

**CROSS-VALIDATION AND THE BOOTSTRAP:
ESTIMATING THE ERROR RATE OF A PREDICTION RULE
BY**

BRADLEY EFRON and ROBERT TIBSHIRANI

TECHNICAL REPORT NO. 176

MAY 1995

**PREPARED UNDER THE AUSPICES
OF**

**PUBLIC HEALTH SERVICE GRANTS
5 R01 CA59039-20 AND 5 R01 CA55325**

**DIVISION OF BIOSTATISTICS
STANFORD UNIVERSITY
STANFORD, CALIFORNIA**



**Cross-Validation and the Bootstrap:
Estimating the Error Rate of a Prediction Rule**

By

Bradley Efron and Robert Tibshirani

Technical Report No. 176

May 1995

Prepared Under the Auspices

Of

Public Health Service Grants

5 R01 CA59039-20 and 5 R01 CA55325

Also supported by National Science Foundation Grant DMS92-04864 and issued
as Technical Report No. 477, Department of Statistics, Stanford University

Division of Biostatistics

Stanford University

Stanford, California

**Cross-Validation and the Bootstrap:
Estimating the Error Rate of a Prediction Rule**

Bradley Efron and Robert Tibshirani

Abstract

A training set of data has been used to construct a rule for predicting future responses. What is the error rate of this rule? The traditional answer to this question is given by cross-validation. The cross-validation estimate of prediction error is nearly unbiased, but can be highly variable. This article discusses bootstrap estimates of prediction error, which can be thought of as smoothed versions of cross-validation. A particular bootstrap method, the 632+ rule, is shown to substantially outperform cross-validation in a catalog of 24 simulation experiments. Besides providing point estimates, we also consider estimating the variability of an error rate estimate. All of the results here are nonparametric, and apply to any possible prediction rule. The simulations include "smooth" prediction rules like Fisher's Linear Discriminant Function, and unsmooth ones like Nearest Neighbors.

1. Introduction

This article concerns estimating the error rate of a prediction rule that has been constructed from a training set of data. The training set $\mathbf{x} = (x_1, x_2, \dots, x_n)$ consists of n observations $x_i = (t_i, y_i)$, with t_i being the predictor or feature vector and y_i being the response. On the basis of \mathbf{x} the statistician constructs a prediction rule $r_{\mathbf{x}}(t)$, and wishes to estimate the error rate of this rule if it were used to predict future responses from their predictor vectors.

Cross-validation, the traditional method of choice for this problem, provides a nearly unbiased estimate of the future error rate. However the low bias of cross-validation is often paid for by high variability. We will show that suitably defined bootstrap procedures can substantially reduce the variability of error rate predictions. The gain in efficiency in our catalog of simulations is roughly equivalent to a 60% increase in the size of the training set. The bootstrap procedures are nothing more than smoothed versions of cross-validation, with some adjustments made to correct for bias.

We will be mainly interested in the situation when the response is dichotomous. This is illustrated in Figure 1, where the $n = 20$ observations $x_i = (t_i, y_i)$ in the training set \mathbf{x} each consist of a bivariate feature vector t_i and a 0 – 1 response y_i ; 12 of the points are labeled “0” and 8 are labelled “1”. Two different prediction rules are indicated. The left panel shows the prediction rule based on Fisher’s linear discriminant function (LDF) as in Efron (1983). The rule $r_{\mathbf{x}}(t)$ will predict $y = 0$ if t lies to the lower left of the LDF boundary, and $y = 1$ if t lies to the upper right. The right panel shows the nearest neighbor (NN) rule, in which future t vectors will have y predicted according to the label of the nearest observation in the training set. We wish to estimate the error rates of the two prediction rules.

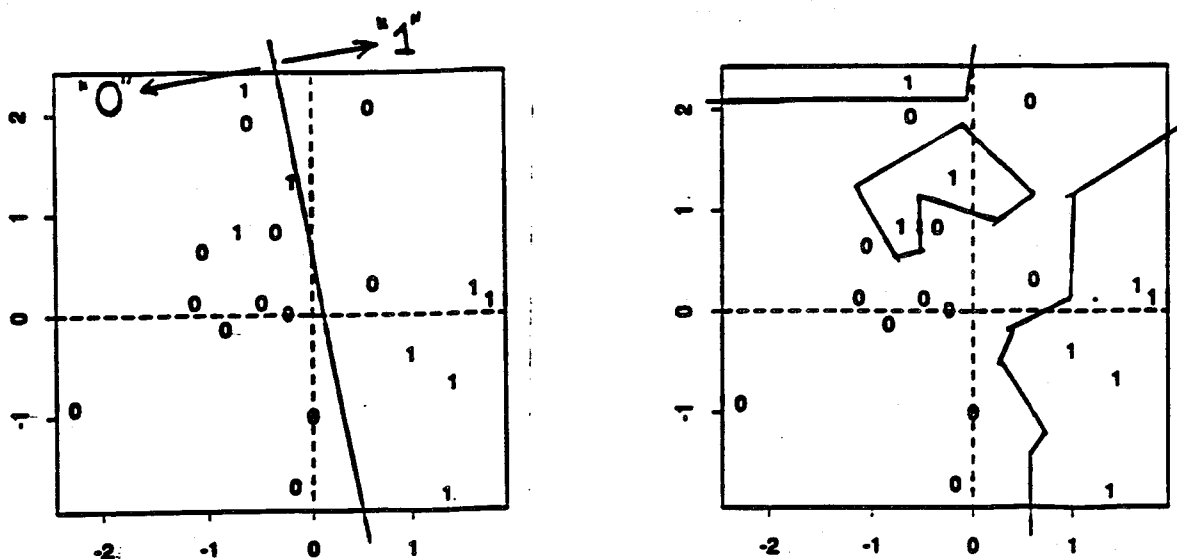


Figure 1. A training set consisting of $n = 20$ observations, 12 labelled “0” and 8 labelled “1”.
Left Panel: Linear discriminant function predicts “0” to lower left of solid line, “1” to upper right.
Right Panel: Nearest neighbor rule predicts “1” in the 3 indicated islands, “0” elsewhere.

The data shown in Figure 1 was generated as part of the extensive simulation experiments described in Section 4. In this case the y_i were selected randomly while the t_i were bivariate normal vectors whose means depended on y_i ,

$$y_i = \begin{cases} 0 \\ 1 \end{cases} \quad \text{prob} \quad \begin{matrix} \frac{1}{2} \\ \frac{1}{2} \end{matrix} \quad \text{and} \quad t_i|y_i \sim N_2 \left(\begin{pmatrix} y_i - \frac{1}{2} \\ 0 \end{pmatrix}, I \right), \quad (1.1)$$

independently for $i = 1, 2, \dots, n = 20$.

Table 1 shows results from the simulations. Cross-validation is compared with the bootstrap-based estimator 632+ described in Section 3. Cross-validation is nearly unbiased as an estimator of the true error rate for both rules, LDF and NN, but the bootstrap-based estimator has root-mean-squared error (RMS) only 80% as large. These results are fairly typical of the 24 simulation experiments reported in Section 4. The bootstrap estimator in these experiments were run with only 50 bootstrap replications per training set, but this turns out to be enough for most purposes, as the *internal variance* calculations of Section 2 show.

	LDF		NN	
	Exp	RMS	Exp	RMS
CV1:	.362	.123	.419	.123
632+:	.357	.096	.380	.099
True:	.357		.418	

Table 1. Error rate estimation for situation (1.1); CV1 is the cross-validation estimate based on omitting one observation at a time from the training set; 632+ is the bootstrap-based estimator described in Section 3. The table shows the expectation and root-mean-squared error of the two estimates, for both the LDF and NN prediction rules. In both cases, $RMS(632+)/RMS(CV1)$ is about 80%.

The bootstrap has other important advantages besides providing more accurate point estimates for prediction error. The bootstrap replications also provide a direct assessment of variability for estimated parameters in the prediction rule. For example Efron and Gong (1983) discuss the stability of the “significant” predictor variables chosen by a complicated step-wise logistic regression program. Section 5 provides another use for the bootstrap replications: to estimate the *variance* of a point estimate of prediction error.

Section 2 begins with a discussion of *bootstrap smoothing*, a general approach to reducing the variability of nonparametric point estimators. When applied to the prediction problem, bootstrap smoothing gives a smoothed version of cross-validation having considerably reduced variability, but with an upward bias. An ad hoc bias correction discussed in Section 3 results in the 632+ estimates of Table 1. The 632+ estimator is shown to substantially outperform ordinary cross-validation in the catalog of 24 sampling experiments described in Section 4. Section 5 shows how

the same bootstrap replications that provide a point estimate of prediction error can also provide an assessment of variability for that estimate. The distance argument underlying the 632+ rule is discussed in Section 6, along with other bias-correction techniques.

There has been considerable work in the literature on cross-validation and the bootstrap for error rate estimation. A good general discussion can be found in McLachlan (1992). Key references for cross-validation are Stone (1974, 1977) and Allen (1974). Efron (1983) proposed a number of bootstrap estimates of prediction error, including the optimism and the .632 estimate. Efron (1986) studied estimates of the “in-sample” prediction error problem including generalized cross-validation (Wahba, 1980) and the Cp statistic of Mallows (1973). The use of cross-validation and the bootstrap for model selection was studied by Breiman (1992), Breiman and Spector (1992), Shao (1993) and Zhang (1992). Breiman and Spector demonstrated that leave-one-out cross-validation has high variance if the prediction rule is unstable: the reason is the leave-one-out training sets are too similar to the full dataset. Five or ten fold cross-validation displayed lower variance in this case. A study of cross-validation and bootstrap methods for tree-structured models was carried out by Crawford (1989). There has also been substantial work on the prediction error problem in the machine learning and pattern recognition fields: see for example the simulation studies of Jain, Dubes & Chen (1987), and Chernick, Murthy, and Nealy (1985, 1986). Kohavi (1995) performed a particularly interesting study that renewed our interest in this problem.

2. Cross-Validation and the Leave-One-Out Bootstrap

This section discusses a bootstrap smoothing of cross-validation that reduces the variability of error-rate estimates. The notation $Q[y, r]$ will indicate the discrepancy between a predicted value r and the actual response y . We are particularly interested in the dichotomous situation where both y and r are either 0 or 1, with

$$Q[y, r] = \begin{cases} 0 & \text{if } r = y \\ 1 & \text{if } r \neq y. \end{cases} \quad (2.1)$$

We will also employ the shorter notation

$$Q(x_0, \mathbf{x}) = Q[y_0, r_{\mathbf{x}}(t_0)] \quad (2.2)$$

to indicate the discrepancy between the predicted value and response for a test point $x_0 = (t_0, y_0)$, when using the rule $r_{\mathbf{x}}$ based on training set \mathbf{x} .

Suppose that the observations $x_i = (t_i, y_i)$ in the training set are a random sample from some distribution F ,

$$x_1, x_2, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} F, \quad (2.3)$$

and that $x_0 = (t_0, y_0)$ is another independent draw from F , called a test point. The *true error rate* of the rule $r_{\mathbf{x}}$ is

$$\text{Err} = \text{Err}(\mathbf{x}, F) = E_{oF} Q(x_0, \mathbf{x}) = E_{oF} Q[y_0, r_{\mathbf{x}}(t_0)], \quad (2.4)$$

the notation E_{oF} indicating that only $x_0 = (t_0, y_0)$ is random in (2.4), \mathbf{x} and $r_{\mathbf{x}}$ being fixed.

