



Random Forests

LEO BREIMAN

Statistics Department, University of California, Berkeley, CA 94720

Editor: Robert E. Schapire

Abstract. Random forests are a combination of tree predictors which each tree depends on the value of a random vector sampled independently and at the same distribution for all trees in the forest. The generalization error for forests converges as well to a limit as the number of trees in the forest becomes large. The generalization error of a forest of tree classifiers depends on the strength of the individual trees in the forest and the correlation between them. Using a random selection of features to split each node yield error rates which compare favorably to AdaBoost (Y. Freund & R. Schapire, *Machine Learning: Proceedings of the Thirteenth International conference*, ***, 148–156), but are more robust with respect to noise. Internal variable importance, strength, and correlation and the error are estimated by bootstrap sampling increasing the number of features used in the splitting. Internal variable importance are also estimated to measure variable importance. These ideas are also applicable to regression.

Keywords: classification, regression, ensemble

1. Random forests

1.1. Introduction

Significant improvements in classification accuracy have led from grouping an ensemble of trees and leaving them out for the most popular classifier. In order to group the ensemble, often random forests are generated that govern the growth of each tree in the ensemble. An early example is bagging (Breiman, 1996), where to group each tree a random selection (with replacement) is made from the example in the training set.

Another example is random splitting (Dietterich, 1998), where at each node the split is selected at random from among the K best splits. Breiman (1999) generates new training sets by randomizing the output in the original training set. Another approach is to select the training set from a random set of weights on the example in the training set. Ho (1998) has written a number of papers on the random subspace method which does a random selection of a subset of features to split each tree.

In an important paper on written character recognition, Amit and Geman (1997) define a large number of geometric features and search over a random selection of the features for the best split at each node. This latter paper has been influential in my thinking.

The common element in all of the procedures is that for the k th tree, a random forest Θ_k is generated, independent of the previous random forest $\Theta_1, \dots, \Theta_{k-1}$ but with the same distribution; and a tree is grown using the training set and Θ_k , resulting in a classifier $h(\mathbf{x}, \Theta_k)$ where \mathbf{x} is an input vector. For instance, in bagging the random forest Θ is

generated at the cost in N by selecting from N drawn at random at the bootstrap, where N is number of example in the training set. In random pilot election Θ consists of a number of independent random integer between 1 and K . The nature and dimensionality of Θ depend on it's three condition.

After a large number of tree is generated, the one for the most popular class. We call the procedure **random forests**.

Definition 1.1. A random forest is a classifier consisting of a collection of tree-structured classifiers $\{h(\mathbf{x}, \Theta_k), k = 1, \dots\}$ where the $\{\Theta_k\}$ are independent identically distributed random vector and each tree can a node for the most popular class input \mathbf{x} .

1.2. Outline of paper

Section 2 gives some theoretical background for random forest. Use of the Strong Law of Large Number shows that the average of the bootstrap is not a problem. We give a simplified and extended version of the Amit and Geman (1997) analysis to show that the accuracy of a random forest depends on the strength of the individual tree classifier and a measure of the dependence between them (see Section 2 for definition).

Section 3 introduces forests using the random election of feature at each node to determine the pilot. An important question is how many features to elect at each node. For guidance, internal estimate of the generalization error, classifier strength and dependence are computed. The are called out-of-bag estimate and are reviewed in Section 4. Section 5 and 6 give empirical results for two different form of random feature. They run random election from the original input; the second is random linear combination of input. The results compare favorably to Adaboost.

The results aim to be independent of the number of features elected to pilot each node. Usually, electing one or two features give near optimal results. To explore this and relate it to strength and correlation, an empirical study is carried out in Section 7.

Adaboost has no random elements and group an ensemble of tree bootstrap weights of the training set where the current weight depend on the path or the ensemble formation. By just a deterministic random number generator can give a good imitation of randomness, my belief is that in later stage Adaboost is emulating a random forest. Evidence for this conjecture is given in Section 8.

Important recent problem, i.e., medical diagnosis and documentation retrieval, often have the property that there are many input variables, often in the hundred or so and, with each one containing only a small amount of information. A single tree classifier will then have accuracy only slightly better than a random choice of classifier. By combining three grouping random feature can produce improved accuracy. In Section 9 we perform on a imbalanced data set with 1,000 input variable, 1,000 example in the training set and a 4,000 example test set. Accuracy comparable to the base rate is achieved.

In many applications, understanding of the mechanism of the random forest black box is needed. Section 10 make a start on this by computing internal estimate of variable importance and binding the together better known.

Section 11 looks at random forests for regression. A bound for the mean squared generalization error is derived which shows the decrease in error from the individual classifier in the forest depends on the correlation between individual and the mean squared error of the individual classifier. Empirical results for regression are in Section 12. Concluding remarks are given in Section 13.

2. Characterizing the accuracy of random forests

2.1. Random forests converge

Given an ensemble of classifiers $h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_K(\mathbf{x})$, and with the training set drawn at random from the distribution of the random vector \mathbf{Y}, \mathbf{X} , define the margin function as

$$mg(\mathbf{X}, Y) = av_k I(h_k(\mathbf{X}) = Y) - \max_{j \neq Y} av_k I(h_k(\mathbf{X}) = j).$$

where $I(\cdot)$ is the indicator function. The margin measures the extent to which the average number of votes at \mathbf{X}, Y for the right class exceed the average vote for another class. The larger the margin, the more confidence in the classification. The generalization error is given by

$$PE^* = P_{\mathbf{X}, Y}(mg(\mathbf{X}, Y) < 0)$$

where the bracket \mathbf{X}, Y indicates that the probability is over the \mathbf{X}, Y space.

In random forests, $h_k(\mathbf{X}) = h(\mathbf{X}, \Theta_k)$. For a large number of trees, it follows from the Strong Law of Large Numbers and theorems such that:

Theorem 1.2. *As the number of trees increases, for almost surely all sequences Θ_1, \dots, PE^* converges to*

$$P_{\mathbf{X}, Y}(P_\Theta(h(\mathbf{X}, \Theta) = Y) - \max_{j \neq Y} P_\Theta(h(\mathbf{X}, \Theta) = j) < 0). \quad (1)$$

Proof: See Appendix I. □

This result implies that random forests do not overfit as more trees are added, but produce a limiting value of the generalization error.

2.2. Strength and correlation

For random forests, an upper bound can be derived for the generalization error in terms of two parameters that are measures of how accurate the individual classifiers are and of the dependence between them. The interplay between these two gives the foundation for understanding the working of random forests. We build on the analysis in Amit and Geman (1997).

Definition 2.1. The margin function for a random forest is

$$mr(\mathbf{X}, Y) = P_\Theta(h(\mathbf{X}, \Theta) = Y) - \max_{j \neq Y} P_\Theta(h(\mathbf{X}, \Theta) = j) \quad (2)$$

and the strength of the effect of classifier $\{h(\mathbf{x}, \Theta)\}$ is

$$s = E_{\mathbf{X}, Y} mr(\mathbf{X}, Y). \quad (3)$$

Averaging $s \geq 0$, Chebyshev's inequality gives

$$PE^* \leq \text{ar}(mr)/s^2 \quad (4)$$

A more revealing expression for the variance of mr is derived in the following: Let

$$\hat{j}(\mathbf{X}, Y) = \arg \max_{j \neq Y} P_\Theta(h(\mathbf{X}, \Theta) = j)$$

so

$$\begin{aligned} mr(\mathbf{X}, Y) &= P_\Theta(h(\mathbf{X}, \Theta) = Y) - P_\Theta(h(\mathbf{X}, \Theta) = \hat{j}(\mathbf{X}, Y)) \\ &= E_\Theta[I(h(\mathbf{X}, \Theta) = Y) - I(h(\mathbf{X}, \Theta) = \hat{j}(\mathbf{X}, Y))]. \end{aligned}$$

Definition 2.2. The raw margin function is

$$rmg(\Theta, \mathbf{X}, Y) = I(h(\mathbf{X}, \Theta) = Y) - I(h(\mathbf{X}, \Theta) = \hat{j}(\mathbf{X}, Y)).$$

Then, $mr(\mathbf{X}, Y)$ is the expectation of $rmg(\Theta, \mathbf{X}, Y)$ with respect to Θ . For an f function f the identity

$$[E_\Theta f(\Theta)]^2 = E_{\Theta, \Theta'} f(\Theta) f(\Theta')$$

holds where Θ, Θ' are independent with the same distribution, implying that

$$mr(\mathbf{X}, Y)^2 = E_{\Theta, \Theta'} rmg(\Theta, \mathbf{X}, Y) rmg(\Theta', \mathbf{X}, Y) \quad (5)$$

