

University of California, Berkeley

From the SelectedWorks of Maya Petersen

September 2007

Cross-validated Bagged Learning

Contact
Author

Start Your Own
SelectedWorks

Notify Me
of New Work



Available at: http://works.bepress.com/maya_petersen/1



Cross-validated bagged learning

Maya L. Petersen^a, Annette M. Molinaro^{b,1}, Sandra E. Sinisi^a,
Mark J. van der Laan^{a,*}

^a*Division of Biostatistics, University of California, Berkeley School of Public Health, Earl Warren Hall 7360, Berkeley, CA 94720-7360, USA*

^b*School of Public Health, Yale University, 60 College Street, New Haven, CT 06520-8034, USA*

Received 20 September 2005

Available online 19 July 2007

Abstract

Many applications aim to learn a high dimensional parameter of a data generating distribution based on a sample of independent and identically distributed observations. For example, the goal might be to estimate the conditional mean of an outcome given a list of input variables. In this prediction context, bootstrap aggregating (bagging) has been introduced as a method to reduce the variance of a given estimator at little cost to bias. Bagging involves applying an estimator to multiple bootstrap samples and averaging the result across bootstrap samples. In order to address the curse of dimensionality, a common practice has been to apply bagging to estimators which themselves use cross-validation, thereby using cross-validation within a bootstrap sample to select fine-tuning parameters trading off bias and variance of the bootstrap sample-specific candidate estimators. In this article we point out that in order to achieve the correct bias variance trade-off for the parameter of interest, one should apply the cross-validation selector *externally* to candidate bagged estimators indexed by these fine-tuning parameters. We use three simulations to compare the new cross-validated bagging method with bagging of cross-validated estimators and bagging of non-cross-validated estimators.

© 2007 Elsevier Inc. All rights reserved.

AMS 2000 subject classification: 62G08; 62G09; 68T05

Keywords: Bootstrap aggregation; Data-adaptive regression; Resistant HIV; Deletion/Substitution/Addition algorithm

* Corresponding author. Fax: +1 510 643 5163.

E-mail address: laan@stat.berkeley.edu (M.J. van der Laan)

URL: <http://www.stat.berkeley.edu/~aan> (M.J. van der Laan).

¹ Fax: +1 203 785 6912.

1. Introduction and motivation

Many applications aim to use a learning data set from a particular data generating distribution to construct a predictor of an outcome as a function of a collection of input variables. One can define an optimal predictor as a parameter of the data generating distribution by defining it as the function of input variables which minimizes the expectation of a particular loss function (of the experimental unit and the candidate regression) with respect to the true data generating distribution. If one selects the squared error loss function (i.e., the square of the difference between the outcome and predicted value), then this optimal predictor is the conditional mean of the outcome given the input variables. In the statistical literature such a location parameter of the conditional distribution of the outcome given the input variables is often referred to as a regression.

In many applications the number of input variables is very large. As a consequence, assuming a fully parameterized regression model and minimizing the empirical mean of the loss function is likely to yield poor estimators. For example, assuming a linear regression model with all main terms and minimizing the sum of squared residual errors is likely to yield a poor estimator since the number of main terms will typically be too large resulting in over-fitting, and other functional forms of the input variables should be considered. In other words, many current applications demand nonparametric regression estimators. Because of the curse of dimensionality, minimizing the empirical mean of the loss function, i.e., the empirical risk, over all allowed regression functions results in an over-fit of the data.

As a consequence, many estimators follow the sieve loss-based estimation strategy. A sequence of subspaces indexed by fine-tuning parameters is selected, the empirical risk over each subspace is minimized or locally minimized to obtain a subspace-specific minimum empirical risk estimator, and the fine-tuning parameter (i.e., the subspace) is selected using an appropriate method to trade-off between bias and variance. Examples of fine-tuning parameters indexing constraints on the space of regression functions include initial dimension reduction, the number of terms in the regression model, and the complexity of the allowed functional forms (for example, the basis functions used). Each specification of the fine-tuning parameters corresponds to a candidate estimator of the true underlying regression. In order to select among these candidate estimators, or in other words to select the fine-tuning parameters, most algorithms either minimize a penalized empirical risk or minimize the cross-validated risk.