We will compare error rate estimators in terms of their ability to predict Err . Section 4 briefly discusses estimating instead the *expected true error*

$$\mu = \mu(F) = E_F\{\text{Err}\} = E_F E_{oF} Q(x_0, \mathbf{x}). \quad (2.5)$$

The results in this case are somewhat more favorable to the bootstrap estimator.

The *Apparent error rate* (or *resubstitution rate*), is

$$\overline{\text{err}} = \text{Err}(\mathbf{x}, \hat{F}) = E_{o\hat{F}} Q(x_0, \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n Q[y_i, r_{\mathbf{x}}(t)] \quad (2.6)$$

with \hat{F} indicating the *empirical distribution* that puts probability $1/n$ on each observation x_1, x_2, \dots, x_n ; $\overline{\text{err}}$ tends to be biased downward as an estimate of Err because the training set \mathbf{x} has been used twice, both to construct the rule and to test it.

Cross-validation (Stone 1974, Geisser 1975) avoids this problem by removing the data point to be predicted from the training set. The ordinary cross-validation estimate of prediction error is

$$\widehat{\text{Err}}^{(cv1)} = \frac{1}{n} \sum_{i=1}^n Q[y_i, r_{\mathbf{x}_{(i)}}(t)] = \frac{1}{n} \sum_{i=1}^n Q(x_i, \mathbf{x}_{(i)}), \quad (2.7)$$

where $\mathbf{x}_{(i)}$ is the training set with the i th observation removed. $\widehat{\text{Err}}^{(cv1)}$ is *leave-one-out* cross-validation; the k -fold version $\widehat{\text{Err}}^{(cvk)}$ partitions the training set into k parts, predicting in turn the observations in each part from the training sample formed from all of the remaining parts.

The statistic $\widehat{\text{Err}}^{(cv1)}$ is a discontinuous function of the training set \mathbf{x} when $Q[y, r]$ itself is discontinuous as in (2.1). *Bootstrap smoothing* is a way of reducing the variance of such functions by averaging. Suppose that $Z(\mathbf{x})$ is an unbiased estimate of a parameter of interest, say

$$\zeta(F) = E_F\{Z(\mathbf{x})\}. \quad (2.8)$$

By definition the nonparametric maximum likelihood estimate (MLE) of the same parameter ζ is

$$\hat{\zeta} = \zeta(\hat{F}) = E_{\hat{F}}\{Z(\mathbf{x}^*)\}. \quad (2.9)$$

Here \mathbf{x}^* is a random sample from \hat{F} ,

$$x_1^*, x_2^*, \dots, x_n^* \stackrel{\text{i.i.d.}}{\sim} \hat{F}, \quad (2.10)$$

i.e., a bootstrap sample. The bootstrap expectation in (2.9) smooths out discontinuities in $Z(\mathbf{x})$, usually reducing its variability. However $\hat{\zeta}$ may now be biased as an estimate of ζ . Breiman (1994) introduces a very similar idea under the sobriquet “bagging”.

Now consider applying bootstrap smoothing to $Z_i(\mathbf{x}) = Q(x_i, \mathbf{x}_{(i)})$, with x_i fixed. The non-parametric MLE of $E_F Z_i(\mathbf{x})$ is $\hat{\zeta}_i = E_{\hat{F}_{(i)}}\{Q(x_i, \mathbf{x}_{(i)}^*)\}$, where $\mathbf{x}_{(i)}^*$ is a bootstrap sample from the empirical distribution on $\mathbf{x}_{(i)}$,

$$\hat{F}_{(i)} : \text{probability } \frac{1}{n-1} \text{ on } x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n. \quad (2.11)$$

It might be argued that the bootstrap samples $\mathbf{x}_{(i)}^*$ should be of size $n-1$ instead of n , but there is no advantage to this. In what follows we will take bootstrap samples from $\hat{F}_{(i)}$ to be of size n , and indicate them by \mathbf{x}^* rather than $\mathbf{x}_{(i)}^*$, so

$$\hat{\zeta}_i = E_{\hat{F}_{(i)}}\{Q(x_i, \mathbf{x}^*)\}. \quad (2.12)$$

Notice that an \mathbf{x}^* sample drawn from $\hat{F}_{(i)}$ never contains the point x_i .

Applying (2.10), (2.11) to each case i in turn leads to the *leave-one-out bootstrap*,

$$\widehat{\text{Err}}^{(1)} = \frac{1}{n} \sum_{i=1}^n E_{\hat{F}_{(i)}}\{Q(x_i, \mathbf{x}^*)\}, \quad (2.13)$$

a smoothed version of $\widehat{\text{Err}}^{(cv1)}$. This estimate predicts the error at point i only from bootstrap samples that do not contain the point i .

The actual calculation of $\widehat{\text{Err}}^{(1)}$ is a straightforward bootstrap exercise. Ordinary bootstrap samples $\mathbf{x}^* = (x_1^*, x_2^*, \dots, x_n^*)$ are generated as in (2.10), so \mathbf{x}^* is a random draw of size n , with replacement, from $\{x_1, x_2, \dots, x_n\}$. A total of B such samples are independently drawn, say $\mathbf{x}^{*1}, \mathbf{x}^{*2}, \dots, \mathbf{x}^{*B}$, with $B = 50$ in our simulations, as discussed later. Let N_i^b be the number of times x_i is included in the b th bootstrap sample, and define

$$I_i^b = \begin{cases} 1 & \text{if } N_i^b = 0 \\ 0 & \text{if } N_i^b > 0. \end{cases} \quad (2.14)$$

Also define

$$Q_i^b = Q(x_i, \mathbf{x}^{*b}) = Q[y_i, r_{\mathbf{x}^{*b}}(t_i)]. \quad (2.15)$$

Then

$$\widehat{\text{Err}}^{(1)} = \frac{1}{n} \sum_{i=1}^n \hat{E}_i \quad \text{where} \quad \hat{E}_i = \sum_b I_i^b Q_i^b / \sum_b I_i^b. \quad (2.16)$$

This definition agrees with (2.13) because a bootstrap sample that has $I_i^b = 1$ is the same as a bootstrap sample from $\hat{F}_{(i)}$, see Efron (1992). A slightly different definition appears in Efron (1983), (where $\widehat{\text{Err}}^{(1)}$ is $\hat{e}^{(0)}$), namely $\sum_i \sum_b I_i^b Q_i^b / \sum_i \sum_b I_i^b$, but the two definitions agree as $B \rightarrow \infty$, and produced nearly the same results in our simulations.

There is another way to view cross-validation and $\widehat{\text{Err}}^{(1)}$, as estimates of the average error $\mu(F)$. Direct application of the bootstrap gives the plug-in estimate $\mu(\hat{F}) = E_{\hat{F}} E_{o\hat{F}} Q(x_0, \mathbf{x})$. This

estimate, discussed in Section 6, tends to be biased downward. The reason is that \hat{F} is being used twice: as the population say F_0 from which bootstrap training sets \mathbf{x}^* are drawn, and as the population F_1 from which test points X_0 are drawn. Let us write $\mu(F)$ explicitly as a function of both F_1 and F_0 :

$$\mu(F_1, F_0) = E_{F_1}\{E_{F_0}Q[Y_0, r_{\mathbf{x}}(T_0)]\} = E_{F_0}E_{F_1}Q[Y_0, r_{\mathbf{x}}(T_0)] \quad (2.17)$$

where for convenience we have switched the order of expectation in the second expression. We assume that in the unknown true state of affairs, $F_1 = F_0 = F$. Plugging in \hat{F} for the test distribution F_0 gives

$$\mu(F_1, \hat{F}) = \frac{1}{n} \sum_{i=1}^n E_{F_1}Q[Y_i, r_{\mathbf{x}}(T_i)] \quad (2.18)$$

The remaining task is to estimate the training sample distribution F_1 . Ideally we would take $F_1 = F$. Notice that for continuous populations F the probability of the test point $X_0 = x_i$ appearing in a training sample drawn from $F_1 = F$ is zero. The plug-in estimate $\mu(\hat{F}) = \mu(\hat{F}, \hat{F})$ uses \hat{F} for F_1 . With this choice, the probability that $X_0 = x_i$ appears in the training sample is $1 - (1 - 1/n)^n \approx .632$. Hence $\mu(\hat{F})$ uses training samples that are too close to the test points, leading to potential underestimation of the error rate. Cross-validation uses the leave-one training samples to ensure that the training samples do not contain the test point. That is, cross-validation estimates $E_{F_1}Q[Y_i, r_{\mathbf{x}}(T_i)]$ by

$$\hat{E}_{F_1}Q[Y_i, r_{\mathbf{x}}(T_i)] = Q[Y_i, r_{\mathbf{x}_{(i)}}(T_i)] \quad (2.19)$$

On the other hand, use of the estimate $\hat{F}_1 = \hat{F}_{(i)}$ in each term $E_{F_1}Q[Y_i, r_{\mathbf{x}}(T_i)]$, gives the leave-one-out bootstrap estimate $\widehat{\text{Err}}^{(1)}$.

The efficacy of bootstrap smoothing is shown in Figure 2, where the solid line plots the standard deviation ratio for the 24 sampling experiments of Section 4. The horizontal axis is the expected true error μ for each experiment, (2.5). We see that $\widehat{\text{Err}}^{(1)}$ always has smaller standard deviation than $\widehat{\text{Err}}^{(cv1)}$, the median ratio over the 24 experiments being 0.79. Going from $\widehat{\text{Err}}^{(cv1)}$ to $\widehat{\text{Err}}^{(1)}$ is roughly equivalent to multiplying the size n of the training set by $1/.79^2 = 1.60$. The improvement of the bootstrap estimators over cross-validation is due mainly to the effect of smoothing. Cross-validation and the bootstrap are closely related as the argument in Section 2 of Efron (1983) shows. In smoother prediction problems, for example when y and r are continuous and $Q[y, r] = (y - r)^2$, we would expect little difference between $\widehat{\text{Err}}^{(cv1)}$ and $\widehat{\text{Err}}^{(1)}$.

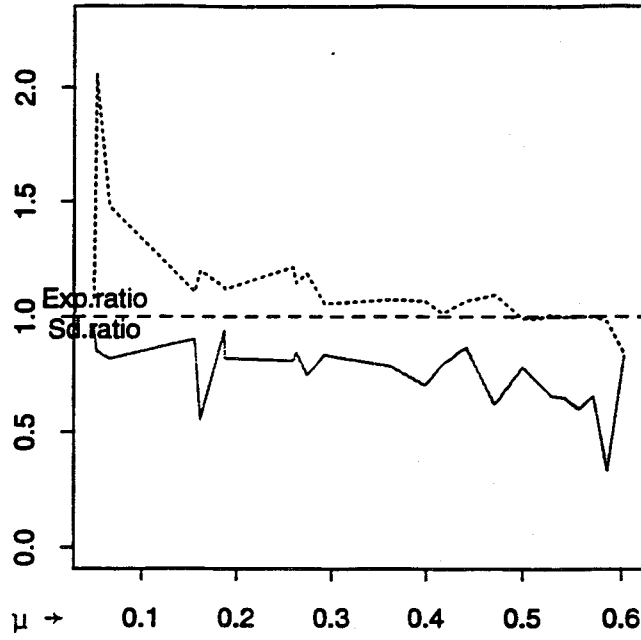


Figure 2. Ratio of standard deviations and expectations for the leave-one-out bootstrap $\widehat{\text{Err}}^{(1)}$ compared to cross-validation $\widehat{\text{Err}}^{(cv1)}$; for the 24 sampling experiments described in Section 4. Median sd ratio for the 24 experiments was .79; median expectation ratio 1.07. Plotted versus expected true error μ , (2.5).