Using (5) gives

$$\begin{aligned} \text{ar}(mr) &= E_{\Theta, \Theta'} (\text{cov}_{\mathbf{X}, Y} rmg(\Theta, \mathbf{X}, Y) rmg(\Theta', \mathbf{X}, Y)) \\ &= E_{\Theta, \Theta'} (\rho(\Theta, \Theta') sd(\Theta) sd(\Theta')) \end{aligned} \quad (6)$$

where $\rho(\Theta, \Theta')$ is the correlation between $rmg(\Theta, \mathbf{X}, Y)$ and $rmg(\Theta', \mathbf{X}, Y)$ holding Θ, Θ' fixed and $sd(\Theta)$ is the standard deviation of $rmg(\Theta, \mathbf{X}, Y)$ holding Θ fixed. Then,

$$\begin{aligned} \text{ar}(mr) &= \bar{\rho}(E_\Theta sd(\Theta))^2 \\ &\leq \bar{\rho} E_\Theta \text{ar}(\Theta) \end{aligned} \quad (7)$$

where $\bar{\rho}$ is the mean value of the correlation; that is,

$$\bar{\rho} = E_{\Theta, \Theta'}(\rho(\Theta, \Theta')sd(\Theta)sd(\Theta'))/E_{\Theta, \Theta'}(sd(\Theta)sd(\Theta'))$$

Write

$$\begin{aligned} E_{\Theta} \text{ar}(\Theta) &\leq E_{\Theta}(E_{\mathbf{X}, Y} \text{rmg}(\Theta, \mathbf{X}, Y))^2 - s^2 \\ &\leq 1 - s^2. \end{aligned} \quad (8)$$

Putting (4), (7), and (8) together yield:

Theorem 2.3. *An upper bound for the generalization error is given by*

$$PE^* \leq \bar{\rho}(1 - s^2)/s^2.$$

Although the bound is likely to be loose, it follows the same generalization for random forests as the VC-type bound does for other types of classifiers. It shows that the two ingredients involved in the generalization error for random forests are the strength of the individual classifier in the forest, and the correlation between them in terms of the raw margin function. The $c/2$ ratio is the correlation divided by the square of the strength. Understanding the functioning of random forests, this ratio will be a helpful guide. The smaller it is, the better.

Definition 2.4. The $c/2$ ratio for a random forest is defined as

$$c/s^2 = \bar{\rho}/s^2.$$

There are implications in the two classification. The margin function is

$$mr(\mathbf{X}, Y) = 2P_{\Theta}(h(\mathbf{X}, \Theta) = Y) - 1$$

The requirement that the strength is positive (see (4)) becomes similar to the familiar weak learning condition $E_{\mathbf{X}, Y} P_{\Theta}(h(\mathbf{X}, \Theta) = Y) > .5$. The raw margin function is $2I(h(\mathbf{X}, \Theta) = Y) - 1$ and the correlation $\bar{\rho}$ is between $I(h(\mathbf{X}, \Theta) = Y)$ and $I(h(\mathbf{X}, \Theta') = Y)$. In particular, if the values for Y are taken to be $+1$ and -1 , then

$$\bar{\rho} = E_{\Theta, \Theta'}[\rho(h(\cdot, \Theta), h(\cdot, \Theta'))]$$

so that $\bar{\rho}$ is the correlation between two different members of the forest averaged over the Θ, Θ' distribution.

For more than two classes, the measure of strength defined in (3) depends on the forest as well as the individual classifier since it is the forest that determines $\hat{j}(\mathbf{X}, Y)$. Another approach

is possible. Write

$$\begin{aligned} PE^* &= P_{\mathbf{X}, Y}(P_{\Theta}(h(\mathbf{X}, \Theta) = Y) - \max_{j \neq Y} P_{\Theta}(h(\mathbf{X}, \Theta) = j) < 0) \\ &\leq \sum_j P_{\mathbf{X}, Y}(P_{\Theta}(h(\mathbf{X}, \Theta) = Y) - P_{\Theta}(h(\mathbf{X}, \Theta) = j) < 0). \end{aligned}$$

Define

$$s_j = E_{\mathbf{X}, Y}(P_{\Theta}(h(\mathbf{X}, \Theta) = Y) - P_{\Theta}(h(\mathbf{X}, \Theta) = j))$$

To be the strength of the effect of classifier $\{h(\mathbf{x}, \Theta)\}$ relative to class j . Note that this definition of strength does not depend on the forecaster. Using Chebyshev's inequality, among all $s_j > 0$ lead to

$$PE^* \leq \sum_j \text{ar}(P_{\Theta}(h(\mathbf{X}, \Theta) = Y) - P_{\Theta}(h(\mathbf{X}, \Theta) = j)) s_j^2 \quad (9)$$

and using identity similar to that used in deriving (7), the variance in (9) can be expressed in terms of average correlation. I did not estimate of the quantity in (9) in our empirical study but I think it could be interesting in a multiple classification problem.

3. Using random features

Some random forests reported in the literature have consistently lower generalization error than others. For instance, random plate selection (Dietterich, 1998) does better than bagging. Breiman's introduction of random noise in the output (Breiman, 1998c) also does better. But none of the existing three forests do as well as Adaboost (Freund & Schapire, 1996) or other algorithms that work by adapting reweighting (arcing) of the training set (see Breiman, 1998b; Dietterich, 1998; Bauer & Kohavi, 1999).

To improve accuracy, the random forest injects random noise to minimize the correlation ρ_L while maintaining strength. The forest has died here consisting of random selected inputs or combination of inputs at each node to grow each tree. The resulting forest gives accuracy that compares favorably with Adaboost. This class of procedure has desirable characteristics:

- i) It is accurate and a good as Adaboost and sometimes better.
- ii) It is reliable robust to outliers and noise.
- iii) It is faster than bagging or boosting.
- i) It gives a clear internal estimate of error, strength, correlation and variable importance.
- It is simple and easily parallelized.

Amid and Geman (1997) gave a hallmark tree for handwritten character recognition using random selection from a large number of geometrically derived features to determine the plates at each node. Although many implementation is different and no problem specific, it is a hallmark that popularized the idea.

3.1. Using out-of-bag estimates to monitor error, strength, and correlation

In most experiments with random forests, bagging is used in tandem with random feature selection. Each new training set is drawn with replacement from the original training set. Then a tree is grown on the new training set using random feature selection. The tree grown are not pruned.

There are two reasons for using bagging. They are (i) that the use of bagging seems to enhance accuracy when random features are used. The second is that bagging can be used to give an ongoing estimate of the generalization error (PE*) of the combined ensemble of trees, a useful estimate for the strength and correlation. These estimates are done out-of-bag, which is explained below.

A simple method for constructing a classifier from an α -training set. Given a specific training set T , form bootstrap training set T_k , construct classifier $h(\mathbf{x}, T_k)$ and let the one to form the bagged predictor. For each y, \mathbf{x} in the training set, aggregate the one only over those classifiers for which T_k does not contain y, \mathbf{x} . Call this the out-of-bag classifier. Then the out-of-bag estimate for the generalization error is the error rate of the out-of-bag classifier on the training set.

Tibshirani (1996) and Wolpert and Macready (1996), proposed using out-of-bag estimates as an ingredient in estimating generalization error. Wolpert and Macready worked on regression type problem and proposed a number of methods for estimating the generalization error of bagged predictor. Tibshirani used out-of-bag estimate of variance to estimate generalization error for arbitrary classifier. The idea of error estimate for bagged classifier in Breiman (1996b), gives empirical evidence to show that the out-of-bag estimate is accurate enough as an estimate of the same idea on the training set. Therefore, using the out-of-bag error estimate removes the need for a separate estimate.

In each bootstrap training set, about one-third of the instances are left out. Therefore, the out-of-bag estimates are based on combining only about one-third of the main classifier in the ongoing main combination. Since the error rate decreases as the number of combination increases, the out-of-bag estimate will tend to overestimate the current error rate. To get unbiased out-of-bag estimate, it is necessary to run parallel points where the error is averaged. But unlike cross-validation, where bias is present but is often unknown, the out-of-bag estimate are unbiased.

Strength and correlation can also be estimated using out-of-bag method. This gives internal estimates that are helpful in understanding classifier accuracy and how to improve it. The details are given in Appendix II. Another application is to the measure of variable importance (see Section 10).

