias and variance using the synthetic threennorm data. Recall that bias and variance for CART are 1.4% and 20.9%. For LDA they are 4.0% and 2.9%. The problem in LDA is usually bias--when it is wrong, it is consistently wrong, and with a simple model there is no hope of low bias across a variety of complex data sets.

## 7 Improved Bagging

The aggregate classifier depends on the distribution  $P$  that the samples are selected from and the number  $N$  selected. Letting the dependence on  $N$  be implicit, denote  $C_A = C_A(x, P)$ . As mentioned in 3.1, bagging replaces  $C_A(x, P)$  by  $C_A(x, P^{(B)})$  with the hope that this approximation is good enough to produce variance reduction. Now  $P^{(B)}$ , at best, is a discrete estimate for a distribution  $P$  that is usually smoother and more spread out than  $P^{(B)}$ . An interesting question is what a better approximation to  $P$  might produce.

To check this possibility, we used the four simulated data sets described in section 3. Once a training set was drawn from one of these distributions, we replaced each  $x_n$  by a spherical normal distribution centered at  $x_n$ . The bootstrap training set  $T^{(B)}$  is iid drawn from this smoothed distribution. Two or three values were tried for the sd of the normal smoothing and the best one adopted. The results are given in Table 9.

Table 9 Smoothed P-Estimate Bagging--Test Set Errors(%)

<u>Data Set</u>	Bagging	Bagging (smoothed)	Arcing
waveform	19.8	18.4	17.8
twonorm	7.4	5.5	4.8
threennorm	20.4	18.6	18.8
ringnorm	11.0	8.7	6.9

The PE values for the smoothed P-estimates show that the better the approximation to  $P$ , the lower the variance. But there are limits to how well we can estimate the unknown underlying distribution from the training set. The aggregated classifiers based on the smoothed approximations had variances significantly above zero, and we doubt that efforts to refine the  $P$  estimates will push them much lower. But note that even with the better  $P$  approximation bagging does not do as well as arcing.

## 8 More on Arcing

### 8.1 Training Set Error

Arcing is much less transparent than bagging. Freund and Schapire [1995] designed arc-fs to drive training set error rapidly to zero, and it does remarkably well at this. But the context in

which arc-fs was designed gives no clues as to its ability to reduce test set error. For instance suppose we run arc-fs but exit the construction loop when the training set error becomes zero. The test set errors and number of steps to exit the loop (averaged over 100 iterations of leave out 10% for the smaller data sets) are given in Table 10 and compared to the stop at k=50 results from Table 2.

For the smaller data sets, we also kept track of the last step in at which the training set error was non-zero. Their averages are given in the figures in parentheses in the last column of table 10. These values show that once the training set error drops to zero, it almost invariably stays at zero. We also ran bagging on the first six data sets in Table 5, exiting the loop when the training error was zero, and kept track of the average number of steps to exit and the average test set error over 100 repetitions of leave out 10%. These numbers are given in Table 11 (soybean was not used because of restarting problems).

Table 10 Test Error(%) and Exit Times for Arc-fs

<u>Data Set</u>	stop: k=50	stop: error=0	steps to exit
heart	1.1	5.3	3.0 (3.1)
breast cancer	3.2	4.9	3.0 (3.0)
ionosphere	6.4	9.1	3.0 (3.1)
diabetes	26.6	28.6	5.0 (5.1)
glass	22.0	28.1	4.9 (5.1)
-----			
letters	3.4	7.9	5
satellite	8.8	12.6	5
shuttle	.007	.014	3
DNA	4.2	6.4	5

Table 11 Test Error(%)and Exit Times for Bagging

<u>Data Set</u>	stop: k=50	stop: error=0	steps to exit
heart	2.8	3.0	15
breast cancer	3.7	4.1	55
ionosphere	7.9	9.2	38
diabetes	23.9	24.7	45
glass	23.2	25.0	22

These results delineate the differences between efficient reduction in training set error and test set accuracy. Arc-fs reaches zero training set error very quickly, after an average of 5 tree constructions (at most). But the accompanying test set error is higher than that of bagging, which takes longer to reach zero training set error. To produce optimum reductions in test set error, arc-fs must be run far past the point of zero training set error.

## 8.2 Non-random Arcing

An important question in understanding arcing is what happens if the randomness is taken out. That is, in the definition of arc-fs in Section 4.1, everything remains the same except that instead of randomly sampling from the probabilities  $\{p(n)\}$ , these probabilities are fed directly into the classifier such that  $p(n)$  is the weight for the case  $(y_n, x_n)$ . CART was modified to use weights and run using a minimum node size of 10 on all data sets and the same test set error estimation methods used to get the Table 5 results. The results are given in Table 12 as compared to random arcing.