The dotted curve in Figure 2 is the expectation ratio $E\{\widehat{\text{Err}}^{(1)}\}/E\{\widehat{\text{Err}}^{(cv1)}\}$. We see that $\widehat{\text{Err}}^{(1)}$ is biased *upwards* relative to the nearly unbiased estimate $\widehat{\text{Err}}^{(cv1)}$. This is not surprising. Let μ_n indicate the expected true error (2.5) when \mathbf{x} has sample size n . Ordinary cross-validation produces an unbiased estimate of μ_{n-1} , while k -fold cross-validation estimates μ_{n-k} . Since smaller training sets produce bigger prediction errors, larger k gives bigger upward bias $\mu_{n-k} - \mu_n$. The amount of bias will depend on the slope of the error curve μ_n at sample size n . Bootstrap samples are typically supported on about $.632n$ of the original sample points, so we might expect $\widehat{\text{Err}}^{(1)}$ to be estimating $\mu_{.632n}$. The more precise calculations in Section 8 of Efron (1983) show that $\widehat{\text{Err}}^{(1)}$ closely agrees with *half-sample cross-validation* (where \mathbf{x} is repeatedly split into equal-sized training and test sets), and that the expectation of $\widehat{\text{Err}}^{(1)}$ is $\mu_{n/2}$ to second order. The next section concerns an ad hoc bias-corrected version of $\widehat{\text{Err}}^{(1)}$ called the .632+ rule, that reduces the upward bias.

The choice of $B = 50$ bootstrap replications in our simulation experiments was based on an assessment of *internal error*, the Monte Carlo error due to using B instead of infinity replications (Efron (1992)). The same bootstrap replications that give $\widehat{\text{Err}}^{(1)}$ also give a jackknife estimate of its internal error. Let $q_i^b = I_i^b Q_i^b$, $q_i^+ = \sum_{b=1}^B q_i^b$ and $I_i^+ = \sum_{b=1}^B I_i^b$. Then the estimate of $\widehat{\text{Err}}^{(1)}$

with the b th bootstrap replication removed is

$$\widehat{\text{Err}}_{(b)}^{(1)} = \frac{1}{n} \sum_{i=1}^n \hat{E}_{i(b)} \quad \text{where} \quad \hat{E}_{i(b)} = \frac{q_i^+ - q_i^b}{I_i^+ - I_i^b}. \quad (2.17)$$

The jackknife estimate of internal standard deviation for $\widehat{\text{Err}}^{(1)}$ is then

$$\widehat{\text{sd}}_{\text{int}} = \left[\frac{B-1}{B} \sum_b (\widehat{\text{Err}}_{(b)}^{(1)} - \widehat{\text{Err}}_{(\cdot)}^{(1)})^2 \right]^{1/2}, \quad (2.18)$$

$$\widehat{\text{Err}}_{(\cdot)}^{(1)} = \sum_b \widehat{\text{Err}}_{(b)}^{(1)} / B.$$

In our simulations $\widehat{\text{sd}}_{\text{int}}$ was typically about .02 for $B = 50$. The *external* standard deviation of $\widehat{\text{Err}}^{(1)}$, i.e. the standard deviation due to randomness of \mathbf{x} , was typically .10. (Section 5 discusses the estimation of external error, also using the same set of bootstrap replications.) This gives corrected external standard deviation $[\cdot 10^2 - \cdot 02^2]^{1/2} = \cdot 098$, indicating that $B = 50$ is sufficient here.

NOTE: Definition (2.4) of prediction error, $\text{Err} = E_{oF} Q[y_0, r_{\mathbf{x}}(x_0)]$, might be called *extra-sample error* since the test point (t_0, y_0) is chosen randomly from F without reference to the training sample \mathbf{x} . Efron (1986) investigated a more restrictive definition of prediction error. For dichotomous problems let $\pi(t) = \text{Prob}_F\{y = 1|t\}$ and $\pi_i = \pi(t_i)$. The *in-sample error* of a rule $r_{\mathbf{x}}$ is defined to be

$$\text{err} = \frac{1}{n} \sum_{i=1}^n E_{o\pi_i} \{Q[y_{oi}, r_{\mathbf{x}}(t_i)]\}, \quad (2.19)$$

the notation $E_{o\pi_i}$ indicating that only $y_{oi} \sim \text{Binomial}(1, \pi_i)$ is random, \mathbf{x} and $r_{\mathbf{x}}(t_i)$ being fixed. This situation is similar to a standard regression problem in that the predictors t_i are treated as fixed at their observed values, rather than as random. In-sample error prediction is mathematically simpler than the out-sample case, and leads to quite different solutions for the error rate prediction problem, see Efron (1986).

3. The 632+ Estimator

Efron (1983) proposed the *632 estimator*

$$\widehat{\text{Err}}^{(632)} = \cdot 368 \cdot \overline{\text{err}} + \cdot 632 \cdot \widehat{\text{Err}}^{(1)}, \quad (3.1)$$

designed to correct the upward bias in $\widehat{\text{Err}}^{(1)}$ by averaging it with the downwardly biased estimate $\widehat{\text{Err}}^{(632)}$. The coefficients $\cdot 368 = e^{-1}$ and $\cdot 632$ were suggested by an argument based on the fact that bootstrap samples are supported on approximately $\cdot 632n$ of the original data points. In Efron's paper $\widehat{\text{Err}}^{(632)}$ performed better than all competitors, but the simulation studies in the 1983 paper did not include highly overfit rules like Nearest Neighbors, where $\overline{\text{err}} = 0$. Such

statistics make $\widehat{\text{Err}}^{(632)}$ itself be downwardly biased. For example if y equals 0 or 1 with probability 1/2, independently of the (useless) predictor vector t , then $\text{Err} = .50$ for any prediction rule, but the expected value of $\widehat{\text{Err}}^{(632)}$ for the Nearest Neighbor rule is $.632 \cdot .5 = .316$. Both $\widehat{\text{Err}}^{(1)}$ and $\widehat{\text{Err}}^{(cv1)}$ have the correct expectation .50 in this case. Breiman, Friedman, Olshen and Stone (1984) suggested this example.

This section proposes a new estimator $\widehat{\text{Err}}^{(632+)}$, designed to be a less biased compromise between $\overline{\text{err}}$ and $\widehat{\text{Err}}^{(1)}$. The 632+ rule puts greater weight on $\widehat{\text{Err}}^{(1)}$ in situations where the amount of overfitting, as measured by $\widehat{\text{Err}}^{(1)} - \overline{\text{err}}$, is large. In order to correctly scale the amount of overfitting, we first need to define the *no-information error rate* γ that would apply if t and y were independent, as in the example of the previous paragraph.

Let F_{ind} be the probability distribution on points $x = (t, y)$ having the same t and y marginals as F , but with y independent of t . As in (2.4), define

$$\gamma = E_{oF_{\text{ind}}} Q(x_0, \mathbf{x}) = E_{oF_{\text{ind}}} Q[y_0, r_{\mathbf{x}}(t_0)], \quad (3.2)$$

the expected prediction error for rule $r_{\mathbf{x}}$ given a test point $x_0 = (t_0, y_0)$ from F_{ind} . An estimate of γ is obtained by permuting the responses y_i and predictors t_j ,

$$\hat{\gamma} = \sum_{i=1}^n \sum_{j=1}^n Q[y_i, r_{\mathbf{x}}(t_j)] / n^2. \quad (3.3)$$

For the dichotomous classification problem (2.1), let \hat{p}_1 be the observed proportion of responses y_i equalling 1, and \hat{q}_1 be the observed proportion of predictions $r_{\mathbf{x}}(t_j)$ equalling 1. Then

$$\hat{\gamma} = \hat{p}_1(1 - \hat{q}_1) + (1 - \hat{p}_1)\hat{q}_1. \quad (3.4)$$

With a rule like nearest neighbors for which $\hat{q}_1 = \hat{p}_1$ the value of $\hat{\gamma}$ is $2\hat{p}_1(1 - \hat{p}_1)$. The multicategory generalization of (3.4) is $\hat{\gamma} = \sum_{\ell} \hat{p}_{\ell}(1 - \hat{q}_{\ell})$.

The *relative overfitting rate* is defined to be

$$\hat{R} = \frac{\widehat{\text{Err}}^{(1)} - \overline{\text{err}}}{\hat{\gamma} - \overline{\text{err}}}, \quad (3.5)$$

a quantity that ranges from 0 if there is no overfitting ($\widehat{\text{Err}}^{(1)} = \overline{\text{err}}$) to 1 if the overfitting equals the no-information value $\hat{\gamma} - \overline{\text{err}}$. The “distance” argument of Section 6 suggests a less biased version of (3.1) where the weights on $\overline{\text{err}}$ and $\widehat{\text{Err}}^{(1)}$ depend on \hat{R} ,

$$\widehat{\text{Err}}^{(632+)} = (1 - \hat{w}) \cdot \overline{\text{err}} + \hat{w} \cdot \widehat{\text{Err}}^{(1)}, \quad \left[\hat{w} = \frac{.632}{1 - .368\hat{R}} \right]. \quad (3.6)$$

The weight w ranges from .632 if $\hat{R} = 0$ to 1 if $\hat{R} = 1$, so $\widehat{\text{Err}}^{(632+)}$ ranges from $\widehat{\text{Err}}^{(632)}$ to $\widehat{\text{Err}}^{(1)}$. We can also express (3.6) as

$$\widehat{\text{Err}}^{(632+)} = \widehat{\text{Err}}^{(632)} + (\widehat{\text{Err}}^{(1)} - \overline{\text{err}}) \frac{.368 \cdot .632 \cdot \hat{R}}{1 - .368\hat{R}}, \quad (3.7)$$

emphasizing that $\widehat{\text{Err}}^{(632+)}$ exceeds $\widehat{\text{Err}}^{(632)}$ by an amount depending on \hat{R} .

It may happen that $\hat{\gamma} \leq \overline{\text{err}}$ or $\overline{\text{err}} < \hat{\gamma} \leq \widehat{\text{Err}}^{(1)}$, in which case \hat{R} can fall outside of $[0, 1]$. To take care of this possibility we modify the definitions of $\widehat{\text{Err}}^{(1)}$ and \hat{R} :

$$\widehat{\text{Err}}^{(1)'} = \min(\widehat{\text{Err}}^{(1)}, \hat{\gamma}) \quad \text{and} \quad \hat{R}' = \begin{cases} (\widehat{\text{Err}}^{(1)} - \overline{\text{err}})/(\hat{\gamma} - \overline{\text{err}}) & \text{if } \widehat{\text{Err}}^{(1)}, \hat{\gamma} > \overline{\text{err}} \\ 0 & \text{otherwise.} \end{cases} \quad (3.8)$$

The 632+ rule used in the simulation experiments of Section 4 was

$$\widehat{\text{Err}}^{(632+)} = \widehat{\text{Err}}^{(632)} + (\widehat{\text{Err}}^{(1)'} - \overline{\text{err}}) \frac{.368 \cdot .632 \cdot \hat{R}'}{1 - .368\hat{R}'}, \quad (3.9)$$

Figure 3 shows that $\widehat{\text{Err}}^{(632+)}$ was a reasonably successful compromise between the upwardly biased $\widehat{\text{Err}}^{(1)}$ and the downwardly biased $\widehat{\text{Err}}^{(632)}$. The plotted values are the relative bias in each of the 24 experiments, measured by $(\text{mean} - \mu)/\mu$.