4. Random forests using random input selection

The implementation of random forests with random feature selection is formed by electing at random, at each node, a small group of input variables to split on. Given three existing CART methodology to make more efficient and do not prune. Denote this procedure as Forest-R. The size F of the group is fixed. Two values of F were tried. They used only one randomly selected

Table 1. Data sets summary.

| Data set | Train size | Test size | Inputs | Classes |
|---------------|------------|-----------|--------|---------|
| Glass | 214 | | 9 | 6 |
| Breast cancer | 699 | | 9 | 2 |
| Diabetes | 768 | | 8 | 2 |
| Sonar | 208 | | 60 | 2 |
| Vowel | 990 | | 10 | 11 |
| Ionosphere | 351 | | 34 | 2 |
| Vehicle | 846 | | 18 | 4 |
| Soybean | 685 | | 35 | 19 |
| German credit | 1000 | | 24 | 2 |
| Image | 2310 | | 19 | 7 |
| Ecoli | 336 | | 7 | 8 |
| Vote | 435 | | 16 | 2 |
| Liver | 345 | | 6 | 2 |
| Lever | 15000 | 5000 | 16 | 26 |
| Sat-image | 4435 | 2000 | 36 | 6 |
| Zip-code | 7291 | 2007 | 256 | 10 |
| Waferform | 300 | 3000 | 21 | 3 |
| Twonorm | 300 | 3000 | 20 | 2 |
| Threenorm | 300 | 3000 | 20 | 2 |
| Ringnorm | 300 | 3000 | 20 | 2 |

variable, i.e., $F = 1$. The second took F to be the next integer less than $\log_2 M + 1$, where M is the number of inputs.

The experiments used 13 smaller sized data sets from the UCI repository, 3 larger separated into training and test and 4 machine data sets. There were 10 entries because I had added them in post search. Table 1 gives a brief summary.

On each of the 13 smaller sized data sets, the following procedure was used: a random 10% of the data was set aside. On the remaining data, random forests were grown, growing and combining 100 trees, once with $F = 1$, and the second time with $F = \min(\log_2 M + 1)$. The set aside 10% was then predicted on each forest to get a set error for both. These set errors were combined to the lower value of the out-of-bag estimate in the 100 runs. This was repeated 100 times and the set error averaged. The same procedure was followed for the Adaboost runs which are based on combining 50 trees.

The size for 100 trees in random forest and 50 for Adaboost come from the source. The out-of-bag estimates are based on only about a third of the data points. To get reliable estimates I opted for 100 trees. The second consideration is that growing random forests is much faster than growing the tree based on all inputs needed in Adaboost. Growing the 100 trees in random forest was considerably quicker than the 50 trees for Adaboost.

Table 2. Test error (%).

| Data set | Adaboost | Selection | Forests-RF | single imp. | One-tree |
|---------------|----------|-----------|------------|-------------|----------|
| Glas | 22.0 | 20.6 | 21.2 | 36.9 | |
| Breast cancer | 3.2 | 2.9 | 2.7 | 6.3 | |
| Diabetes | 26.6 | 24.2 | 24.3 | 33.1 | |
| Sonar | 15.6 | 15.9 | 18.0 | 31.7 | |
| Vowel | 4.1 | 3.4 | 3.3 | 30.4 | |
| Ionosphere | 6.4 | 7.1 | 7.5 | 12.7 | |
| Vehicle | 23.2 | 25.8 | 26.4 | 33.1 | |
| German credit | 23.5 | 24.4 | 26.2 | 33.3 | |
| Image | 1.6 | 2.1 | 2.7 | 6.4 | |
| Ecoli | 14.8 | 12.8 | 13.0 | 24.5 | |
| Vote | 4.8 | 4.1 | 4.6 | 7.4 | |
| Liver | 30.7 | 25.1 | 24.7 | 40.6 | |
| Lever | 3.4 | 3.5 | 4.7 | 19.8 | |
| Sat-image | 8.8 | 8.6 | 10.5 | 17.2 | |
| Zip-code | 6.2 | 6.3 | 7.8 | 20.6 | |
| Waferform | 17.8 | 17.2 | 17.3 | 34.0 | |
| Tionorm | 4.9 | 3.9 | 3.9 | 24.7 | |
| Threenorm | 18.8 | 17.5 | 17.5 | 38.4 | |
| Ringnorm | 6.9 | 4.9 | 4.9 | 25.7 | |

In the run on the larger data sets, the random forest results for the two data sets were based on combining 100 trees; the ip-code procedure combined 200. For Adaboost, 50 trees were combined for the three data sets and 100 for ip-code. The synthetic data was described in Breiman (1996) and also used in Schapire et al. (1997). There were 50 runs. In each run, a new training set of size 300 and test set of size 3000 were generated. In random forests 100 trees were combined in each run, 50 in Adaboost. The results of the runs are given in Table 2.

The second column are the results obtained from the two group-wise mean of leave-one-out-of-bag error. The third column is the average error using just one random feature to grow the tree. The fourth column contains the out-of-bag estimate of the generalization error of the individual tree in the forest computed for the best splitting (single or selection). This estimate is computed by fitting the left-over instance at each node in each tree grown and averaging the results over all trees in the forest.

The error rate using random impurity election compares favorably with Adaboost. The comparison might be even more favorable if the search is over more alternatives of F instead of the pre-specified. By the procedure is no longer unique in the choice of F . The average absolute difference between the error rate using $F = 1$ and the higher value of F is less than 1%. The difference is more pronounced on the three large data sets.

Table 3. Test error (%).

| Data set | Adaboost | Forest-RC | | |
|---------------|----------|-----------|-------|----------|
| | | Selection | Trees | One tree |
| Glass | 22.0 | 24.4 | 23.5 | 42.4 |
| Breast cancer | 3.2 | 3.1 | 2.9 | 5.8 |
| Diabetes | 26.6 | 23.0 | 23.1 | 32.1 |
| Sonar | 15.6 | 13.6 | 13.8 | 31.7 |
| Vowel | 4.1 | 3.3 | 3.3 | 30.4 |
| Ionosphere | 6.4 | 5.5 | 5.7 | 14.2 |
| Vehicle | 23.2 | 23.1 | 22.8 | 39.1 |
| German credit | 23.5 | 22.8 | 23.8 | 32.6 |
| Image | 1.6 | 1.6 | 1.8 | 6.0 |
| Ecoli | 14.8 | 12.9 | 12.4 | 25.3 |
| Vote | 4.8 | 4.1 | 4.0 | 8.6 |
| Liner | 30.7 | 27.3 | 27.2 | 40.3 |
| Letter | 3.4 | 3.4 | 4.1 | 23.8 |
| Sat-image | 8.8 | 9.1 | 10.2 | 17.3 |
| Zip-code | 6.2 | 6.2 | 7.2 | 22.7 |
| Waferform | 17.8 | 16.0 | 16.1 | 33.2 |
| Twonorm | 4.9 | 3.8 | 3.9 | 20.9 |
| Threennorm | 18.8 | 16.8 | 16.9 | 34.8 |
| Ringnorm | 6.9 | 4.8 | 4.6 | 24.6 |

8.5%, on the lower data to 3.0%, but the zip-code test error did not decrease. Acting on an informed hunch, I tried Forest-RC with $F = 25$. The zip-code test error dropped to 5.8%. The same technique achieved so far achieved on the other three data sets three ensemble.

5.1. Categorical variables

Some or all of the input variable may be categorical and it is very hard to define additive combination of variable, we need to define how categorical will be treated or the can be combined with numerical variable. My approach is that each time a categorical variable is selected to split on a node, to select a random subset of the categories of the variable, and define a binary variable that is one when the categorical value of the variable is in the subset and zero otherwise.

Since a categorical variable with I values can be coded into $I - 1$ dummy 0-1 variables, we make the variable $I - 1$ time a probabilistic numeric variable to be selected in node splitting. When many of the variable are categorical, taking a low value of F results in low correlation, but also low strength. F may be increased to about two-thirds time in $(\log_2 M + 1)$ to get enough strength to provide good predictive accuracy.

For instance, on the DNA data set having 60 features and 2,000 examples in the training set and 1,186 in the test set, using Forecast-RIL with $F = 20$ gave a 1% error rate of 3.6% (4.2% for Adaboost). The soybean data has 685 examples, 35 variables, 19 classes, and 15 categorical variables. Using Forecast-RIL with $F = 12$ gives a 1% error of 5.3% (5.8% for Adaboost). Using Forecast-RCU with combination of 3 and $F = 8$ gives an error of 5.5%.

One advantage of this approach is that it generalizes to domains with categorical variables. In the two-class problem, this can be avoided by using the decision proposed in Breiman et al. (1985), which reduces the search for the best categorical splits to an $O(I)$ computation. For more classes, the search for the best categorical splits is an $O(2^{I-1})$ computation. In the random forest implementation, the computation for an categorical variable involves only the election of a random subset of the categories.

6. Empirical results on strength and correlation

The purpose of this section is to look at the effect of strength and correlation on the generalization error. Another aspect that we wanted to get more understanding of was the lack of sensitivity in the generalization error to the group F . To conduct an empirical study of the effect of strength and correlation in a series of data sets, a 10-fold-cross-validation of the strength and correlation, as described in Section 3.1, were used.