Table 12 Test Set Error, Arc-fs (%)

<u>Data Set</u>	<u>random</u>	<u>non-random</u>
heart	1.1	.33
breast cancer	3.2	3.0
ionosphere	6.4	5.7
diabetes	26.6	25.7
glass	22.0	22.6
soybean	5.8	5.7
-----		
letters	3.4	2.8
satellite	8.8	9.0
shuttle	.007	.014
DNA	4.2	4.5
digit	6.2	7.1

These results show that randomness is not an important element in arcing's ability to reduce test set error. This sharply differentiates arcing from bagging. In bagging the random sampling is critical, and the weights remain constant and equal.

### 8.3 Remarks

The arcing classifier is not expressible as an aggregated classifier based on some approximation to  $P$ . The distributions from which the successive training sets are drawn change constantly as the procedure continues. For the arc-fs algorithm, the successive  $\{p(n)\}$  form a multivariate Markov chain. Suppose they have a stationary distribution  $\pi(d\mathbf{p})$  (which I suspect is true, but a referee doubts). Let  $Q(j | \mathbf{x}, \mathbf{p}) = P_T(C(\mathbf{x}, T)=j)$ , where the probability  $P_T$  is over all training sets drawn from the original training set using the distribution  $\mathbf{p}$  over the cases. Then, in steady-state with unweighted voting, class  $j$  gets vote  $Q(j | \mathbf{x}, \mathbf{p}) \pi(d\mathbf{p})$ .

It is not clear how this steady-state probability structure relates to the error-reduction properties of arcing. But its importance is suggested by our experiments. The results in Table 3 show that arcing takes longer to reach its minimum error rate than bagging. If the error reduction properties of arcing come from its steady-state behavior, then this longer reduction time may reflect the fact that the dependent Markov property of the arc-fs algorithm takes longer to reach steady-state than bagging in which there is independence between the successive bootstrap training sets and the Law of Large Numbers sets in quickly. But how the steady-state behavior of arcing algorithms relates to their ability to drive the training set error to zero in a few iterations is unknown.

Another complex aspect of arcing is illustrated in the experiments done to date. In the diabetes data set it gives higher error rate than a single run of CART. The Freund-Schapire[1996] and Quinlan[1996] experiments used C4.5, a tree structured program similar to CART and compared C4.5 to the arc-fs and bagging classifiers based on C4.5. In 5 of the 39 data sets examined in the two experiments, the arc-fs test set error was over 20% larger than that of C4.5. This did not occur with bagging. Its not understood why arc-fs causes this infrequent degeneration in test set error, usually with smaller data sets. One conjecture is that this may be caused by outliers in the data. An outlier will be consistently misclassified, so that its probability of being sampled will continue to increase as the arcing continues. It will then start appearing multiple times in the resampled data sets. In small data sets, this may be enough to warp the classifiers.



### 7.3 Future Work

Arc-fs and other arcing algorithms can reduce the test set error of methods like CART to the point where they are the most accurate available off-the-shelf classifiers on a wide variety of data sets. The Freund-Schapire discovery of adaptive resampling as embodied in arc-fs is a creative idea which should lead to interesting research and better understanding of how classification works. The arcing algorithms have a rich probabilistic structure and it is a challenging problem to connect this structure to their error reduction properties. It is not clear what an optimum arcing algorithm would look like. Arc-fs was devised in a different context and arc-x4 is ad-hoc. While the concepts of bias and variance seem suitable as an explanation for bagging, and while arcing does effectively reduce variance, the bias-variance setting does not give a convincing explanation for why arcing works. Better understanding of how arcing functions will lead to further improvements.

### 9. Acknowledgments

I am indebted to Yoav Freund for giving me the draft papers referred to in this article and to both Yoav Freund and Robert Schapire for informative email interchanges and help in understanding the boosting context: To Trevor Hastie for making available the preprocessed US Postal Service data: to Harris Drucker who responded generously to my questioning at NIPS95 and whose subsequent work on comparing arc-fs to bagging convinced me that arcing needed looking into: to Tom Dietterich for his comments on the first draft of this paper; and to David Wolpert for helpful discussions about boosting. An Associate Editor and the referees made interesting and constructive comments that resulted in an improved paper. In particular, I'm indebted to one of the referees for suggesting the experiment with non-random arcing reported in Section 8.2