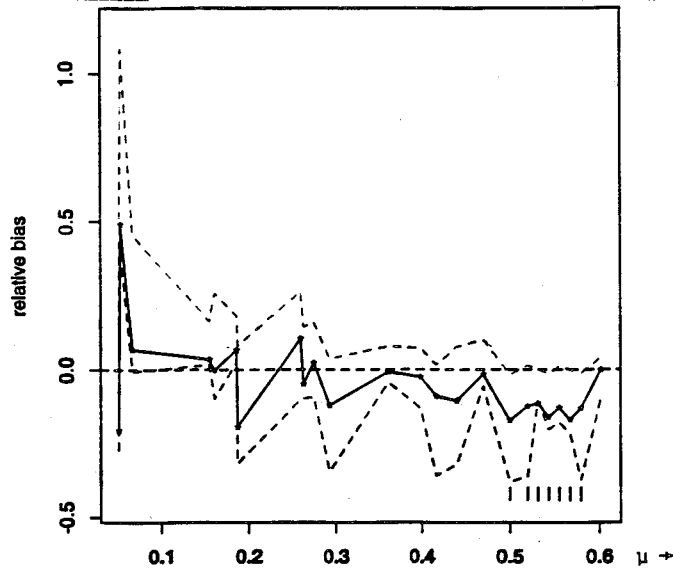


Figure 3. Relative bias of $\widehat{\text{Err}}^{(632+)}$ for the 24 experiments (solid curve), compared to $\widehat{\text{Err}}^{(1)}$ top curve, and $\widehat{\text{Err}}^{(632)}$ bottom curve; plotted values are $(\text{mean} - \mu)/\mu$. Dashes indicate the seven no-information experiments, where $\text{Err} = .50$.

The dashes in Figure 3 indicate the seven “no-information” experiments, where y was independent of t and $\text{Err} \equiv .50$. (The μ values for these seven experiments have been spread out for clarity.) In these cases the definitions in (3.8) effectively truncate $\widehat{\text{Err}}^{(632+)}$ at or near .50. This almost always gives a more accurate estimate of Err on a case by case basis, but yields a downward bias overall. To put things the other way, we could also improve the accuracy of $\widehat{\text{Err}}^{(cv1)}$ by truncating it at $\hat{\gamma}$, but then $\widehat{\text{Err}}^{(cv1)}$ would no longer be nearly unbiased.

Better bias adjustments of $\widehat{\text{Err}}^{(1)}$ are available, as discussed in Section 6. However in reducing the bias of $\widehat{\text{Err}}^{(1)}$ they lose about half of the reduction in variance enjoyed by $\widehat{\text{Err}}^{(632+)}$, and so offer less dramatic improvements over $\widehat{\text{Err}}^{(cv1)}$.

Note that $\widehat{\text{Err}}^{(632+)}$ requires no additional applications of the prediction rule, after computation of $\widehat{\text{Err}}^{(1)}$. Since 50 bootstrap samples are often sufficient, both $\widehat{\text{Err}}^{(1)}$ and $\widehat{\text{Err}}^{(632+)}$ can be less costly than $\widehat{\text{Err}}^{(cv1)}$ if n is large.

4. Twenty-Four Sampling Experiments

Table 2 describes the 24 sampling experiments performed for this study. Each experiment involved the choice of a training set size n , a probability distribution F giving \mathbf{x} as in (2.3), a dimension p for the prediction vectors t , and a prediction rule $r_{\mathbf{x}}$. 21 of the experiments were dichotomous while 3 involved 4 or more classes; 0 – 1 error (2.1) was used in all 24 experiments. The first 12 experiments each comprised 200 Monte Carlo simulations, that is 200 independent choices of the training set \mathbf{x} , while the last 12 experiments comprised 50 simulations each. The bootstrap samples were *balanced* in the sense that the indices of the bootstrap data points were obtained by randomly permuting a string of B copies of the integers 1 to n , see Davison et al. (1988), but the balancing had little effect on our results.

Here are some explanatory comments concerning the 24 experiments:

- In experiments #1–18, the response y_i equals 0 or 1 with probability .50, while the conditional distribution of $t_i|y_i$ is multivariate normal. For example experiment #3 is as described in (1.1), $t_i|y_i \sim N_2(y_i - .5, 0), I$, but #4 has the no-information form $t_i|y_i \sim N_2((0, 0), I)$, so that y_i is independent of t_i and every prediction rule has $\text{Err} = .50$. In experiment #19 the response y_i equals 1, 2, 3, or 4 with probability .25, and $t_i|y_i \sim N_i(\xi_i, I)$ with $\xi_1 = (-.5, -.5)$, $\xi_2 = (-.5, .5)$, $\xi_3 = (.5 - .5)$, $\xi_4 = (.5, .5)$.

- Experiments #13–16 are taken from Friedman (1994). There are two classes in 10 dimensions and 100 training observations. The predictors in class 1 are independent standard normal; those in class 2 are independent normal with mean $\sqrt{j}/2$ and variance $1/j$, for $j = 1, 2, \dots, 10$. All predictors are useful here, but the ones with higher index j are more so.

- Experiments #20–24 refer to real data sets taken from the machine learning archive at UC Irvine. The dimensions of these datasets are (846,19), (683,10) and (562,18) respectively (after

#	n	p	classes	μ	rule	F
1	14	5	2	.24	LDF	$\text{Normal}_5 \pm (1, 0, 0, 0, 0)$
2	14	5	2	.50	LDF	$\text{Normal}_5 \pm (0, 0, 0, 0, 0)$
3*	20	2	2	.34	LDF	$\text{Normal}_2 \pm (.5, 0)$
4	20	2	2	.50	LDF	$\text{Normal}_2 \pm (0, 0)$
5	14	5	2	.28	NN	$\text{Normal}_5 \pm (1, 0, 0, 0, 0)$
6	14	5	2	.50	NN	$\text{Normal}_5 \pm (0, 0, 0, 0, 0)$
7*	20	2	2	.41	NN	$\text{Normal}_2 \pm (.5, 0)$
8	20	2	2	.50	NN	$\text{Normal}_2 \pm (0, 0)$
9	14	5	2	.26	3NN	$\text{Normal}_5 \pm (1, 0, 0, 0, 0)$
10	14	5	2	.50	3NN	$\text{Normal}_5 \pm (0, 0, 0, 0, 0)$
11	20	2	2	.39	3NN	$\text{Normal}_2 \pm (.5, 0)$
12	20	2	2	.50	3NN	$\text{Normal}_2 \pm (0, 0)$

#	n	p	classes	μ	rule	F
13	100	10	2	.15	LDF	$N_{10}(0, I)$ versus
14				.18	NN	$N_{10}(\frac{\sqrt{I}}{2}, \frac{1}{j})$
15				.17	TREES	
16				.05	QDA	
17	20	2	2	.18	LDF	$N_2 \pm (1, 0)$
18	14	12	2	.50	LDF	$N_{12}(0, I)$
19	20	2	4	.60	LDF	$\text{Normal}_2 \pm (.5, 0)$
20	100	19	4	.26	LDF	Vehicle
21				.07	NN	Data
22	36	10	2	.07	LDF	Breast Cancer
23				.04	NN	Data
24	80	15	15	.47	3NN	Soybean Data

Table 2. The 24 Sampling Experiments; further described in text. Results for experiments #1–4 in Table 3; #5–8 in Table 4; #4–12 in Table 5; #13–16 in Table 6; #17–19 in Table 7; #20–24 in Table 8. *Experiments #3 and #7 appear in Table 1 and Figure 1.

removing incomplete observations) with 4, 2 and 15 classes respectively. We have followed Kohavi (1995) and chosen a random subset of the data to act as the training set, the training set size n chosen so that μ_n was still sloping downward. The idea is that if n is so large that the error curve is flat, then the error rate estimation problem is too easy since the potential biases arising from changing the training set size will be small. We chose training set sizes 100, 36, and 80 respectively. In the soybean data there are actually 35 categorical predictors, many having more than 2 possible values. To keep the computations manageable we used only the 15 binary predictors.

• The prediction rules are: LDF, Fisher’s linear discriminant analysis; 1-NN and 3-NN, Nearest Neighbor and 3-nearest neighbor classifier; trees- a classification tree using the `tree` function in S-PLUS; QDF; quadratic discriminant function, i.e., estimation of a separate mean and covariance in each class and then use of the Gaussian log-likelihood for classification.

Tables 3–8 report the performance of several error rate estimators in the 24 sampling experiments. In each table the Exp and Sd columns give means and standard deviations; RMS is the square root of the average squared error for estimating $\text{Err}(\mathbf{x}, F)$, (2.4). The bootstrap estimators all used $B = 50$ bootstrap replications per simulation.

The error rate estimators include $\widehat{\text{Err}}^{(632)}$, $\widehat{\text{Err}}^{(632+)}$, $\widehat{\text{Err}}^{(1)}$, and four different cross-validation

rules: $cv1$, $cv5f$, and $cv10f$ are leave-one-out, five-fold and ten-fold cross-validations, while $cv5fr$ is five-fold cross-validation averaged over 10 random choice of the partition (making the total number of recomputations equal to 50, the same as for the bootstrap rules). Also shown are other bias-corrected versions of $\widehat{Err}^{(1)}$ called $bootop$, $bc1$, $bc2$, and $\widehat{Err}^{(2)}$, see Section 6. The tables also give statistical summaries for Err , \overline{err} , and \hat{R}' , eq. (3.8).

Table 3: Results for LDF classifier

	1: (14,5)			2: (14,5,ind)			3: (20,2)			4: (20,2,ind)		
	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS
Err	0.259	0.063	0.000	0.501	0.011	0.000	0.357	0.051	0.000	0.500	0.010	0.000
Errhat1	0.327	0.116	0.147	0.500	0.115	0.115	0.388	0.101	0.104	0.502	0.087	0.088
632	0.232	0.095	0.117	0.393	0.106	0.150	0.343	0.093	0.093	0.448	0.081	0.097
632+	0.286	0.116	0.133	0.416	0.086	0.121	0.357	0.092	0.096	0.443	0.073	0.094
cv-1	0.269	0.144	0.156	0.501	0.176	0.175	0.362	0.130	0.123	0.505	0.135	0.135
cv5f	0.303	0.158	0.173	0.485	0.158	0.158	0.379	0.129	0.123	0.515	0.137	0.139
cv5fr	0.299	0.125	0.141	0.497	0.135	0.134	0.371	0.114	0.109	0.506	0.117	0.117
bootop	0.182	0.105	0.147	0.375	0.135	0.183	0.345	0.107	0.106	0.459	0.102	0.110
bc1	0.275	0.153	0.164	0.499	0.163	0.163	0.376	0.115	0.113	0.499	0.109	0.109
bc2	0.180	0.231	0.250	0.498	0.259	0.258	0.355	0.151	0.147	0.493	0.161	0.161
Errhat2	0.256	0.118	0.136	0.458	0.142	0.147	0.358	0.109	0.107	0.472	0.103	0.107
errbar	0.070	0.075	0.214	0.209	0.104	0.310	0.267	0.090	0.128	0.355	0.084	0.168
R	0.671	0.234		0.931	0.136		0.657	0.300		0.936	0.165	