We begin by running Forecast-RIL on the sonar data (60 inputs, 208 examples) using from 1 to 50 inputs. In each iteration, 10% of the data was split off as a test set. Then F , the number of random inputs selected at each node, varied from 1 to 50. For each value of F , 100 trees were grown to form a random forest and the terminal value of the error, strength, correlation, etc. recorded. Eight iterations were done, each time removing a random 10% of the data for a new set, and all results averaged over 80 repetitions. Altogether, 400,000 trees were grown in this experiment.

The top graph of Figure 1 plots the value of strength and correlation vs. F . The result is fascinating. For about $F = 4$, the strength remains constant; adding more inputs does not help. Below the correlation continues to increase. The second graph plots the test set error and the 10-fold-cross-validation of the generalization error again vs. F . The 10-fold-cross-validation are more volatile. Both however show the same behavior: a small drop from $F = 1$ to $F = 4$, and then a general gradual increase. This increase in error starts at the beginning of the constant region for the strength.

Figure 2 has plots for similar runs on the breast cancer data set where feature combining of random combination of three inputs are used. The number of features varied from 1 to 25. Again, the correlation has a local rise, while the strength varies in all constant, so that the minimum error is at $F = 1$. The surprise is the exceptionality of the relationship between strength and correlation. Since the correlation are local but gradually increasing, the local error occurs when only a few inputs or features are used.

Since the larger data set seemed to have a different behavior than the smaller, we ran a similar experiment on the aethane data set. The number of features, each combining of a random sum of two inputs, varied from 1 to 25, and for each, 100 classifiers were combined. The results are shown in Figure 3. The results differ from those on the smaller



Figure 1. Effect of number of input variables on random data.

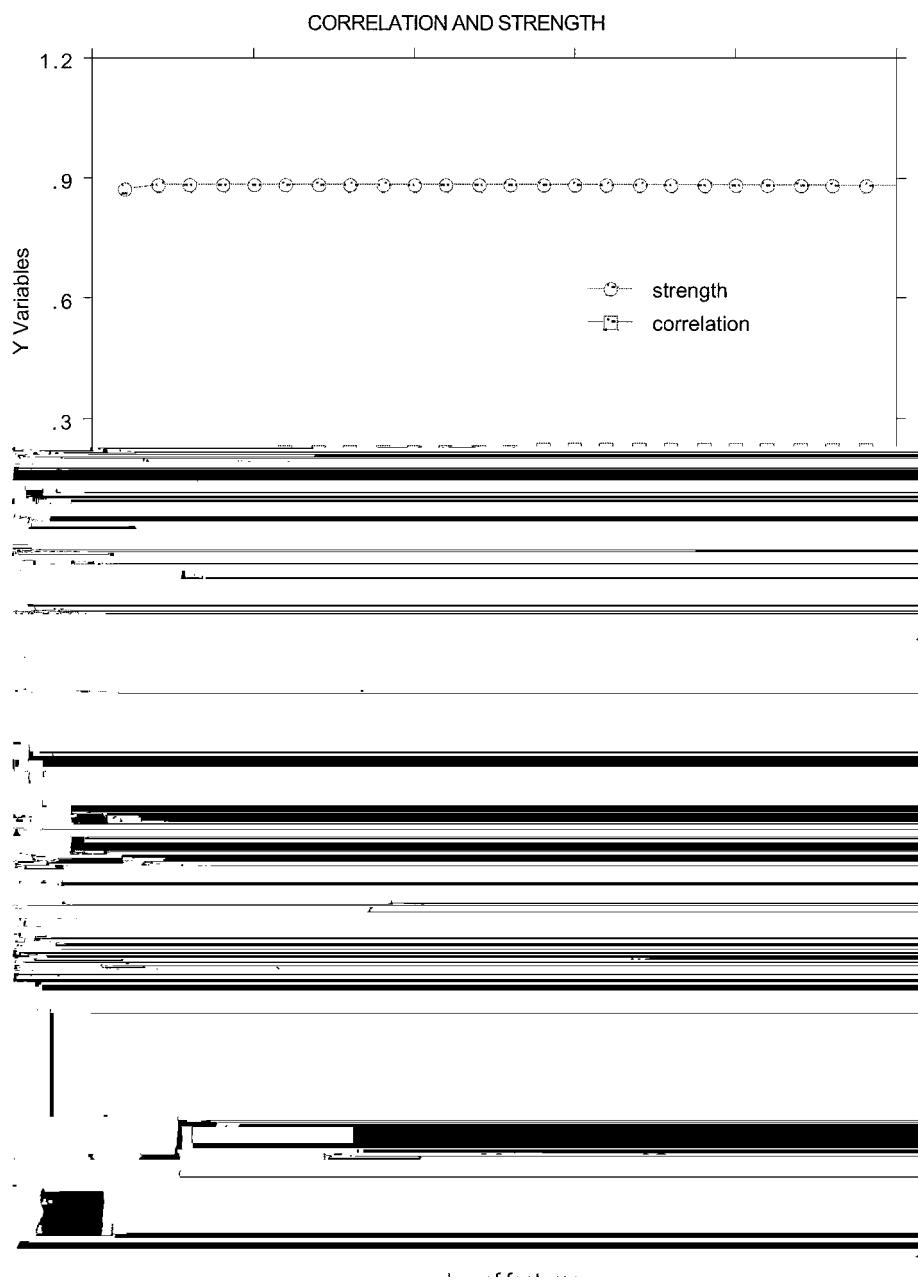


Figure 2. Effect of the number of features on the breast data set.

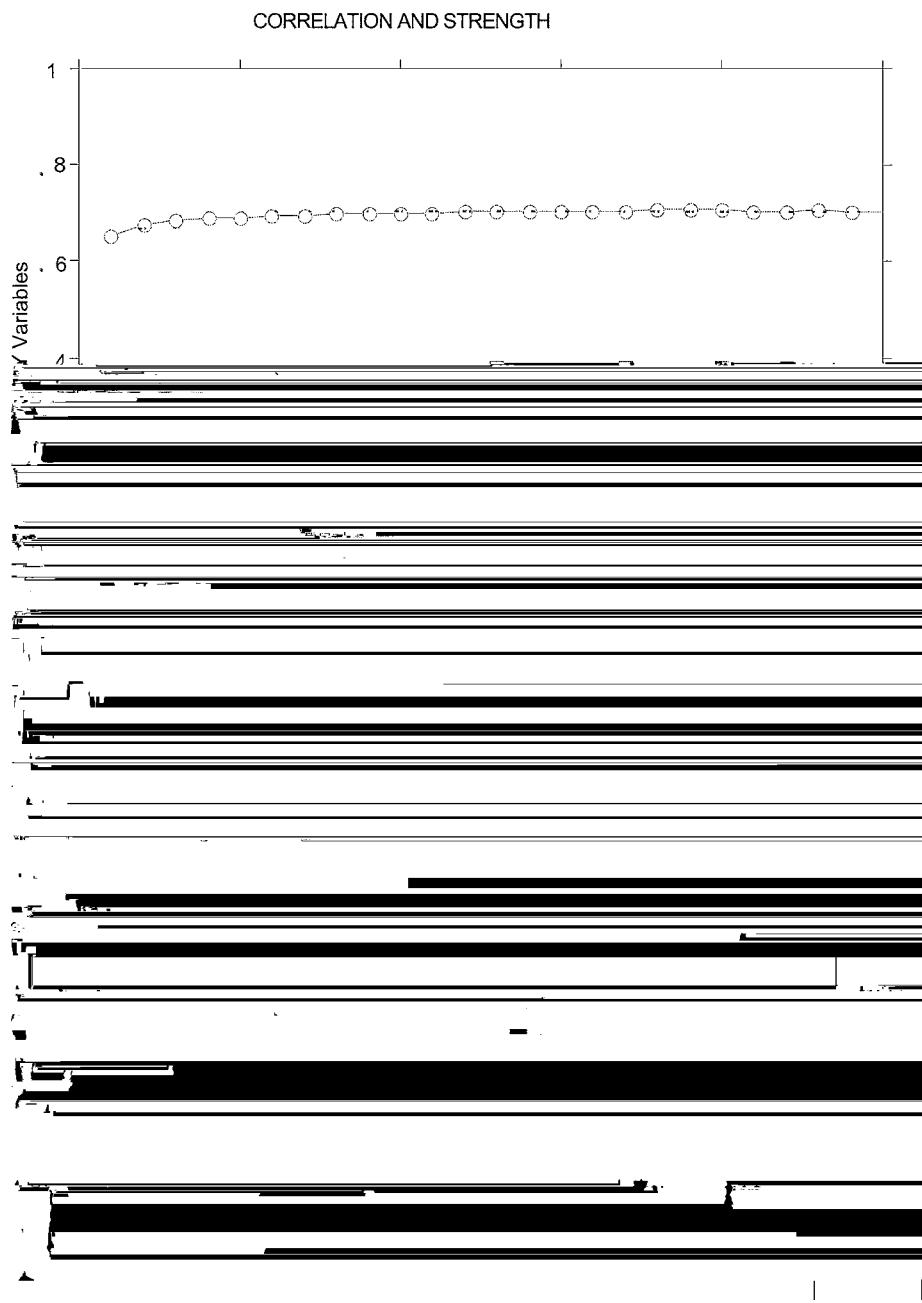


Figure 3. Effect of number of features on aellipte data.

data set. Both the correlation and strength have a small bias toward increase. The error rate has a slight decrease. We conjecture that with larger and more complex data sets, the strength continues to increase longer before it plateaus.