### References

- Ali, K. [1995] Learning Probabilistic Relational Concept Descriptions, Thesis, Computer Science, University of California, Irvine
- Breiman, L. [1996a] Bagging predictors, *Machine Learning* 26, No. 2, pp. 123-140
- Breiman, L. [1996b] The heuristics of instability in model selection, *Annals of Statistics*, 24, pp. 2350-2383
- Breiman, L., Friedman, J., Olshen, R., and Stone, C. [1984] *Classification and Regression Trees*, Chapman and Hall
- Dietterich, T.G. and Kong, E. B. [1995] Error-Correcting Output Coding Corrects Bias and Variance, *Proceedings of the 12th International Conference on Machine Learning* pp. 313-321.
- Drucker, H. and Cortes, C. [1995] Boosting decision trees, *Advances in Neural Information Processing Systems* Vol. 8, pp. 479-485.
- Freund, Y. and Schapire, R. [1995] A decision-theoretic generalization of on-line learning and an application to boosting. to appear, *Journal of Computer and System Sciences*.
- Freund, Y. and Schapire, R. [1996] Experiments with a new boosting algorithm, *Machine Learning: Proceedings of the Thirteenth International Conference*, pp. 148-156

- Friedman, J. H. [1996] On Bias, Variance, 0/1-loss, and the Curse of Dimensionality, to appear, Journal of Knowledge Discovery and Data Mining.
- Geman, S., Bienenstock, E., and Doursat, R.[1992] Neural networks and the bias/variance dilemma. Neural Computations 4, pp. 1-58
- Hastie, T. and Tibshirani, R. [1994] Handwritten digit recognition via deformable prototypes, (unpublished ms. ftp stat.stanford.edu/pub/hastie/zip.ps.Z)
- Kearns, M. and Valiant, L.G.[1988] Learning Boolean Formulae or Finite Automata is as Hard as Factoring, Technical Report TR-14-88, Harvard University Aiken Computation Laboratory
- Kearns, M. and Valiant, L.G.[1989] Cryptographic Limitations on Learning Boolean Formulae and Finite Automata. Proceedings of the Twenty-First Annual ACM Symposium on Theory of Computing , ACM Press, pp. 433-444.
- Kohavi, R. and Wolpert, D.H.[1996] Bias Plus Variance Decomposition for Zero-One Loss Functions, Machine Learning: Proceedings of the Thirteenth International Conference, pp. 275-283
- Le Cun, Y. Boser, B., Denker, J., Henderson, D., Howard, R.,Hubbard, W. and Jackel, L. [1990], Handwritten digit recognition with a back-propagation network, Advances in Neural Information Processing Systems, Vol.2, pp. 396-404
- Michie, D., Spiegelhalter, D. and Taylor, C. [1994] Machine Learning, Neural and Statistical Classification, Ellis Horwood, London
- Quinlan, J.R.[1996] Bagging, Boosting, and C4.5, Proceedings of AAAI'96 National Conference, on Artificial Intelligence, pp. 725-730.
- Schapire, R.[1990] The Strength of Weak Learnability, Machine Learning, 5, pp. 197-227
- Simard, P., Le Cun, Y., and Denker, J., [1993] Efficient pattern recognition using a new transformation distance, Advances in Neural Information Processing Systems Vol. 5,pp.50-58
- Tibshirani, R [1996] Bias, Variance, and Prediction Error for Classification Rules, Technical Report, Statistics Department, University of Toronto
- Vapnik, V. [1995] The Nature of Statistical Learning Theory, Springer

### **Appendix The Boosting Context of Arc-fs**

Freund and Schapire [1995] designed arc-fs to drive the training error rapidly to zero. They connected this training set property with the test set behavior in two ways. The first was based on structural risk minimization (see Vapnik[1995]). The idea here is that bounds on the test set error can be given in terms of the training set error where these bounds depend on the VC-dimension of the class of functions used to construct the classifiers. If the bound is tight this approach has a contradictory consequence. Since stopping as soon as the training error is zero gives the least complex classifier with the lowest VC dimension, then the test set error corresponding to this stopping rule should be lower than if we continue to combine classifiers. Table 10 shows that this does not hold.

The second connection is through the concept of boosting. Freund and Schapire [1995] devised arc-fs in the context of boosting theory (see Schapire[1990]) and named it Adaboost. We follow Freund[1995] in setting out the definitions: Assume that there is an input space of vectors  $\mathbf{x}$  and an unknown function  $C_0(\mathbf{x}) \in \{0,1\}$  defined on the space of input vectors  $\mathbf{x}$  that assigns a class label to each input vector. The problem is to "learn"  $C_0$ .