Table 4: Results for 1NN classifier

	5: (14,5)			6: (14,5,ind)			7: (20,2)			8: (20,2,ind)		
	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS
Err	0.293	0.056	0.000	0.500	0.011	0.000	0.418	0.047	0.000	0.500	0.011	0.000
Errhat1	0.303	0.134	0.122	0.491	0.132	0.132	0.424	0.105	0.095	0.507	0.097	0.097
632	0.192	0.085	0.129	0.310	0.083	0.207	0.268	0.067	0.162	0.320	0.062	0.190
632+	0.257	0.127	0.120	0.413	0.094	0.128	0.380	0.101	0.099	0.439	0.068	0.092
cv-1	0.287	0.161	0.151	0.496	0.169	0.168	0.419	0.133	0.123	0.513	0.136	0.136
cv5f	0.297	0.167	0.155	0.490	0.162	0.163	0.423	0.144	0.134	0.508	0.139	0.139
cv5fr	0.297	0.144	0.133	0.496	0.138	0.138	0.420	0.122	0.110	0.509	0.117	0.117
bootop	0.107	0.047	0.194	0.172	0.046	0.331	0.150	0.037	0.271	0.180	0.035	0.322
bc1	0.307	0.139	0.128	0.490	0.143	0.143	0.423	0.109	0.100	0.506	0.102	0.101
bc2	0.313	0.158	0.149	0.486	0.179	0.179	0.421	0.131	0.124	0.503	0.126	0.126
Errhat2	0.197	0.088	0.126	0.319	0.087	0.201	0.274	0.069	0.157	0.327	0.063	0.185
errbar	0.000	0.000	0.298	0.000	0.000	0.500	0.000	0.000	0.420	0.000	0.000	0.500
R	0.641	0.252		0.922	0.134		0.853	0.169		0.949	0.099	

Table 5: Results for 3NN classifier

	9: (14,5)			10: (14,5,ind)			11: (20,2)			12: (20,2,ind)		
	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS
Err	0.273	0.065	0.000	0.500	0.011	0.000	0.399	0.062	0.000	0.501	0.011	0.000
Errhat1	0.314	0.116	0.131	0.494	0.112	0.113	0.427	0.097	0.091	0.507	0.083	0.083
632	0.245	0.099	0.110	0.400	0.100	0.142	0.346	0.084	0.093	0.412	0.074	0.115
632+	0.277	0.113	0.122	0.421	0.087	0.119	0.388	0.091	0.090	0.437	0.066	0.091
cv-1	0.263	0.154	0.154	0.496	0.173	0.173	0.401	0.139	0.126	0.509	0.138	0.138
cv5f	0.273	0.154	0.155	0.491	0.161	0.161	0.405	0.133	0.124	0.511	0.143	0.143
cv5fr	0.290	0.133	0.139	0.495	0.144	0.145	0.411	0.123	0.110	0.509	0.117	0.117
bootop	0.237	0.122	0.127	0.412	0.135	0.162	0.359	0.106	0.106	0.431	0.101	0.123
bc1	0.329	0.125	0.146	0.495	0.128	0.128	0.431	0.114	0.106	0.505	0.098	0.098
bc2	0.355	0.187	0.213	0.499	0.210	0.209	0.438	0.181	0.175	0.502	0.177	0.176
Errhat2	0.250	0.120	0.124	0.425	0.131	0.152	0.369	0.103	0.099	0.441	0.099	0.115
errbar	0.126	0.091	0.175	0.239	0.106	0.283	0.207	0.078	0.208	0.250	0.080	0.263
R	0.604	0.265		0.943	0.121		0.814	0.218		0.946	0.119	

Table 6: Results for (100,10) problem

	13: LDF			14: 1 nearest neighbor			15: trees			16: QDA		
	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS
Err	0.155	0.013	0.000	0.188	0.017	0.000	0.161	0.033	0.000	0.053	0.012	0.000
Errhat1	0.180	0.035	0.045	0.203	0.028	0.032	0.203	0.032	0.061	0.111	0.022	0.063
632	0.157	0.033	0.036	0.128	0.018	0.063	0.145	0.024	0.042	0.074	0.016	0.029
632+	0.160	0.033	0.037	0.151	0.025	0.044	0.161	0.028	0.041	0.079	0.018	0.034
cv-1	0.163	0.039	0.042	0.181	0.034	0.034	0.169	0.057	0.057	0.054	0.026	0.028
cv5f	0.169	0.040	0.044	0.193	0.036	0.036	0.180	0.046	0.056	0.066	0.028	0.032
cv10	0.164	0.035	0.038	0.185	0.034	0.034	0.171	0.042	0.046	0.060	0.026	0.029
bootop	0.157	0.037	0.039	0.073	0.010	0.116	0.116	0.025	0.058	0.049	0.015	0.019
bc1	0.167	0.038	0.042	0.202	0.030	0.033	0.188	0.042	0.059	0.063	0.026	0.029
bc2	0.142	0.046	0.049	0.201	0.039	0.041	0.160	0.066	0.075	-0.024	0.044	0.088
errbar	0.118	0.032	0.050	0.000	0.000	0.188	0.047	0.019	0.119	0.011	0.010	0.045
R	0.162	0.047		0.406	0.056		0.346	0.065		0.204	0.037	

Table 7: LDF for (20,2,+), (14,12,ind) and (20,2,4) problems

	17: (20,2, +)			18: (14,12,ind)			19: (20, 2, 4)		
	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS
Err	0.187	0.028	0.000	0.502	0.012	0.000	0.602	0.046	0.000
Errhat1	0.221	0.088	0.094	0.496	0.067	0.069	0.627	0.083	0.096
632	0.191	0.082	0.082	0.315	0.044	0.193	0.546	0.083	0.109
632+	0.199	0.087	0.088	0.438	0.065	0.093	0.602	0.098	0.106
cv-1	0.196	0.093	0.095	0.507	0.203	0.204	0.748	0.100	0.186
cv5f	0.207	0.104	0.106	0.492	0.147	0.148	0.763	0.101	0.195
cv5fr	0.204	0.092	0.094	0.504	0.093	0.094	0.730	0.085	0.162
bootop	0.188	0.091	0.092	0.179	0.035	0.326	0.554	0.113	0.129
bc1	0.203	0.093	0.094	0.492	0.105	0.107	0.609	0.105	0.112
bc2	0.169	0.115	0.115	0.485	0.193	0.195	0.577	0.157	0.162
errbar	0.141	0.078	0.092	0.005	0.021	0.498	0.406	0.101	0.224
R	0.271	0.196		0.967	0.074		0.743	0.189	

Table 8: Results for real data examples

	20: veh/LDF			21: veh/1-NN			22: breast/lda			23: breast/1-NN			24: soybean/3-NN		
	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS	Exp	Sd	RMS
Err	0.262	0.022	0.000	0.442	0.023	0.000	0.067	0.025	0.000	0.050	0.018	0.000	0.471	0.053	0.000
Errhat1	0.300	0.035	0.057	0.476	0.050	0.067	0.098	0.040	0.055	0.058	0.046	0.041	0.519	0.041	0.081
632	0.236	0.033	0.049	0.301	0.032	0.147	0.067	0.030	0.038	0.037	0.029	0.029	0.445	0.041	0.065
632+	0.249	0.034	0.044	0.395	0.054	0.077	0.072	0.033	0.040	0.040	0.034	0.032	0.463	0.040	0.062
cv-1	0.262	0.041	0.047	0.446	0.058	0.065	0.066	0.050	0.051	0.054	0.048	0.042	0.474	0.067	0.077
cv5f	0.275	0.052	0.058	0.458	0.062	0.068	0.084	0.053	0.056	0.056	0.053	0.046	0.511	0.063	0.082
cv10	0.269	0.044	0.047	0.452	0.059	0.067	0.073	0.057	0.058	0.060	0.054	0.047	0.495	0.069	0.079
bootop	0.222	0.043	0.065	0.172	0.018	0.271	0.048	0.029	0.042	0.021	0.016	0.034	0.443	0.050	0.064
errbar	0.128	0.035	0.142	0.000	0.000	0.442	0.013	0.019	0.062	0.000	0.000	0.054	0.309	0.044	0.170
R	0.280	0.036		0.640	0.067		0.202	0.081		0.134	0.106		0.367	0.062	

The results vary considerably from experiment to experiment, but in terms of rms error the 632+ rule is an overall winner. In Figure 4 the solid line graphs $rms\{\widehat{Err}^{(632+)}\}/rms\{\widehat{Err}^{(cv1)}\}$ versus the true expected error μ , (2.5). The median value of the ratio for the 24 experiments was .78. The dotted line is the rms ratio for estimating, μ rather than Err , a measure that is slightly more favorable to the 632+ rule, the median ratio now being .72.

Simulation results must be viewed with caution, especially in an area as broadly defined as the prediction problem. The smoothing argument of Section 2 strongly suggests that it *should* be possible to improve on cross-validation. With this in mind, Figure 4 says that $\widehat{Err}^{(632+)}$ has made full use of the decreased standard deviation seen in Figure 2. However the decrease in rms is less dependable than the decrease in standard deviation, and part of the rms decrease is due to the truncation at $\hat{\gamma}$ in definitions (3.8), (3.9). The truncation effect is particularly noticeable in the seven no-information experiments.

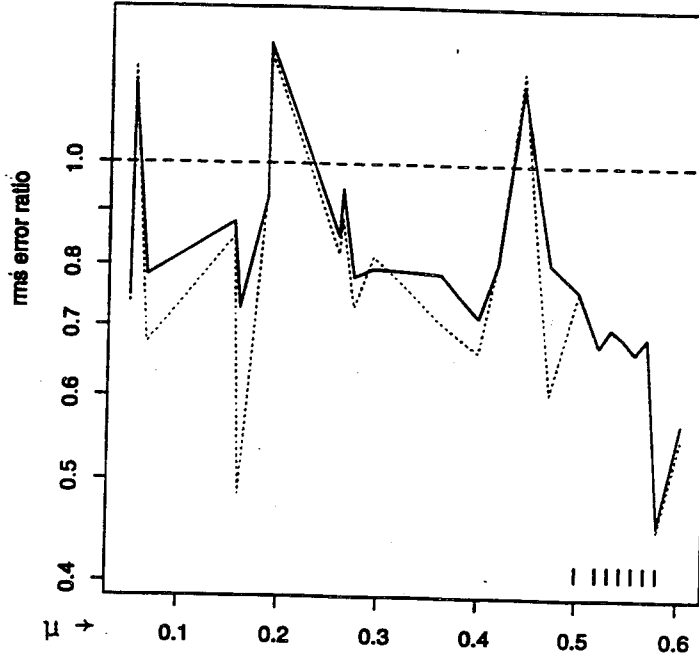


Figure 4. Solid line is $rms\{\widehat{Err}^{(632+)}\}/rms\{\widehat{Err}^{(cv1)}\}$, rms ratio from Table 3-8; plotted versus expected true error μ . Dotted line is rms ratio for estimating μ instead of Err . Dashes indicate the seven no-information experiments. The vertical axis is plotted logarithmically.

5. Estimating the Standard Error of $\widehat{Err}^{(1)}$

The same set of bootstrap replications that gives a point estimate of prediction error can also be used to assess the variability of that estimate. The method presented here, called "delta-method-after-bootstrap" in Efron (1992), works well for estimators like $\widehat{Err}^{(1)}$ that are smooth

functions of \mathbf{x} . We will discuss estimating the usual *external* standard deviation of $\widehat{\text{Err}}^{(1)}$, meaning the variability in $\widehat{\text{Err}}^{(1)}$ caused by the random choice of \mathbf{x} . The *internal* variability, due to the random choice of the B bootstrap samples (as at the end of Section 2), is also discussed since it affects the assessment of external variability. Finally we discuss estimating the standard deviation of the *difference* $\widehat{\text{Err}}^{(1)}(\text{rule 1}) - \widehat{\text{Err}}(\text{rule 2})$, where rule 1 and rules 2 are two different prediction rules applied to the same set of data.