Our results indicate that better (lower generalization error) random forests have lower correlation between classifier and higher strength. The random forest in tree construction has no aim for low correlation ρ_L while maintaining reasonable strength. This conclusion has been hinted at in previous work. Dietterich (1998) has measured performance of different ensemble methods and noted that more accurate ensemble has larger dimension. Freund (personal communication) believes that one reason why Adaboost works so well is that at each step it tries to decompose the feature space into smaller from the current one. Amit et al. (1999) give an analysis showing that the Adaboost algorithm is aimed at keeping the covariance between classifiers small.

7. Conjecture: Adaboost is a random forest

Various classifiers can be modified to use both a training set and a set of weights on the training set. Consider the following random forest: a large collection of K different sets of non-negative m-one weights on the training set is defined. Denote the weights by $\mathbf{w}(1), \mathbf{w}(2), \dots, \mathbf{w}(K)$. Corresponding to the weights are probabilities $p(1), p(2), \dots, p(K)$, whose sum is one. Draw from the integer $1, \dots, K$ according to the probabilities. The outcome is Θ . If $\Theta = k$ draw the classifier $h(\mathbf{x}, \Theta)$ using the training set with weight $\mathbf{w}(k)$.

In the original version, Adaboost (Freund & Schapire, 1996) is a deterministic algorithm that selects the weight on the training set for input \mathbf{x} by the nearest classifier based on the minimum classification distance to the preexisting classifier. In our experiment, random forests were produced as follows: Adaboost runs 75 times on a data set producing sets of non-negative m-one weights $\mathbf{w}(1), \mathbf{w}(2), \dots, \mathbf{w}(50)$ (here $r = 25$ were discarded). The probability for the k th classifier of weight i is proportional to $Q(\mathbf{w}_k) = \log[(1 - \text{error}(k))/\text{error}(k)]$, where $\text{error}(k)$ is the $\mathbf{w}(k)$ -weighted training set error of the k th classifier. Then the forest is run 250 times.

This was repeated 100 times on a fixed data set, each time learning only 10% of the set and then averaging the error rate. On each data set, the Adaboost error rate was closer to the random forest error rate. A typical result is on the Wisconsin Breast Cancer data set where Adaboost produced an average of 2.91% error and the random forest produced 2.94%.

In the Adaboost algorithm, $\mathbf{w}(k+1) = \phi(\mathbf{w}(k))$, where ϕ is a function determined by the base classifier. Denote the k th classifier by $h(\mathbf{x}, \mathbf{w}_k)$. The outcome of the k th classifier is weighted by $Q(\mathbf{w}_k)$ or the normalized one for classifier j if \mathbf{x} is equal to

$$\sum_k I(h(\mathbf{x}, \mathbf{w}_k) = j) Q(\mathbf{w}_k) / \sum_k Q(\mathbf{w}_k). \quad (10)$$

For an function f defined on the weight space, define the operator $\mathbf{T}f(\mathbf{w}) = f(\phi(\mathbf{w}))$. We conjecture that \mathbf{T} is ergodic with invariant measure $\pi(d\mathbf{w})$. Then (10) will converge to $E_{Q\pi}[I(h(\mathbf{x}, \mathbf{w}) = j)]$, where the distribution $Q\pi(d\mathbf{w}) = Q(\mathbf{w})\pi(d\mathbf{w})/\int Q(\mathbf{v})\pi(d\mathbf{v})$. If this conjecture is true, then Adaboost is equivalent to a random forest where the weights on the training set are selected at random from the distribution $Q\pi$.

If γ is high, then all old or plain Adaboost does not overfit, as more trees are added to the ensemble, and some performance factors have been proposed. There is some experimental evidence that Adaboost may overfit if γ is too large (Grove & Schuurman, 1998), but the experiments were done using a simple base classifier and many more classifiers than trees, as the base classifier. The experience running Adaboost on up to 1,000,000 trees on a number of datasets indicates that the error converges to an asymptotic value.

The truth of this conjecture does not solve the problem of how Adaboost selects the favorable distribution on the weight space that it does. Note that the distribution of the weights will depend on the training set. In the final random forest, the distribution of the random vector does not depend on the training set.

8. The effects of output noise

Dietterich (1998) hypothesized that when a fraction of the output label in the training set are randomly altered, the accuracy of Adaboost degrades, while bagging and random playout are more immune to the noise. Since some noise in the output is often present, robustness with respect to noise is a desirable property. Following Dietterich (1998), the following experiments were done, which changed about one in ten class label (injecting 5% noise).

For each dataset in the experiment, 10% of random instances are left off at a time. Two runs are made on the remaining training set. The runs interchange the training set and the test set. The second run is on a noisy version of the training set. The noise level is given by changing, at random, 5% of the class label into an alternate class label chosen uniformly from the other label.

This is repeated 50 times using Adaboost (deterministic version), Forest-RI and Forest-RC. These three runs are averaged over the 50 repetition and the percent increased error due to noise is computed. In both random forests, the number of features is fixed at 10. The noise error in the Section 5 and 6 experiments. Because of the length of the run, only the 9 smallest datasets are used. Table 4 gives the increase in error rate due to the noise.

Table 4. Increase in error rate due to noise (%).

| Dataset | Adaboost | Forest-RI | Forest-RC |
|---------------|----------|-----------|-----------|
| Gla | 1.6 | .4 | -.4 |
| Breast cancer | 43.2 | 1.8 | 11.1 |
| Diabetes | 6.8 | 1.7 | 2.8 |
| Sonar | 15.1 | -6.6 | 4.2 |
| Ionosphere | 27.7 | 3.8 | 5.7 |
| Soybean | 26.9 | 3.2 | 8.5 |
| Ecoli | 7.5 | 7.9 | 7.8 |
| Vote | 48.9 | 6.3 | 4.6 |
| Liver | 10.3 | -.2 | 4.8 |

Adaboost deteriorate markedly with 5% noise while the random forest procedure generally has small change. The effect on Adaboost is critical depending on the noise level. The Adaboost algorithm iteratively increases the weight on the instance most recently misclassified. In practice having incorrect class label will prevent it from being misclassified. Then, Adaboost will concentrate increasing weight on the noisy instance and become warped. The random forest procedure does not concentrate weight on an individual of the instance and the noise effect is maller.

9. Data with many weak inputs

Data sets with many weak inputs are becoming more common, i.e. in medical diagnosis, document retrieval, etc. The common characteristic is no single input or small group of inputs can distinguishing between the classes. This type of data is difficult for the classifier, neural network and tree.

To see if there is a possibility that Forest-Random method can work, the following 10 classes, 1,000 binary inputs data was generated: (rnd is a uniform random number, elected and each time it appears)

```

do j=1,10
do k=1,1000
p(j,k)=.2*rnd+.01
end do
end do

do j=1,10
do i=1, min(400*rnd) !min=nearest integer
k=min(1000*rnd)
p(j,k)=p(j,k)+.4*rnd
end do
end do

do n=1,N
j=min(10*rnd)
do m=1,1000
if (rnd<p(j,m)) then
(m,n)=1
else
(m,n)=0
end if
(n)=j ! (n) is the class label of the next example
end do
end do

```

This code generates a set of probabilities $\{p(j, m)\}$ where j is the class label and m is the input number. Then the inputs for a class example are a string of M binary variables with the m th variable having probability $p(j, m)$ of being one.

For the training set, $N = 1,000$. A 4,000 example set was also considered using the same $\{p(j, k)\}$. Examination of the code shows that each class has higher underlying probability at certain location. By the total of all classes of the location is about 2,000, where it is significantly overlapping. Among one knows all of the $\{p(j, k)\}$, the base error rate for the particular $\{p(j, m)\}$ computed is around 1.0%.

Since the inputs are independent of each other, the Naïve Bayes classifier which estimates the $\{p(j, k)\}$ from the training data is probably optimal and has an error rate of 6.2%. This is not an endorsement of Naïve Bayes, since it would be easier to create a dependence between the inputs which would increase the Naïve Bayes error rate. In addition, this example because the Bayes error rate and the Naïve Bayes error rate are equal to compute.