A classifying method is called a weaklearner if there exist  $\epsilon > 0, \delta > 0$  and integer  $N$  such that given a training set  $T$  consisting of  $x_1, x_2, \dots, x_N$  drawn at random from any distribution  $P(\mathbf{dx})$  on input space together with the corresponding  $y_n = C_0(x_n), n=1, \dots, N$  and the classifier  $C(\mathbf{x}, T)$  constructed, then the probability of a  $T$  such that  $P(C(\mathbf{X}, T) \neq C_0(\mathbf{X}) | T) < 0.5 - \epsilon$  is greater than  $\delta$ , where  $\mathbf{X}$  is a random vector having distribution  $P(\mathbf{dx})$ .

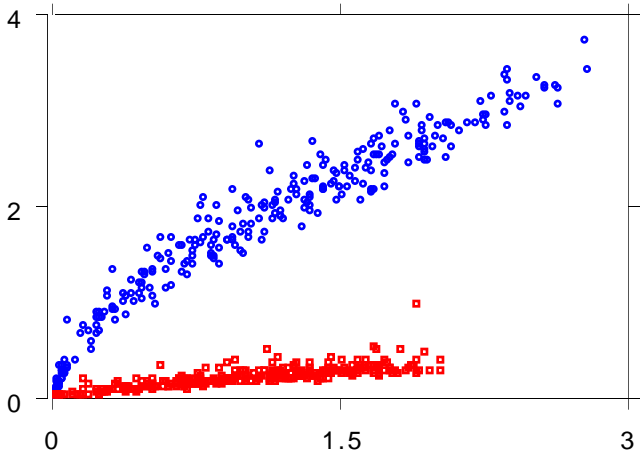
A classifying method is called a stronglearner if for any  $\epsilon > 0, \delta > 0$  there is an integer  $N$  such that if it is given a training set  $T$  consisting of  $x_1, x_2, \dots, x_N$  drawn at random from any distribution  $P(\mathbf{dx})$  on input space together with the corresponding  $y_n = C_0(x_n), n=1, \dots, N$ , and the classifier  $C(\mathbf{x}, T)$  constructed, then the probability of a  $T$  such that  $P(C(\mathbf{X}, T) \neq C_0(\mathbf{X}) | T) > \epsilon$  is less than  $\delta$ , where  $\mathbf{X}$  is a random vector having distribution  $P(\mathbf{dx})$ .

Note that a stronglearner has low error over the whole input space, not just the training set--i.e. it has small test set error. The concept of weak learning was introduced by Kearns and Valiant[1988], [1989], who left open the question of whether weak and strong learnability are equivalent. The question was termed the *boosting problem* since equivalence requires the method to boost the low accuracy of a weaklearner to the high accuracy of a stronglearner. Schapire[1990] proved that boosting is possible. A boosting algorithm is a method that takes a weaklearner and converts it into a stronglearner. Freund [1995] proved that an algorithm similar to arc-fs is boosting. Freund and Schapire [1995] apply the results in Freund[1995] and conclude that Adaboost is boosting.

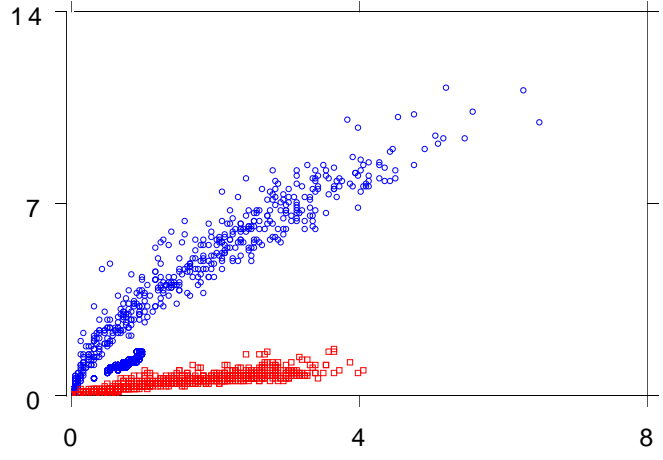
The boosting assumptions are restrictive. For instance, if there is any overlap between classes (if the Bayes error rate is positive) then there are no weak or strong learners. Even if there is no overlap between classes, it is easy to give examples of input spaces and  $C_0$  such that there are no weak learners. The boosting theorems really say "if there is a weak learner, then..." but in virtually all of the real data situations in which arcing or bagging is used, there is overlap between classes and no weak learners exist. Thus the Freund[1995] boosting theorem is not applicable. In particular, it is not applicable in all of the examples of simulated data used in this paper and most, if not all, of the examples of real data sets used in this paper, in Freund and Schapire[1996], and in Quinlan[1996].