The nonparametric delta method estimate of standard error applies to symmetrically defined statistics, those that are invariant under permutation of the points x_i in $\mathbf{x} = (x_1, x_2, \dots, x_n)$. In this case we can write the statistic as a function $S(\hat{F})$ of the empirical distribution. The form of S can depend on n , but $S(\hat{F})$ must be smoothly defined in the following sense: let $\hat{F}_{\epsilon,i}$ be a version of \hat{F} that puts extra probability on x_i ,

$$\hat{F}_{\epsilon,i} : \text{probability} \begin{cases} \frac{1-\epsilon}{n} + \epsilon & \text{on } x_i \\ \frac{1-\epsilon}{n} & \text{on } x_j \text{ for } j \neq i. \end{cases} \quad (5.1)$$

Then we need the derivatives $\partial S(\hat{F}_{\epsilon,i})/\partial \epsilon$ to exist at $\epsilon = 0$. Defining

$$\hat{D}_i = \frac{1}{n} \frac{\partial S(\hat{F}_{\epsilon,i})}{\partial \epsilon} \Big|_0, \quad (5.2)$$

the nonparametric delta method standard error estimate for $S(\hat{F})$ is

$$\widehat{\text{se}}_{\text{del}}(S) = \left[\sum_{i=1}^n \hat{D}_i^2 \right]^{1/2} \quad (5.3)$$

see Section 5 of Efron (1992). The vector $\hat{\mathbf{D}} = (\hat{D}_1, \hat{D}_2, \dots, \hat{D}_n)$ is $1/n$ times the *empirical influence function* of S .

If the prediction rule $r_{\mathbf{x}}(t)$ is a symmetric function of the points x_i in \mathbf{x} , as it usually is, then $\widehat{\text{Err}}^{(1)}$ is also symmetrically defined. The expectation in (2.12) guarantees that $\widehat{\text{Err}}^{(1)}$ will be smoothly defined in \mathbf{x} , so we can apply formulas (5.2)–(5.3).

We first consider the *ideal case* where $\widehat{\text{Err}}^{(1)}$ is based on all $B = n^n$ possible bootstrap samples $\mathbf{x}^* = (x_{i_1}, x_{i_2}, \dots, x_{i_n})$, each $i_j \in \{1, 2, \dots, n\}$. Following the notation in (2.13)–(2.15), let

$$q_i^b = I_i^b \cdot Q_i^b \quad \text{and} \quad q_i^b = \frac{1}{n} \sum_{i=1}^n q_i^b. \quad (5.4)$$

In this notation

$$\widehat{\text{Err}}^{(1)} = \frac{1}{n} \sum_{i=1}^n \hat{E}_i \quad \text{where} \quad \hat{E}_i = \sum_{b=1}^B q_i^b / \sum_{b=1}^B I_i^b. \quad (5.5)$$

We also define \hat{C}_i to be the bootstrap covariance between N_i^b and q_i^b ,

$$\hat{C}_i = \frac{1}{B} \sum_{b=1}^B (N_i^b - 1) q_i^b. \quad (5.6)$$

The following lemma is proved at the end of this section:

Lemma. The derivative (5.2) for $S(\hat{F}) = \widehat{\text{Err}}^{(1)}$ is

$$\hat{D}_i = (2 + \frac{1}{n-1}) \frac{\hat{E}_i - \widehat{\text{Err}}^{(1)}}{n} + e_n \hat{C}_i, \quad (5.7)$$

$$e_n \equiv (1 - 1/n)^{-n}.$$

A naive estimate of standard error for $\widehat{\text{Err}}^{(1)}$ would be $[\sum_i (\hat{E}_i - \widehat{\text{Err}}^{(1)})^2 / n^2]^{1/2}$, based on the false assumption that the \hat{E}_i are independent. This amounts to taking $\hat{D}_i = (\hat{E}_i - \widehat{\text{Err}}^{(1)})/n$ in (5.3). The actual influences (5.7) usually result in a larger standard error than the naive estimates.

In practice we only have B bootstrap samples, with B being much smaller than n^n . In that case we can set

$$\hat{D}_i = (2 + \frac{1}{n-1}) \frac{\hat{E}_i - \widehat{\text{Err}}^{(1)}}{n} + \frac{\sum_{b=1}^B (N_i^b - \bar{N}_i) q_i^b}{\sum_{b=1}^B I_i^b}, \quad (5.8)$$

where \hat{E}_i and $\widehat{\text{Err}}^{(1)}$ are as given in (5.5), and $\bar{N}_i = \sum_{b=1}^B N_i^b / B$. [$\bar{N}_i = 1$ for a balanced set of bootstrap samples.] The bootstrap expectation of I_i^b equals e_n^{-1} , so (5.8) goes to (5.7) as we approach the ideal case $B = n^n$.

Finally, we can use the jackknife to assess the internal error of the \hat{D}_i estimates, namely the Monte Carlo error that comes from using only a limited number of bootstrap replications. Let $\hat{D}_{i(b)}$ indicate the value of \hat{D}_i calculated from the $B - 1$ bootstrap samples not including the b th one. Simple computational formulas for $\hat{D}_{i(b)}$ are available along the lines of (2.16). The internal standard error of \hat{E}_i is given by the jackknife formula

$$\hat{\Delta}_i = [\frac{B-1}{B} \sum_b (\hat{D}_{i(b)} - \hat{D}_{i(\cdot)})^2]^{1/2}, \quad (5.9)$$

$\hat{D}_{i(\cdot)} \equiv \sum_b \hat{D}_{i(b)} / B$, the total internal error being

$$\hat{\text{se}}_{\text{int}} = [\sum_{i=1}^n \hat{\Delta}_i^2]^{1/2}. \quad (5.10)$$

This leads to an adjusted standard error estimate for $\widehat{\text{Err}}^{(1)}$,

$$\hat{\text{se}}_{\text{adj}} = [\sum_{i=1}^n \hat{D}_i^2 - \sum_{i=1}^n \hat{\Delta}_i^2]^{1/2}. \quad (5.11)$$

These formulas were applied to a single realization of experiment #3, having $n = 20$ points as in Table 2. $B = 1000$ bootstrap replications were generated, and the standard error formulas (5.3), (5.10), (5.11) calculated from the first 50, the first 100, etc. The results appear in Table 9. Using all $B = 1000$ replications gave $\hat{\text{se}}_{\text{del}} = .100$, nearly the same as $\hat{\text{se}}_{\text{adj}} = .097$. This might be

compared to the actual standard deviation .110 for $\widehat{\text{Err}}^{(1)}$ in experiment #3, though of course we expect any data-based standard error estimate to vary from sample to sample.

B	$\widehat{\text{se}}_{\text{del}}$	$\widehat{\text{se}}_{\text{int}}$	$\widehat{\text{se}}_{\text{adj}}$	B	$\widehat{\text{se}}_{\text{del}}$	$\widehat{\text{se}}_{\text{int}}$	$\widehat{\text{se}}_{\text{adj}}$
				1:00	.131	.063	.115
1:50	.187	.095	.161	101:200	.122	.077	.094
				201:300	.128	.068	.108
1:100	.131	.063	.115	301:400	.118	.068	.097
				401:500	.134	.076	.110
1:200	.119	.049	.109	501:600	.116	.085	.079
				601:700	.126	.060	.111
1:400	.110	.034	.105	701:800	.108	.076	.077
				801:900	.109	.084	.068
1:1000	.100	.023	.097	901:1000	.116	.082	.082
Mean					.121	.074	.0941
(St.dev)					(.009)	(.009)	(.0168)

Table 9. $B = 1000$ bootstrap replications from a single realization of Experiment #3, Table 2; standard error estimates (5.3), (5.10), (5.11) based on portions of the 1000 replications.

The right side of Table 9 shows $\widehat{\text{se}}_{\text{del}}$ and $\widehat{\text{se}}_{\text{adj}}$ for successive groups of 100 bootstrap replications. The values of $\widehat{\text{se}}_{\text{del}}$ are remarkably stable, but biased upward from the answer based on all 1000 replications; the bias-adjusted values $\widehat{\text{se}}_{\text{adj}}$ are less biased but about twice as variable from group to group. In this example, both $\widehat{\text{se}}_{\text{del}}$ and $\widehat{\text{se}}_{\text{adj}}$ gave useful results even for B as small as 100.

The delta method cannot be directly applied to find the standard error of $\widehat{\text{Err}}^{(632+)}$ because the 632+ rule involves $\overline{\text{err}}$, an unsmooth function of \mathbf{x} . A reasonable estimate of standard error for $\widehat{\text{Err}}^{(632+)}$ is obtained by multiplying (5.3) or (5.11) by $\widehat{\text{Err}}^{(632+)}/\widehat{\text{Err}}^{(1)}$. This is reasonable because the coefficient of variation for the two estimators was nearly the same in our experiments.

Finally, suppose that we apply two different prediction rules $r'_\mathbf{x}$ and $r''_\mathbf{x}$ to the same training set, and wish to assess the significance of the difference $\widehat{\text{Diff}} = \widehat{\text{Err}}^{(1)'} - \widehat{\text{Err}}^{(1)''}$ between the error rate estimates. For example $r'_\mathbf{x}$ and $r''_\mathbf{x}$ could be LDF and NN respectively, as in Figure 1. The previous theory goes through if we change the definition of Q_i^b in (5.4) to

$$Q_i^b = Q[y_i, r_\mathbf{x}(t_i)'] - Q[y_i, r_\mathbf{x}(t_i)'']. \quad (5.12)$$

Then the delta-method estimate of standard error for $\widehat{\text{Diff}}$ is $(\sum_i \hat{D}_i^2)^{1/2}$, where

$$\hat{D}_i = (2 + \frac{1}{n-1}) \frac{\widehat{\text{Diff}}_i - \widehat{\text{Diff}}}{n} + e_i \hat{C}_i. \quad (5.13)$$

Here $\widehat{\text{Diff}}_i = \sum_b q_i^b / \sum_b I_i^b$, $q_i^b \equiv I_i^b Q_i^b$, and everything else is as defined in (5.6), (5.7).

Proof of the Lemma. Define $f_{\epsilon,i}(j)$ to be the density function on points x_j corresponding to $\hat{F}_{\epsilon,i}$, (5.1),

$$f_{\epsilon,i}(j) = [1 + (n\delta_{ij} - 1)\epsilon]f_0(j), \quad [j = 1, 2, \dots, n] \quad (5.14)$$

with $f_0(j) = 1/n$ and δ_{ij} equalling 1 or 0 as j does or does not equal i . Also define $g_{\epsilon,i}(b)$ to be the density for the b th of the $B = n^n$ possible bootstrap samples $\mathbf{x}^* = (x_{i_1}, x_{i_2}, \dots, x_n)$,

$$g_{\epsilon,i}(b) = [(1 - \epsilon)^n (1 + \frac{n\epsilon}{1 - \epsilon})^{N_i^b}] g_0(b), \quad [b = 1, 2, \dots, n^n] \quad (5.15)$$

with $g_0(b) = 1/n^n$. It is easy to show that $g_{\epsilon,i}(b)$ is the probability of the b th bootstrap sample if the original points x_i are weighted according to $\hat{F}_{\epsilon,i}$, rather than equally. See Lemma 3 of Efron (1992).