I started with a run of Forecast-RI with $F = 1$. It converged after 1000 and by 2,500 iteration, when it had stopped, it had still not converged. The total error was 10.7%. The strength was .069 and the correlation .012 with a $c/2$ ratio of 2.5. Even though the strength was low, the almost zero correlation means that they were adding small increments of accuracy as the iteration proceeded.

Clearly, what was needed was an increase in strength while keeping the correlation low. Forecast-RI ran again using $F = \min(\log_2 M + 1) = 10$. The results were encouraging. It converged after 2,000 iteration. The total error was 3.0%. The strength was .22, the correlation .045 and $c/2 = .91$. Going with the trend, Forecast-RI ran with $F = 25$ and stopped after 2,000 iteration. The total error was 2.8%. Strength was .28, correlation .065 and $c/2 = .83$.

It is interesting that Forecast-RI could produce error rates no far above the Bayes error rate. The individual classifier are weak. For $F = 1$, the average tree error rate is 80%; for $F = 10$, it is 65%; and for $F = 25$, it is 60%. Forecast seems to have the ability to work with weaker classifiers as their correlation is low. A comparison using Adaboost was tried, but I can't get Adaboost to run on this data because the base classifier are too weak.

10. Exploring the random forest mechanism

A forecast of three is impenetrable as far as interpretation of its mechanism goes. In some applications, analysis of medical experiments for example, it is critical to understand the interaction of variables that is providing the predictive accuracy. A variation problem is made by using internal out-of-bag estimate, and estimation by rerunning only selected variables.

Suppose there are M inputs variables. After each tree is constructed, the value of the m th variable in the out-of-bag example are randomly permuted and the out-of-bag data is run down the corresponding tree. The classification given for each x_n that is out-of-bag is used. This is repeated for $m = 1, 2, \dots, M$. At the end of the run, the probability of out-of-bag classifier for x_n with the m th variable noise is compared with the classifier label of x_n to give a misclassification rate.

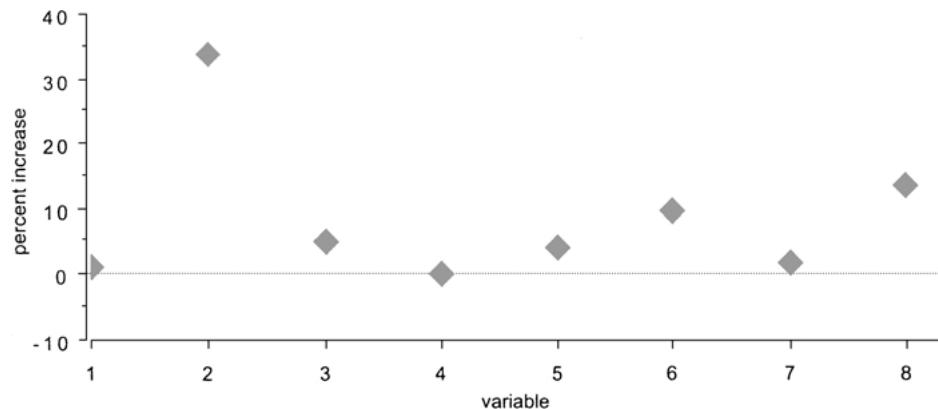


Figure 4. Measure of variable importance, diabetes data.

The top is the percent increase in misclassification rate compared to the out-of-bag rate (with all variable included). We get the estimate by fitting a single run of a forest with 1,000 trees and no replacement. The procedure is illustrated below.

In the diabetes data set, fitting only one variable with $F = 1$, the relative error due to the noise of variable is given in figure 4.

The second variable appears to be far more important followed by variable 8 and variable 6. Running the random forest in 100 repetitions using only variable 2 and leading to 10% each time to mean relative error gave an error of 29.7%, compared with 23.1% using all variables. But when variable 8 is added, the error falls only to 29.4%. When variable 6 is added to variable 2, the error falls to 26.4%.

The reason why variable 6 seems important, even though once variable 2 is entered it is a characteristic of how dependent variable affects prediction error in random forests. Suppose there are two variables x_1 and x_2 which are identical and carry significant predictive information. Because each gets picked up about the same frequency in a random forest, noting each separately will result in the same increase in error rate. But once x_1 is entered as a predictive variable, fitting x_2 in addition will not produce an decrease in error rate. In the diabetes data set, the 8th variable carries some of the same information as the second. So it does not add predictive accuracy when combined with the second.

The relative magnitude of relative error rates are fairly stable with respect to the importance feature. The experiments above also repeated using combination of three inputs with $F = 2$. The results are in figure 5.

Another interesting example is the voting data. This has 435 examples corresponding to 435 Congressmen and 16 variables reflecting their voting record on 16 issues. The class variable is Republican or Democrat. To see which is the most important, we again ran the noisy variable program generating 1,000 trees. The lowest error rate on the original data was given using single input with $F = 5$, so the parameters were fixed in the run. The results are in figure 6.

Variable 4 (and only) the error triple if variable 4 is removed. We reran this data set using only variable 4. The relative error is 4.3%, about the same as if all variables were used. The

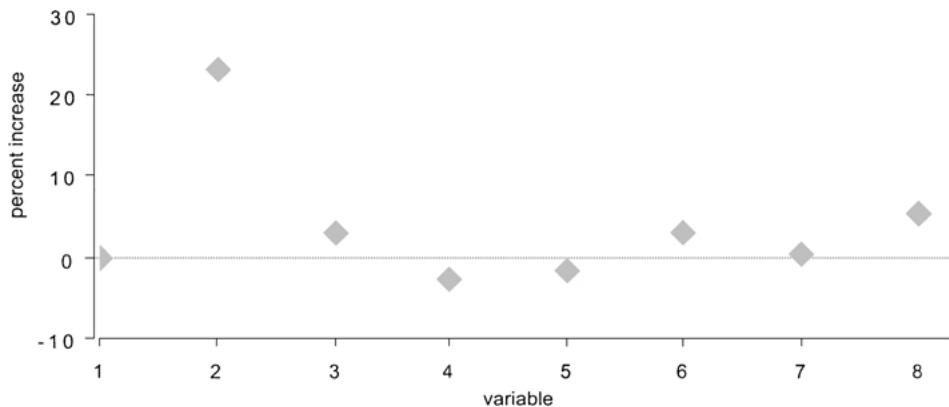


Figure 5. Measure of variable importance-diabetes data.

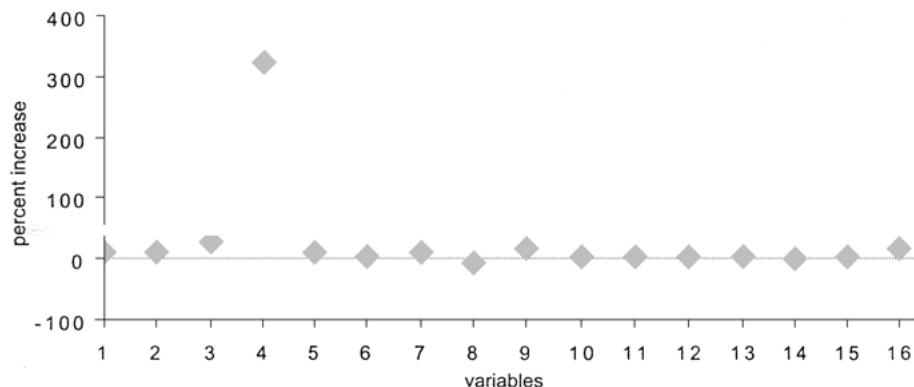


Figure 6. Measure of variable importance-oe data.

one on 4 separate Rep blican from Democrat almost as well as the one on 4 combined with the one on all other 15 i.e.

The approach given in this section is only a beginning. More research will be necessary to understand how to give a more complete picture.

11. Random forests for regression

Random forests for regression are formed by growing trees depending on a random vector Θ such that the predictor $h(\mathbf{x}, \Theta)$ take on numerical values as opposed to class label. The outputs are numerical and we assume that the training set is independently drawn from the distribution of the random vector Y, \mathbf{X} . The mean-squared generalization error for an numerical predictor $h(\mathbf{x})$ is

$$E_{\mathbf{X}, Y} (Y - h(\mathbf{X}))^2 \quad (11)$$

The random forest predictor is formed by taking the average over k of the trees $\{h(\mathbf{x}, \Theta_k)\}$. Similarly, the classification error follows hold:

Theorem 11.1. *As the number of trees in the forest goes to infinity, almost surely,*

$$E_{\mathbf{X}, Y}(Y - av_k h(\mathbf{X}, \Theta_k))^2 \rightarrow E_{\mathbf{X}, Y}(Y - E_{\Theta} h(\mathbf{X}, \Theta))^2. \quad (12)$$

Proof: See Appendix I. □

Denote the right hand side of (12) as $PE^*(\text{forest})$, the generalization error of the forest. Define the average generalization error of a tree as:

$$PE^*(\text{tree}) = E_{\Theta} E_{\mathbf{X}, Y}(Y - h(\mathbf{X}, \Theta))^2$$

Theorem 11.2. *Assume that for all Θ , $EY = E_{\mathbf{X}} h(\mathbf{X}, \Theta)$. Then*

$$PE^*(\text{forest}) \leq \bar{\rho} PE^*(\text{tree})$$

where $\bar{\rho}$ is the weighted correlation between the residuals $Y - h(\mathbf{X}, \Theta)$ and $Y - h(\mathbf{X}, \Theta')$ where Θ, Θ' are independent.