While there may be a connection between the ability of arcing algorithms to rapidly drive training set error to zero and their steady-state test set reduction, it is not rooted in the boosting context.

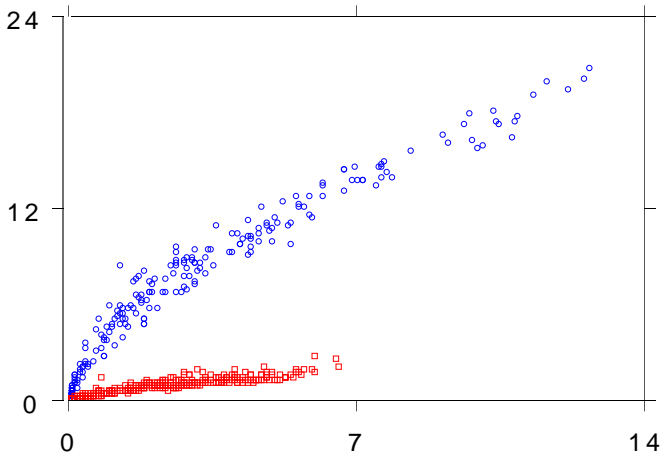
Waveform



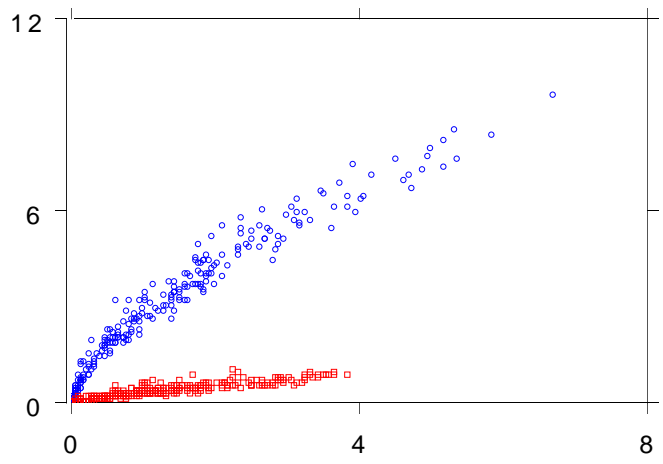
Heart



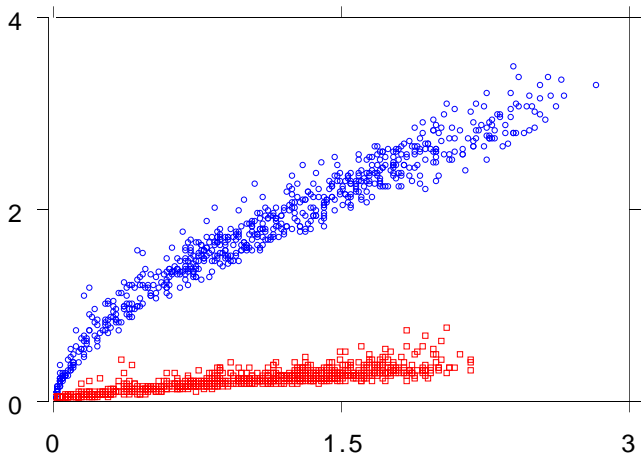
Breast Cancer



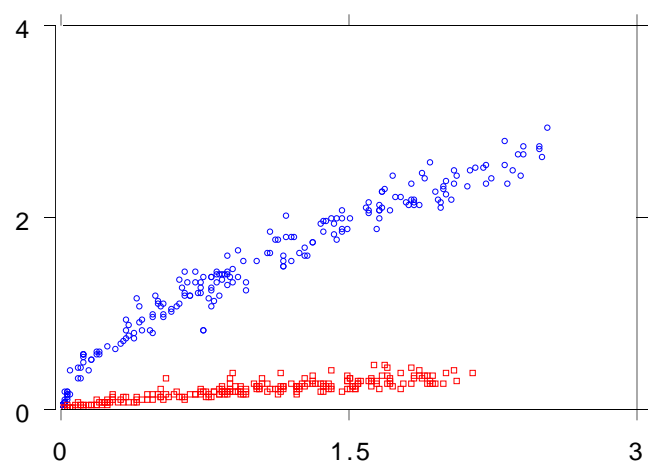
Ionosphere



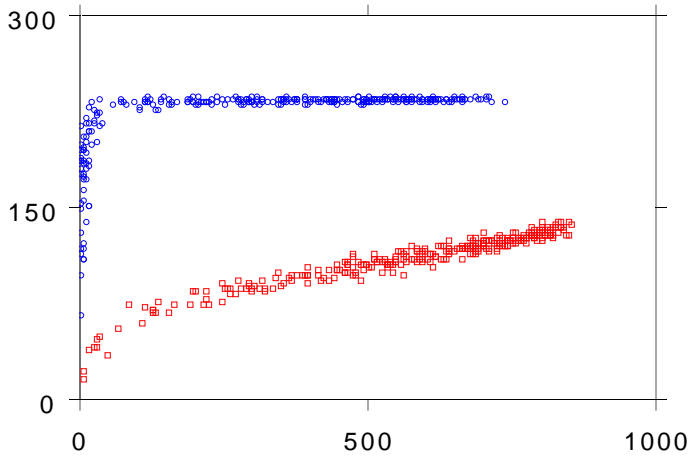
Diabetes



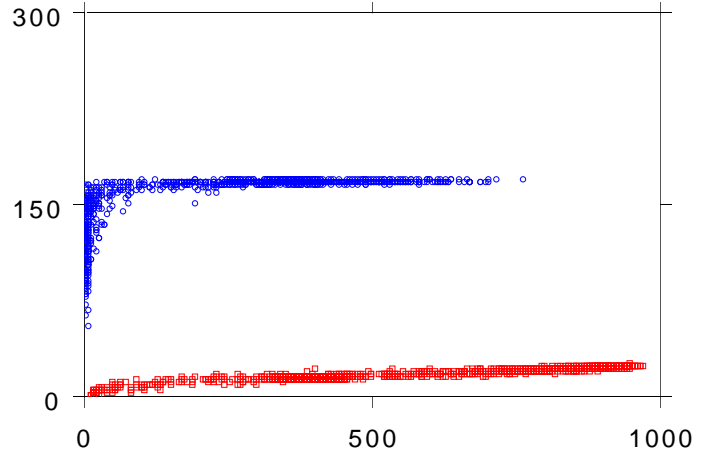