The statistic $S(\hat{F}) = \widehat{\text{Err}}^{(1)}$ (2.16) has value

$$S(\hat{F}_{\epsilon,i}) = \sum_{j=1}^n f_{\epsilon,i}(j) \frac{\sum_b I_j^b Q_j^b g_{\epsilon,i}(b)}{\sum_b I_j^b g_{\epsilon,i}(b)} = \sum_j f_{\epsilon,i}(j) \frac{A_{\epsilon,i}(j)}{B_{\epsilon,i}(j)} \quad (5.16)$$

at $\hat{F}_{\epsilon,i}$, where $A_{\epsilon,i}(j) = \sum_b I_j^b Q_j^b g_{\epsilon,i}(b)$ and $B_{\epsilon,i}(j) = \sum_b I_j^b g_{\epsilon,i}(b)$. When $\epsilon = 0$, notice that

$$B_0(j) = \text{Prob}\{N_i^b = 0\} = e_n^{-1} \quad \text{and} \quad A_0(j)/B_0(j) = \hat{E}_j,$$

(2.16). Then

$$\frac{\partial}{\partial \epsilon} S(\hat{F}_{\epsilon,i})|_{\epsilon=0} = \sum_j \dot{f}_0(j) \frac{A_0(j)}{B_0(j)} + \sum_j f_0(j) \frac{\dot{A}_0(j)}{B_0(j)} - \sum_j f_0(j) \frac{A_0(j)}{B_0(j)} \frac{\dot{B}_0(j)}{B_0(j)}, \quad (5.17)$$

the dot indicating differentiation with respect to ϵ , with i fixed.

Denote the right side of (5.17) as **I** + **II** + **III**. Since $\dot{f}_0(j) = (n\delta_{ij} - 1)f_0(j)$,

$$\text{I} = \sum_j (\delta_{ij} - 1/n) \frac{A_0(j)}{B_0(j)} = \hat{E}_j - \widehat{\text{Err}}^{(1)}. \quad (5.18)$$

Denoting $q_i^b = I_i^b Q_i^b$ as in (5.4),

$$\dot{A}_0(j) = \sum_b q_j^b \dot{g}_0(b) = \sum_b n(N_i^b - 1) q_j^b g_0(b), \quad (5.19)$$

so

$$\begin{aligned} \text{II} &= \frac{e_n}{n} \sum_j \sum_b n(N_i^b - 1) q_j^b g_0(b) = ne_n \sum_b (N_i^b - 1) q_i^b g_0^{(b)} \\ &= ne_n \hat{C}_i. \end{aligned} \quad (5.20)$$

Similarly

$$\text{III} = -e_n \sum_j \hat{E}_j \{ \sum_b (N_i^b - 1) I_j^b g_0(b) \} = -e_n \sum_j \hat{E}_j \text{Cov}_*(N_i^b, I_j^b), \quad (5.21)$$

Cov_* indicating ordinary bootstrap covariance. But notice that $\text{Cov}_*(N_i^b, I_i^b) = -E_*\{N_i^b\}E_*(I_i^b) = -e_n^{-1}$ [since $I_i^b \cdot N_i^b \equiv 0$], and that $\sum_j \text{Cov}_*(N_i^b, I_j^b) = e_n^{-1}/(n-1)$ [by symmetry and the fact that $\sum_i \text{Cov}_*(N_i^b, I_j^b) = \text{Cov}_*(N_+^b, I_j^b) = 0$.] thus

$$\text{III} = \frac{n}{n-1} (\hat{E}_i - \widehat{\text{Err}}^{(1)}). \quad (5.22)$$

Combining (5.18), (5.19), and (5.21) gives (5.7).

6. Distance and Bias Corrections

One way to understand the biases of the various error rate estimators is in terms of the distance of test points from the training set: $\overline{\text{err}}$, (2.6), is the error rate for test points that are zero distance away from the training set \mathbf{x} , while the true value Err , (2.4), is the error rate for a new test point x_0 that may lie some distance away from \mathbf{x} . Since we expect the error rate of a rule $r_{\mathbf{x}}(t_0)$ to increase with distance from \mathbf{x} , $\overline{\text{err}}$ underestimates Err . Cross-validation has the test points nearly the right distance away from the training set and so is nearly unbiased. The leave-one-out-bootstrap $\widehat{\text{Err}}^{(1)}$, (2.13), has x_0 too far away from the bootstrap training sets \mathbf{x}^* , since these are supported on only about $.632n$ of the points in \mathbf{x} , producing an upward bias.

A quantitative version of the distance argument leads to the 632+ rule (3.6). This section presents the argument, which is really quite rough, and then goes on to discuss other more careful bias-correction methods. However these “better” methods did not produce better estimators in our experiments, the reduced bias being paid for by too great an increase in variance.

Efron (1983) used distance methods to motivate $\widehat{\text{Err}}^{(632)}$, (3.1). Here is a brief review of the arguments in that paper, which will lead up to the motivation for $\widehat{\text{Err}}^{(632+)}$. Given a system of neighborhoods around points $x = (t, y)$, let $S(x, \Delta)$ indicate the neighborhood of x having probability content Δ ,

$$\text{Prob}\{X_0 \in S(x, \Delta)\} = \Delta. \quad (6.1)$$

(In this section capital letters indicate random quantities distributed according to F , e.g., X_0 in (6.1), while lower-case x values are considered fixed.) As $\Delta \rightarrow 0$ we assume that the neighborhood $S(x, \Delta)$ shrinks toward the point x . The distance of test point x_0 from a training set \mathbf{x} is defined by its distance from the nearest point in \mathbf{x} ,

$$\delta(x_0, \mathbf{x}) = \inf_{\Delta} \{x_0 \in \cup_{i=1}^n S(x_i, \Delta)\}. \quad (6.2)$$

Figure 5 shows neighborhoods of probability content $\Delta = .05$ for a realization from the (20, 2) problem. Here we have chosen $S(x, \Delta)$ to be circles in the planes $y = 0$ and $y = 1$. That is, if $x_0 = (t_0, y_0)$ then $S(x_0, \Delta) = \{(t, y) : y = y_0 \text{ and } ||t - t_0|| \leq r\}$ where r is chosen so that (6.1) holds. Notice how the neighborhoods grow smaller near the decision boundary: this occurs because the probability in (6.1) refers not to the distribution of t but to the joint distribution of t and y .

Let

$$\mu(\Delta) = E\{Q(X_0, \mathbf{X}) | \delta(X_0, \mathbf{X}) = \Delta\}, \quad (6.3)$$

the expected prediction error for test points distance Δ from the training set. The true expected error (2.5) is given by

$$\mu = \int_0^1 \mu(\Delta) g(\Delta) d\Delta, \quad (6.4)$$

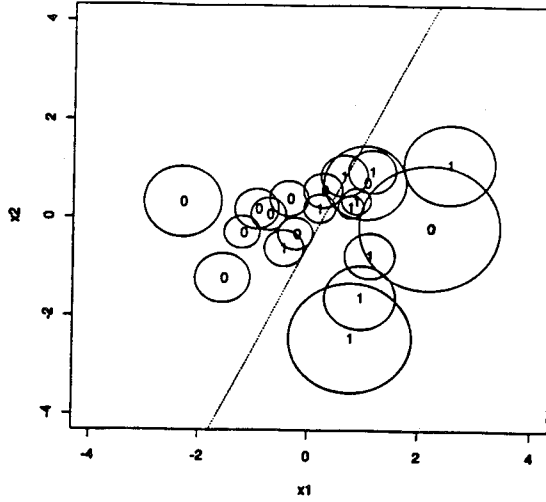


Figure 5. Two Gaussian classes in two dimensions, from the (20, 2) problem. The circles represent neighborhoods of probability content $\Delta = .05$ around each training point. The dotted line represents the LDF decision boundary.

where $g(\Delta)$ is the density of $\delta(X_0, \mathbf{X})$. Under reasonable conditions $g(\Delta)$ approaches the exponential density

$$g(\Delta) \doteq ne^{-n\Delta} \quad (\Delta \in (0, \infty)), \quad (6.5)$$

see the appendix of Efron (1983). Another important fact is that

$$\mu(0) \equiv \nu \equiv E_F\{\overline{\text{err}}\}, \quad (6.6)$$

which is just another way of saying that $\overline{\text{err}}$ is the error rate for test points zero distance away from \mathbf{x} .

We can also define a bootstrap analogue of $\mu(\Delta)$,

$$\mu_*(\Delta) = E\{Q(X_0^*, \mathbf{X}^*) | \delta(X_0^*, \mathbf{X}^*) = \Delta\}, \quad (6.7)$$

the expectation in (6.7) being over the choice of $X_1, X_2, \dots, X_n \sim F$ and then $X_0^*, X_1^*, X_2^*, \dots, X_n^* \sim \hat{F}$. Notice that if $\delta > 0$ then X_0^* must not equal any of the X_i points in \mathbf{X}^* . This and definition (2.13) give

$$\xi \equiv E\{\widehat{\text{Err}}^{(1)}\} = \int_0^1 \mu_*(\Delta) g_*(\Delta) d\Delta, \quad (6.8)$$

where $g_*(\Delta)$ is the density of $\delta^* = \delta(X_0^*, \mathbf{X}^*)$ given that $\delta^* > 0$.

Because bootstrap samples are supported on about .632n of the points in the training set, the same argument that gives $g(\Delta) \doteq ne^{n\Delta}$ also shows that

$$g_*(\Delta) \doteq .632ne^{-.632n\Delta} \quad (6.9)$$

Efron (1983) further supposes that

$$\mu_*(\Delta) \doteq \mu(\Delta) \quad \text{and} \quad \mu(\Delta) \doteq \nu + \beta\Delta \quad (6.10)$$

for some value β , the intercept ν coming from (6.6). Combining these approximations gives

$$\mu \doteq \nu + \frac{\beta}{n} \quad \xi \doteq \nu + \frac{\beta}{.632n}, \quad (6.11)$$

so

$$\mu \doteq \{.318\nu + .632\xi\} \quad (6.12)$$

Substituting $\overline{\text{err}}$ for ν and $\widehat{\text{Err}}^{(1)}$ for ξ results in $\widehat{\text{Err}}^{(.632)}$, (3.1).

Figure 6 shows $\mu(\Delta)$ for experiments #3 and #7, estimated using all of the data from each set of 200 simulations. The linearity assumption $\mu = \nu + \beta\Delta$ in (6.10) is reasonably accurate for the LDF rule of experiment #3, but not for the Nearest Neighbor rule of #7. The expected apparent error $\nu = \mu(0)$ equals 0 for Nearest Neighbors, producing a sharp bend near 0 in the $\mu(\Delta)$ curve.

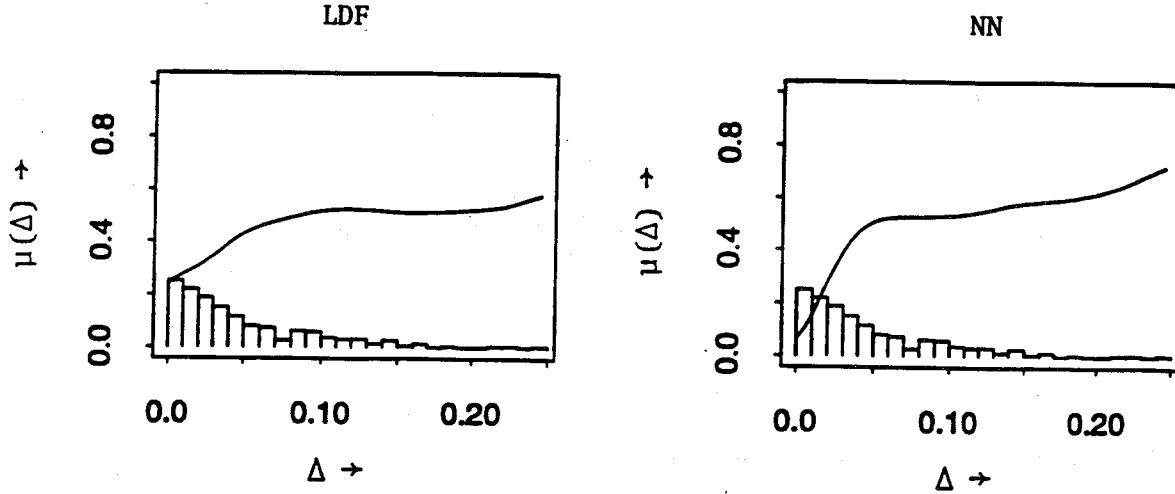


Figure 6. $\mu(\Delta)$ curves, (6.3), for experiment #3 (left) and #7 (right); linearity assumption $\mu(\Delta) = \nu + \beta\Delta$ is reasonably accurate for the LDF rule, but not for NN. The histogram, for distance δ (6.2), supports the exponential approximation (6.5).