Proof:

$$\begin{aligned} PE^*(\text{forest}) &= E_{\mathbf{X}, Y}[E_{\Theta}(Y - h(\mathbf{X}, \Theta))^2] \\ &= E_{\Theta} E_{\Theta'} E_{\mathbf{X}, Y}(Y - h(\mathbf{X}, \Theta))(Y - h(\mathbf{X}, \Theta')) \end{aligned} \quad (13)$$

The term on the right in (13) is a covariance and can be written as:

$$E_{\Theta} E_{\Theta'} (\rho(\Theta, \Theta') sd(\Theta) sd(\Theta'))$$

where $sd(\Theta) = \sqrt{E_{\mathbf{X}, Y}(Y - h(\mathbf{X}, \Theta))^2}$. Define the weighted correlation as:

$$\bar{\rho} = E_{\Theta} E_{\Theta'} (\rho(\Theta, \Theta') sd(\Theta) sd(\Theta')) / (E_{\Theta} sd(\Theta))^2 \quad (14)$$

Then

$$PE^*(\text{forest}) = \bar{\rho} (E_{\Theta} sd(\Theta))^2 \leq \bar{\rho} PE^*(\text{tree}). \quad \square$$

Theorem (11.2) pinpoints the requirement for accurate regression forest: low correlation between individual and low error tree. The random forest decreases the average error of the tree employed by the factor $\bar{\rho}$. The random forest employed no aim at low correlation.

12. Empirical results in regression

In regression forests we random feature selection on top of bagging. Therefore, we can use the monitoring provided by out-of-bag estimation to estimate of $\text{PE}^*(\text{forest})$, $\text{PE}^*(\text{tree})$ and $\bar{\rho}$. These are derived similarly to the estimate in classification. Through this, features formed by a random linear combination of inputs are used. We comment later on how many of these features to use to determine the splits at each node. The more features used, the lower $\text{PE}^*(\text{tree})$ but the higher $\bar{\rho}$. In our empirical study the following data sets are used, see Table 5.

Of the data sets, the Boston Housing, Abalone and Serpent are available at the UCI repository. The Robot Arm data was provided by Michael Jordan. The last three datasets are synthetic. They originated in Friedman (1991) and are also described in Breiman (1998b). These are the same data sets used to compare adaptive bagging to bagging (see Breiman, 1999), except that one synthetic dataset (Peak20), which includes found anomalies both better and worse than the original, is eliminated.

The first three datasets listed are moderate in size and their error can be estimated by leaving out a random 10% of the instances, running on the remaining 90% and using the left-out 10% as a test set. This was repeated 100 times and the average error averaged. The abalone dataset is larger with 4,177 instances and 8 input variables. It originally came with 25% of the instances as a separate set. We ran this dataset by leaving out a random selected 25% of the instances to evaluate it, repeated this 10 times and averaged.

Table 6 gives the average mean-squared error for bagging, adaptive bagging and the random forest. The experiments all run using 25 features, each a random linear combination of two randomly selected inputs, no split at each node, each feature a random combination of two inputs. All runs with all datasets, combined 100 trees. In all datasets, the rules don't split if the node size is < 5 units enforced.

An interesting difference between regression and classification is that the correlation increases quite quickly as the number of features used increases. The major effect is the decrease in $\text{PE}^*(\text{tree})$. Therefore, a relatively large number of features are required to reduce $\text{PE}^*(\text{tree})$ and get near optimal average error.

Table 5. Data summary.

| Data set | Nr. inputs | #Training | #Test |
|----------------|------------|-----------|-------|
| Boston Housing | 12 | 506 | 10% |
| Ozone | 8 | 330 | 10% |
| Serpent | 4 | 167 | 10% |
| Abalone | 8 | 4177 | 25% |
| Robot Arm | 12 | 15,000 | 5000 |
| Friedman#1 | 10 | 200 | 2000 |
| Friedman#2 | 4 | 200 | 2000 |
| Friedman#3 | 4 | 200 | 2000 |

Table 6. Mean-quartered error.

| Data set | Bagging | Adapt. bag | Forest |
|-----------------------------|---------|------------|--------|
| Boston Housing | 11.4 | 9.7 | 10.2 |
| Ozone | 17.8 | 17.8 | 16.3 |
| Servo $\times 10 - 2$ | 24.5 | 25.1 | 24.6 |
| Abalone | 4.9 | 4.9 | 4.6 |
| Robot Arm $\times 10 - 2$ | 4.7 | 2.8 | 4.2 |
| Friedman #1 | 6.3 | 4.1 | 5.7 |
| Friedman #2 $\times 10 + 3$ | 21.5 | 21.5 | 19.6 |
| Friedman #3 $\times 10 - 3$ | 24.8 | 24.8 | 21.6 |

The results in Table 6 are mixed. Random forest-random feature is always better than bagging. In data sets for which adaptive bagging gives sharp decrease in error, the decrease is produced by forests not by pronouncing. In data sets in which adaptive bagging gives no improvement over bagging, forests produce improvement.

For the same number of inputs combined, over a wide range, the error does not change much with the number of features. If the number is too small, PE*(tree) becomes too large and the error goes up. If the number is too large, the correlation goes up and the error again increases. The intermediate range is usually large. In this range, as the number of features goes up, the correlation increases, but PE*(tree) compensates decreasing.

Table 7 gives the tree error, the out-of-bag error estimate, and the OB estimate for PE*(tree) and the correlation.

As expected, the OB Error estimate are consistently high. It is low in the robot arm data, but I believe that this is an artifact caused by separate training and testing, where there is a marginal higher error rate than the training error.

In an experiment, I turned off the bagging and replaced it by randomizing output (Breiman, 1998b). In this procedure, mean-zero Gaussian noise is added to each of the outputs. The standard deviation of the noise is equal to the standard deviation of the

Table 7. Error and OB estimate.

| Data Set | Tree error | OB error | PE*(tree) | Cor. |
|-----------------------------|------------|----------|-----------|------|
| Boston Housing | 10.2 | 11.6 | 26.3 | .45 |
| Ozone | 16.3 | 17.6 | 32.5 | .55 |
| Servo $\times 10 - 2$ | 24.6 | 27.9 | 56.4 | .56 |
| Abalone | 4.6 | 4.6 | 8.3 | .56 |
| Robot Arm $\times 10 - 2$ | 4.2 | 3.7 | 9.1 | .41 |
| Friedman #1 | 5.7 | 6.3 | 15.3 | .41 |
| Friedman #2 $\times 10 + 3$ | 19.6 | 20.4 | 40.7 | .51 |
| Friedman #3 $\times 10 - 3$ | 21.6 | 22.9 | 48.3 | .49 |

Table 8. Mean-quared error.

| Data set | With bagging | With Noise |
|-----------------------------|--------------|------------|
| Boston Housing | 10.2 | 9.1 |
| Ozone | 17.8 | 16.3 |
| Servo $\times 10 - 2$ | 24.6 | 23.2 |
| Abalone | 4.6 | 4.7 |
| Robot Arm $\times 10 - 2$ | 4.2 | 3.9 |
| Friedman #1 | 5.7 | 5.1 |
| Friedman #2 $\times 10 + 3$ | 19.6 | 20.4 |
| Friedman #3 $\times 10 - 3$ | 21.6 | 19.8 |

oip). Similar to the bagging experiments, three configurations were done using 25 features, each a random linear combination of two random selected inputs, to split each node. The results are given in Table 8.

The error rates on these two datasets are the lowest to date. Overall, adding oip noise to work with random feature selection better than bagging. This illustrates the effectiveness of the random forest using various combination of randomness can be added to see what works best.

13. Remarks and conclusions

Random forests are an effective tool in prediction. Because of the Law of Large Number they do not overfit. Injecting the right kind of randomness make them accurate classifier and regressor. Furthermore, the framework in term of strength of the individual predictor and their correlation give insight into the ability of the random forest to predict. Using out-of-bag estimation make concrete the otherwise theoretical idea of strength and correlation.

For awhile, the conventional thinking was that forests could not compete with arcing type algorithm in term of accuracy. Over time, belief, belief led to improving question. Boosting and arcing algorithm have the ability to reduce bias and variance (Schapire et al., 1998). The adaptive bagging algorithm in regression (Breiman, 1999) is designed to reduce bias and operate effectively in classification and regression. But, like arcing, it also changes the training and a little progress.