Glass



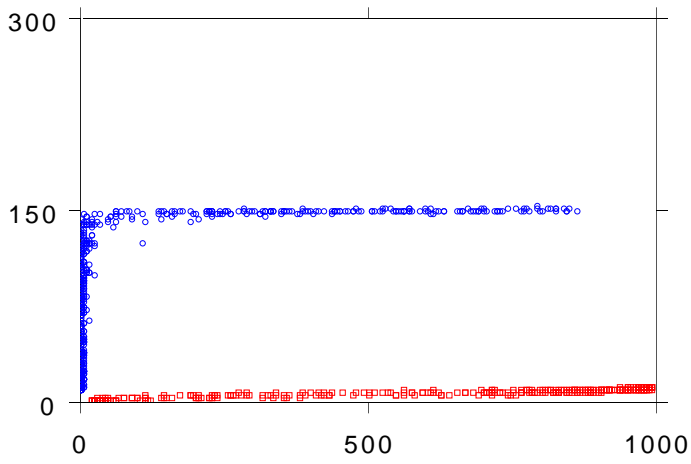
Waveform



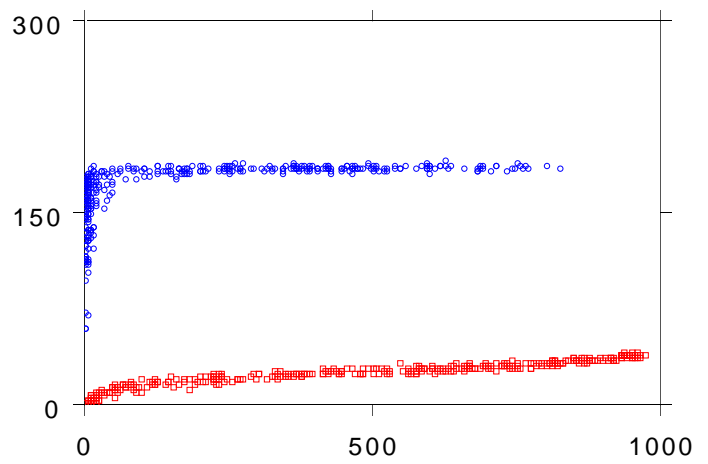
Heart



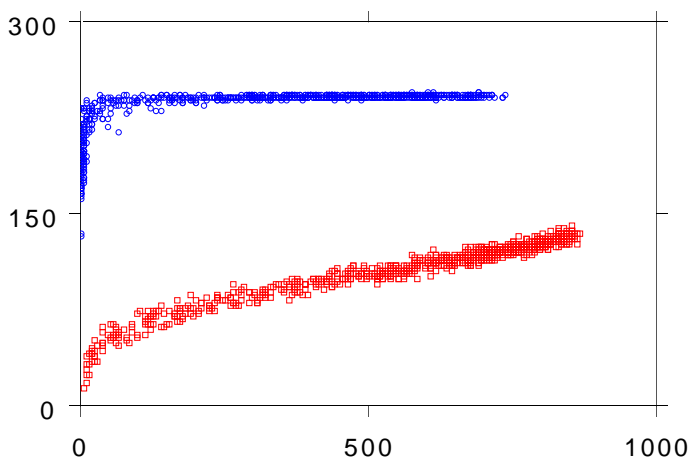
Breast Cancer



Ionosphere



Diabetes



Glass

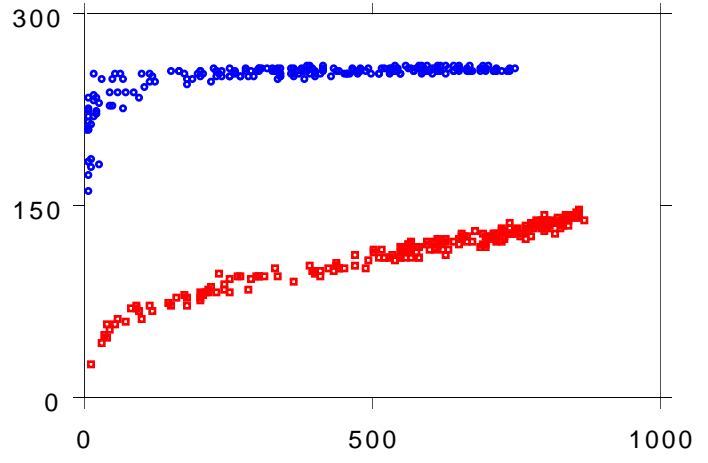