The 632+ rule of Section 3 replaces linearity with an exponential curve, better able to match the form of $\mu(\Delta)$ seen in Figure 5. Now we assume that

$$\mu(\Delta) = \gamma - e^{-(\alpha + \beta\Delta)}. \quad (6.13)$$

Here α and β are positive parameters and γ is an upper asymptote for $\mu(\Delta)$ as Δ gets large. Formulas (6.5) and (6.9), taken literally, produce simple expressions for μ , (6.4), and for ξ , (6.8),

$$\mu = \frac{\beta\gamma + n\nu}{\beta + n} \quad \text{and} \quad \xi = \frac{\beta\gamma + .632n\nu}{\beta + .632n}. \quad (6.14)$$

Combined with $\nu = \gamma - e^{-\alpha}$ from (6.6), (6.14) gives

$$\mu = (1 - w)\nu + w\xi \quad (6.15)$$

where

$$w = \frac{.632}{1 - .368R} \quad \text{and} \quad R = \frac{\xi - \nu}{\gamma - \nu}. \quad (6.16)$$

$\widehat{\text{Err}}^{(632+)}$, (3.6), is obtained by substituting $\overline{\text{err}}$ for ν , $\widehat{\text{Err}}^{(1)}$ for ξ , and $\hat{\gamma}$, (3.4), for γ .

All of this is a reasonable plausibility argument for the 632+ rule, but not much more than that. The assumption $\mu(\Delta) \doteq \mu_*(\Delta)$ in (6.10) is particularly vulnerable to criticism, though in the example shown in Figure 2 of Efron (1983) it is quite accurate. A theoretically more defensible approach to the bias-correction of $\widehat{\text{Err}}^{(1)}$ can be made using the ANOVA decomposition arguments of the (1983) paper.

Define

$$a = -nE\{Q(X_1, \mathbf{X}) - \mu\} \quad \text{and} \quad b = n^2E\{Q(X_1, \mathbf{X}') - \mu\} \quad (6.17)$$

where

$$\mathbf{X} = (X_1, X_2, X_3, \dots, X_n) \quad \text{and} \quad \mathbf{X}' = (X_2, X_2, X_3, \dots, X_n). \quad (6.18)$$

(Both a and b will usually be positive). Then the formal ANOVA decompositions of Section 7, Efron (1983), give

$$E\{\overline{\text{err}}\} = \mu - a/n, \quad E\{\widehat{\text{Err}}^{(cv1)}\} = \mu + \frac{b}{n^2} + O\left(\frac{1}{n^3}\right), \quad E\{\widehat{\text{Err}}^{(1)}\} = \mu + \frac{b}{2n} + \left(\frac{1}{n^2}\right). \quad (6.19)$$

Also, letting $\hat{\mu} \equiv E_{\hat{F}}E_{o\hat{F}}Q(X_0^*, \mathbf{X}^*)$ denote the nonparametric MLE of $\mu = EQ(X_0, \mathbf{X})$,

$$E\{\widehat{\text{Err}}^{(1)} - \hat{\mu}\} = \frac{a}{n} + O\left(\frac{1}{n^2}\right). \quad (6.20)$$

We can combine (6.19) and (6.20) to obtain a bias-corrected version of $\widehat{\text{Err}}^{(1)}$ that, like $\widehat{\text{Err}}^{(cv1)}$, has bias of order $1/n^2$ rather than $1/n$:

$$\widehat{\text{Err}}^{(2)} = \widehat{\text{Err}}^{(1)} - \hat{\mu} + \overline{\text{err}}. \quad (6.21)$$

$\widehat{\text{Err}}^{(2)}$ can be attractively reexpressed in terms of the bootstrap covariances between I_i^b and Q_i^b . Following the notation in (5.4), (5.7), it turns out that

$$\widehat{\text{Err}}^{(2)} = \overline{\text{err}} + \frac{e_n}{n} \sum_{i=1}^n \widehat{\text{Cov}}_i, \quad (6.22)$$

where

$$\widehat{\text{Cov}}_i = \frac{1}{B} \sum_{b=1}^B (I_i^b - \bar{I}_i) \cdot Q_i^b, \quad (6.23)$$

$\bar{I}_i = \sum_b I_i^b / B$. Formula (6.22) says that we can bias-correct the apparent error rate $\overline{\text{err}}$ by adding e_n times the average covariance between I_i^b (2.14), the absence or presence of x_i in \mathbf{x}^* , and Q_i^b (2.15), whether or not $r_{\mathbf{x}^*}(t_i)$ incorrectly predicts y_i . These covariances will usually be positive.

Despite its attractions, $\widehat{\text{Err}}^{(2)}$ was an underachiever in Efron (1983), where it appears in Table 2 as $\hat{\omega}^{(0)}$, and also in the experiments here. Generally speaking, it gains only about half of the rms advantage over $\widehat{\text{Err}}^{(cv1)}$ enjoyed by $\widehat{\text{Err}}^{(1)}$. Moreover its bias-correction powers fail for cases like experiment #7, where the formal expansions (6.19), (6.20) are also misleading.

The estimator called bootop in Tables 3-8 is defined as

$$\widehat{\text{Err}}^{(\text{bootop})} = \overline{\text{err}} - \frac{1}{n} \sum_{i=1}^n \text{Cov}_*(N_i^b, Q_i^b) \quad (6.24)$$

at (2.10) of Efron (1983), bootop standing for “bootstrap optimism”. Here $\text{Cov}_*(N_i^b, Q_i^b) = \sum_b (N_i^b - 1)Q_i^b / B$. Section 7 of the 1983 paper shows that the bootop rule, like $\widehat{\text{Err}}^{(cv1)}$ and $\widehat{\text{Err}}^{(2)}$, has bias of order only $1/n^2$ instead of $1/n$. This does not keep it from being badly biased downward in several of the sampling experiments, particularly for the NN rules.

We also tried to bias-correct $\widehat{\text{Err}}^{(1)}$ using a second level of bootstrapping. For each training set \mathbf{x} , 50 second-level bootstrap samples were drawn, by resampling (one-time each) from the 50 original bootstrap samples. The number of distinct original points x_i appearing in a second-level bootstrap sample is approximately $.502 \cdot n$, compared with $.632 \cdot n$ for a first-level sample. Let $\widehat{\text{Err}}^{(\text{sec})}$ be the $\widehat{\text{Err}}^{(1)}$ statistic (2.16) computed using the second-level samples instead of the first. The rules called bc1 and bc2 in Tables 3-7 are linear combinations of $\widehat{\text{Err}}^{(1)}$ and $\widehat{\text{Err}}^{(\text{sec})}$,

$$\begin{aligned} \widehat{\text{Err}}^{(\text{bc1})} &= 2 \cdot \widehat{\text{Err}}^{(1)} - \widehat{\text{Err}}^{(\text{sec})} \\ \text{and} \\ \widehat{\text{Err}}^{(\text{bc2})} &= 3.83\widehat{\text{Err}}^{(1)} - 2.83\widehat{\text{Err}}^{(\text{sec})}. \end{aligned} \quad (6.25)$$

The first of these is suggested by the usual formulas for bootstrap bias correction. The second is based on linearly extrapolating $\hat{\mu}_{.502n} = \widehat{\text{Err}}^{(\text{sec})}$ and $\hat{\mu}_{.632n} = \widehat{\text{Err}}^{(1)}$ to an estimate for $\hat{\mu}_n$. The bias correction works reasonably in Tables 3-7, but once again at a substantial price in variability.

All in all, we feel that bias was not a major problem for $\widehat{\text{Err}}^{(632+)}$ in our simulation studies, and that attempts to correct bias were too expensive in terms of added variability. At the same time it seems possible that further research might succeed in producing a better compromise between the unbiasedness of cross-validation and the reduced variance of the leave-one-out bootstrap.

Acknowledgements

We would like to thank Ronny Kohavi for reviving our interest in this problem, and Jerry Friedman, Trevor Hastie and Richard Olshen for enjoyable and fruitful discussions. The second author was supported by the Natural Sciences and Engineering Research Council of Canada.

References

- Breiman, L. (1994). Bagging predictors. Technical Report, University of California, Berkeley.
- Breiman, L. and Spector, P. (1992). Submodel selection and evaluation in regression: the x -random case. *Intern. Statist. Rev.* **60**, 291-319.
- Breiman, L., Friedman, J., Olshen, R. and Stone, C. (1984). *Classification and Regression Trees*. Wadsworth.
- Chernick, M., Murthy, V. and Nealy, C. (1985). Application of bootstrap and other resampling methods: evaluation of classifier performance. *Pattern Recognition Letters* **4**, 167-178.
- Chernick, M., Murthy, V. and Nealy, C. (1986). "Correction note to Application of bootstrap and other resampling methods: evaluation of classifier performance". *Pattern Recognition Letters* **4**, 133-142.
- Cosman, P., Perlmutter, K., Perlmutter, S., Olshen, R. and Gray, R. (1991). Training sequence size and vector quantizer performance. In *25th Asilomar Conference on Signals, Systems, and Computers*, Nov. 4-6, 1991, by Ray R. Chen, IEE Computer Society Press, Los Alamitos, CA, pp. 434-448.
- Dawid, A., Hinkley, D. and Schechtman, E. (1986). Efficient bootstrap simulations. *Biometrika* **73** 555-566.
- Efron, B. (1979). Bootstrap methods: another look at the jackknife. *Ann. Statist.* **7**, 1-26.
- Efron (1983). Estimating the error rate of a prediction rule: some improvements on cross-validation. *J. Amer. Statist. Assoc.* **78**, 316-331.
- Efron, B. (1992). Jackknife-after-bootstrap standard errors and influence functions (with Discussion). *J. Royal Statist. Society, B*, 83-111.
- Efron, B. and Tibshirani, R. (1993). *An Introduction to the Bootstrap*. Chapman and Hall.
- Friedman, J. (1994). Flexible metric nearest neighbour classification. Technical Report, Stanford University.
- Geisser, S. (1975). The predictive sample reuse method with applications. *J. Amer. Statist. Assoc.* **70**, 320-328.
- Jain, A., Dubes, R. P. and Chen, C. (1987). Bootstrap techniques for error estimation. *IEEE Trans. On Pattern Analysis and Mach. Intell.* **9**, 628-633.
- Kohavi, R. (1995). A study of cross-validation and bootstrap for accuracy assessment and model selection. Technical Report, Stanford University.
- Mallows, C. (1973). Some comments on cp. *Technometrics*, 661-675.
- McLachlan, G. (1992). *Discriminant Analysis and Statistical Pattern Recognition*. Wiley, New York.

- Stone, M. (1974). Cross-validators choice and assessment of statistical predictions. *J. Roy. Statist. Soc* **36**, 111-147.
- Wahba, G. (1980). Spline bases, regularization, and generalized cross-validation for solving approximation problems with large quantities of noisy data. In *Proceedings of the International Conference on Approximation Theory in Honour of George Lorenz*, Academic Press, Austin, Texas.