Forests give rise to competition with boosting and adaptive bagging, and do not progress effectively change the training set. Their accuracy indicates that the accuracy is reduced. The mechanism for this is not obvious. Random forests may also be based on a Bayesian procedure. Although I do not know if this is a fine line of exploration, if it could explain the bias reduction, I might become more of a Bayesian.

Random input and random feature produce good results in classification, less so in regression. The only use of randomness is in hybridized bagging and random feature. It may well be that other types of injected randomness give better results. For instance, one of the referee has suggested the use of random Boolean combination of features.

An almost obvious question is whether gain in accuracy can be gotten by combining random features with boosting. For the larger data sets, it seems that gains can't be very large. For one run, the error rate goes down 5.1% on the zip-code data, 2.2% on the letter data and 7.9% on the aellie data. The improvements are on the smaller data sets. More work is needed on this; but it does suggest that different injection of randomness can produce better results.

A recent paper (Breiman, 2000) has shown that in distribution pace for the classification problem, random forests are equivalent to a kernel acting on the margin. Arguments are given that random forests (large correlation) enforce the smoothness of the kernel while strength enhances a desirable feature called predictability. Hopefully, this shed light on the dual role of correlation and strength. The theoretical framework given by Kleinberg (2000) for stochastic discrimination may also help understanding.

Appendix I: Almost sure convergence

Proof of theorem 1.2: It suffices to show that there is a constant of probability zero C on the equivalence pace $\Theta_1, \Theta_2, \dots$ such that outside of C , for all \mathbf{x} ,

$$\frac{1}{N} \sum_{n=1}^N I(h(\Theta_n, \mathbf{x}) = j) \rightarrow P_\Theta(h(\Theta, \mathbf{x}) = j).$$

For a fixed training set and fixed Θ , the set of all \mathbf{x} such that $h(\Theta, \mathbf{x}) = j$ is a union of hyper-rectangle. For all $h(\Theta, \mathbf{x})$ there is only a finite number K of such union of hyper-rectangle, denoted by S_1, \dots, S_K . Define $\varphi(\Theta) = k$ if $\{\mathbf{x} : h(\Theta, \mathbf{x}) = j\} = S_k$. Let N_k be the number of times that $\varphi(\Theta_n) = k$ in the trial. Then

$$\frac{1}{N} \sum_{n=1}^N I(h(\Theta_n, \mathbf{x}) = j) = \frac{1}{N} \sum_k N_k I(\mathbf{x} \in S_k)$$

By the Law of Large Number,

$$N_k = \frac{1}{N} \sum_{n=1}^N I(\varphi(\Theta_n) = k)$$

converges a.s. to $P_\Theta(\varphi(\Theta) = k)$. Taking union of all the sets on which convergence does not occur for some value of k gives a set C of zero probability such that outside of C ,

$$\frac{1}{N} \sum_{n=1}^N I(h(\Theta_n, \mathbf{x}) = j) \rightarrow \sum_k P_\Theta(\varphi(\Theta) = k) I(\mathbf{x} \in S_k).$$

The right hand side is $P_\Theta(h(\Theta, \mathbf{x}) = j)$. □

Proof of theorem 9.1: There are at most c of h per-rectangle R_1, \dots, R_K , such that if \bar{y}_k is the average of the training error for all training inputs in R_k when $h(\Theta, \mathbf{x})$ has one of the values $I(\mathbf{x} \in S_k)\bar{y}_k$. The rest of the proof parallel that of Theorem 1.2. \square

Appendix II: Out-of bag estimates for strength and correlation

At the end of a combination run, let

$$Q(\mathbf{x}, j) = \sum_k I(h(\mathbf{x}, \Theta_k) = j; (y, \mathbf{x}) \notin T_{k,B}) / \sum_k I((y, \mathbf{x}) \notin T_{k,B}).$$

Then, $Q(\mathbf{x}, j)$ is the out-of-bag proportion of one case at \mathbf{x} for class j , and is an estimate for $P_\Theta(h(\mathbf{x}, \Theta) = j)$. From Definition 2.1 the strength in the expectation of

$$P_\Theta(h(\mathbf{x}, \Theta) = y) - \max_{j \neq y} P_\Theta(h(\mathbf{x}, \Theta) = j)$$

is obtained by taking $Q(\mathbf{x}, j), Q(\mathbf{x}, y)$ for $P_\Theta(h(\mathbf{x}, \Theta) = j), P_\Theta(h(\mathbf{x}, \Theta) = y)$ in this layer expectation and taking the average over the training set giving the strength estimate.

From Eq. (7),

$$\bar{\rho} = \text{ar}(mr) / (E_\Theta sd(\Theta))^2.$$

The variance of mr is

$$E_{\mathbf{X}, Y} [P_\Theta(h(\mathbf{x}, \Theta) = y) - \max_{j \neq y} P_\Theta(h(\mathbf{x}, \Theta) = j)]^2 - s^2 \quad (\text{A1})$$

where s is the strength. Replacing the first term in (A1) by the average over the training set of

$$(Q(\mathbf{x}, y) - \max_{j \neq y} Q(\mathbf{x}, j))^2$$

and s by the out-of-bag estimate of s gives the estimate of $\text{ar}(mr)$. The standard deviation is given by

$$sd(\Theta) = [p_1 + p_2 + (p_1 - p_2)^2]^{1/2} \quad (\text{A2})$$

where

$$\begin{aligned} p_1 &= E_{\mathbf{X}, Y} (h(\mathbf{X}, \Theta) = Y) \\ p_2 &= E_{\mathbf{X}, Y} (h(\mathbf{X}, \Theta) = \hat{j}(\mathbf{X}, Y)) \end{aligned}$$

After the k th classifier is computed, $Q(\mathbf{x}, j)$ is computed, and the error computed is $\hat{j}(\mathbf{x}, y)$ for every example in the training set. Then, let p_1 be the average over all (y, \mathbf{x}) in the training set by no in the k th bagged training set of $I(h(\mathbf{x}, \Theta_k) = y)$. Then p_2 is the similar average of $I(h(\mathbf{x}, \Theta_k) = \hat{j}(\mathbf{x}, y))$. Substituting the estimate into (A2) to get an estimate of $sd(\Theta_k)$. Average the $sd(\Theta_k)$ over all k to get the final estimate of $sd(\Theta)$.

References

- Amit, Y. & Geman, D. (1997). Shape quantization and recognition with randomised tree. *Neural Computation*, 9, 1545–1588.
- Amit, Y., Blanchard, G., & Wilder, K. (1999). Multiple randomized classifier: MRCL Technical Report, Department of Statistics, University of Chicago.
- Bauer, E. & Kohavi, R. (1999). An empirical comparison of voting classification algorithms. *Machine Learning*, 36(1/2), 105–139.
- Breiman, L. (1996a). Bagging predictor. *Machine Learning*, 26(2), 123–140.
- Breiman, L. (1996b). Out-of-bag estimation, <http://www.berkeley.edu/~breiman/OOBestimation.pdf>.
- Breiman, L. (1998a). Arcing classifier (discussion paper). *Annals of Statistics*, 26, 801–824.
- Breiman, L. (1998b). Randomizing output to increase prediction accuracy. Technical Report 518, Statistical Department, UCB (in press), Machine Learning).
- Breiman, L. 1999. Using adaptive bagging to debias regression. Technical Report 547, Statistical Department, UCB.
- Breiman, L. 2000. Some insights for predictor ensemble. Technical Report 579, Statistical Department, UCB.
- Dieverich, T. (1998). An experimental comparison of three methods for combining ensemble of decision tree: Bagging, boosting and randomization, *Machine Learning*, 1–22.
- Friedman, J. & Schapire, R. (1996). Experiment with a new boosting algorithm, *Machine Learning: Proceedings of the Thirteenth International Conference*, 148–156.
- Grove, A. & Schuurman, D. (1998). Boosting in the limit: Maximizing the margin of learned ensemble. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence (AAAI-98)*.
- Ho, T. K. (1998). The random subspace method for classification forester. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 20(8), 832–844.
- Kleinberg, E. (2000). On the algorithmic implementation of stochastic discrimination. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 22(5), 473–490.
- Schapire, R., Friedman, J., Bartlett, P., & Lee, W. (1998). Boosting the margin: A negative planation for the effectiveness of voting method. *Annals of Statistics*, 26(5), 1651–1686.
- Tibshirani, R. (1996). Bias, variance, and prediction error for classification rule. Technical Report, Statistical Department, University of Toronto.
- Wolpert, D. H. & Macready, W. G. (1997). An efficient method to estimate Bagging's generalization error (in press), Machine Learning).

Received November 30, 1999

Revised April 11, 2001

Accepted April 11, 2001

Final manuscript April 11, 2001