FIGURE 3 Proportion of Times Misclassified for Two Cases

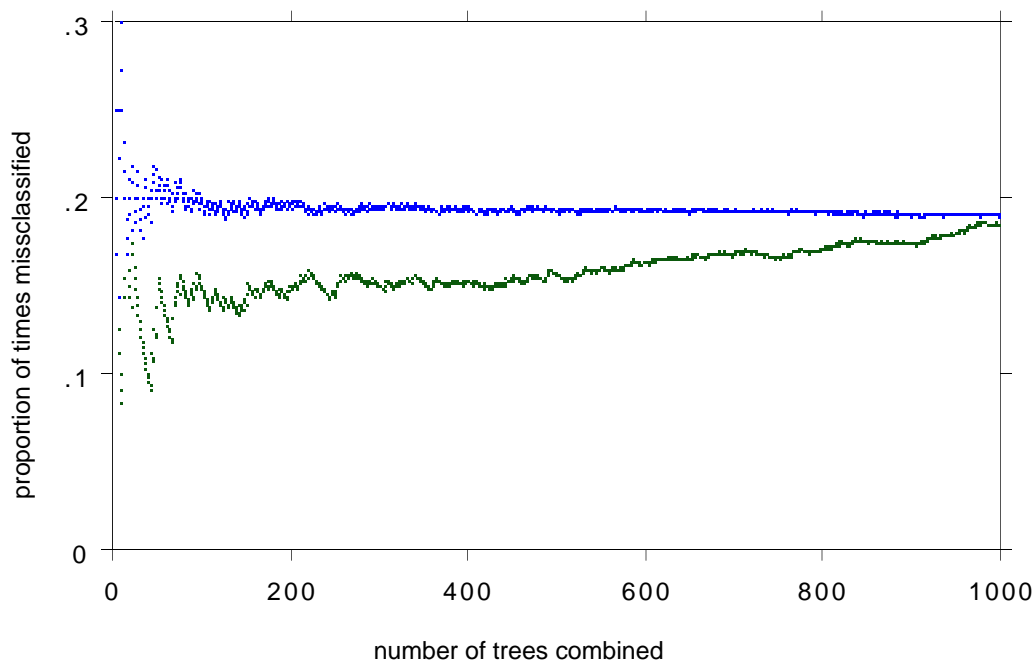
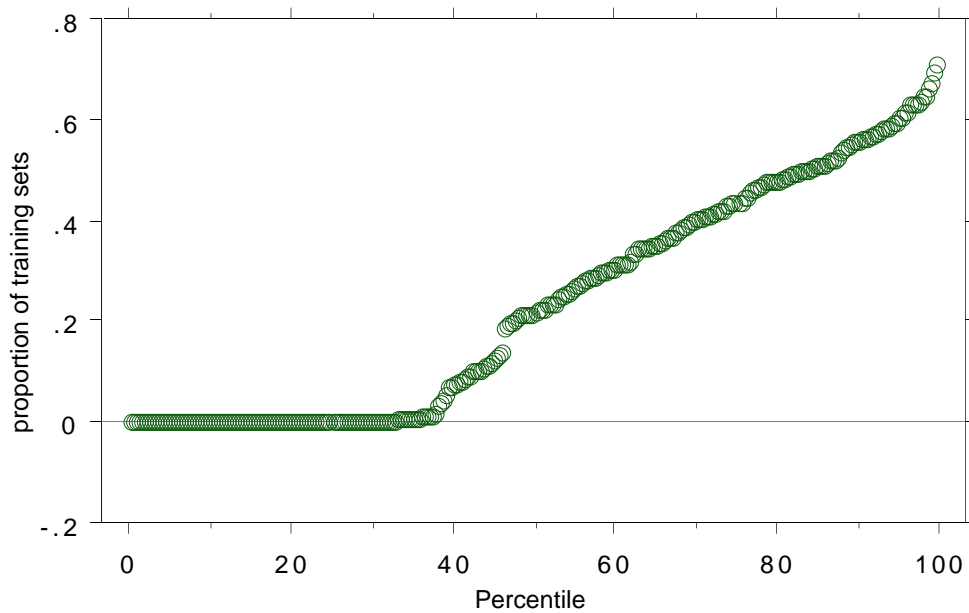
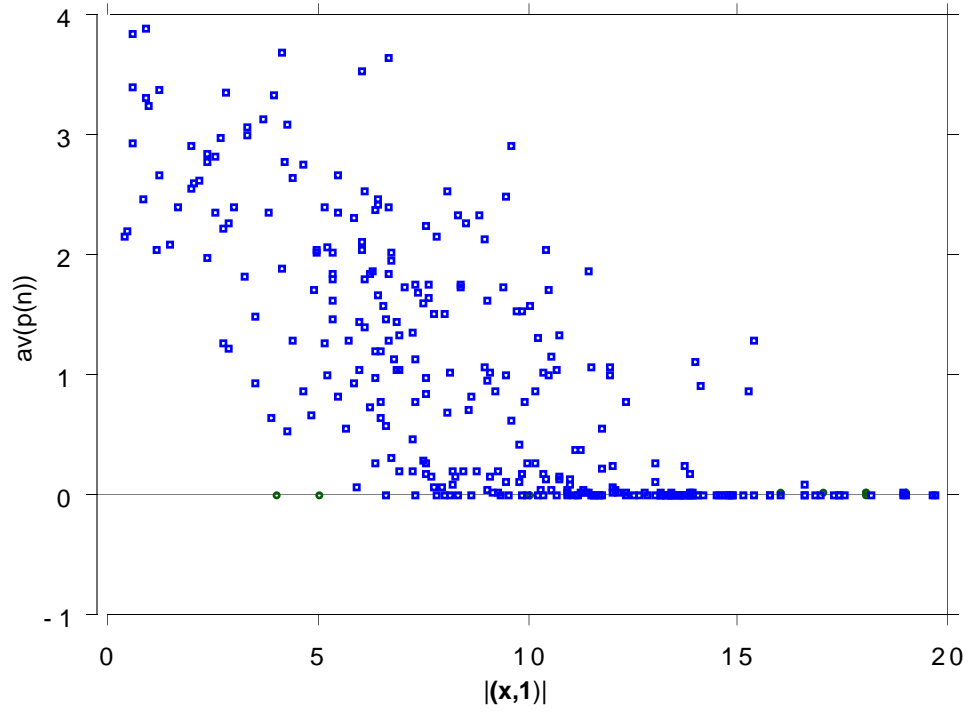


FIGURE 4 Percentile Plot--Proportion of Training Sets that Cases are In



ARC-FS



ARC-X4