Application of such “machine learning” algorithms to a data set commonly results in a very low-dimensional fit. For example, in a recent HIV data application involving prediction of viral replication capacity based on the mutation profile of the virus, in spite of the fact that the employed algorithm searched over a high dimensional space of regression functions, a linear regression with two main terms and a single interaction was selected [1]. Although such an estimator is based on a sensible trade-off between bias and variance, the resulting fit is disappointing from two perspectives. First, in many applications the true regression is believed to be a function of almost all variables, with many variables making small contributions. Second, a practitioner often wishes to obtain a measure of importance for each variable considered, and such a low dimensional fit reflects zero importance for all variables that do not appear in the estimator. It has been common practice to address the second issue by reporting many of the fits the algorithm has searched over, and to summarize these different fits in a particular manner. Initially, we also followed this approach, but came to the conclusion that the statistical interpretation of such a summary measure is unclear.

Based on these concerns, the following statistical challenge can be formulated: construction of nonparametric regression estimators that are high dimensional, so that the majority of variables

contribute to the obtained regression, and yet still correspond with a sensible trade-off between bias and variance and thereby have good asymptotic convergence properties. In the current article, we address this challenge. In order to construct high dimensional learners we employ the existing machine learning method “bootstrap aggregating”, “bagging”, or “aggregate prediction”, as introduced by Breiman [3]. However, in order to establish a sensible trade-off between bias and variance, we provide a fundamental improvement to the current practice of bagging by changing the way cross-validation enters into the methodology.

Breiman suggested bagging as a method to stabilize and thereby improve upon a given highly variable estimator. Specifically, given an estimator, Breiman defined a corresponding bagged estimator as the average across bootstrap samples of the bootstrap sample-specific estimators. Since different bootstrap samples typically result in different regression fits, the bagged estimator is typically a very high dimensional regression. Two applications of bagging are provided in random forest and linear regression [6]. In random forests, the bagged regression estimator is defined as an average of bootstrap-specific classification and regression tree (CART) estimators [8], where in each bootstrap sample CART is applied without cross-validation to obtain a fine partitioning. In the linear regression context, Breiman [3] proposed a bagged estimator as the average of bootstrap-specific cross-validated regression estimators, such as a linear regression estimator using forward selection and cross-validation to data adaptively select the size of the model. To conclude, in the current literature on bagging one either aggregates over-fitted regression estimators or one aggregates cross-validated regression estimators.

We note that the latter type of cross-validation within a bootstrap sample provides the right trade-off between bias and variance for the single sample estimator applied to the bootstrap sample. However, since the bagging operation typically reduces the variance and increases bias, it can result in the wrong trade-off for the corresponding bootstrap aggregated estimators. In this article we propose a cross-validated bagged estimator which (1) acknowledges that each estimator indexed by fine-tuning parameters corresponds with a bagged estimator, and (2) uses (external) cross-validation to select among these candidate bagged estimators, and possibly between these estimators and additional non-bagged estimators. By including non-bagged estimators in the set of candidate estimators, this procedure data adaptively selects between bagged and non-bagged estimators, which can be useful in cases where it is unclear if bagging actually improves prediction performance; some of these are discussed in our overview of the bagging literature below.

In order to assess the performance of the cross-validated bagged estimator proposed in this paper, a second level of cross validation is performed. The cross-validated bagged estimator is applied to a learning sample, and its fit is evaluated on a test sample, across different splits of the data into learning and test samples. Performance assessment thus involves double cross-validation.

The organization of this article is as follows. Section 2 provides a brief review of bagging In Section 3 we present the proposed cross-validated bagged estimator in the context of the general unified loss-based estimation framework as introduced in [17]. That is, our estimator applies to any parameter which can be represented as the minimizer over the parameter space of an expectation of a loss function of the experimental unit at a parameter value, where we allow the loss function to be indexed by an unknown nuisance parameter. In particular, this general framework allows us to define the cross-validated bagged estimator of a conditional density, conditional hazard, or conditional location parameter (e.g., mean, median), based on censored and uncensored data. For example, our framework includes prediction of a survival time when the survival time is subject to right censoring. In Section 4 we compare our proposed cross-validated bagged estimator with the bagged cross-validated and bagged non-cross-validated estimators proposed by Breiman.

2. Brief review of bagging

Bagging, or bootstrap aggregating, was introduced by Breiman [3] as a tool for reducing the variance of a predictor. The general idea is to generate multiple versions of a predictor and then use these to get an aggregated predictor. The multiple predictors are obtained by using bootstrap replicates of the data, and bagging is meant to yield gains in accuracy. Whether or not bagging will improve accuracy is related to the stability of the procedure that constructs each predictor [3]. Breiman [4] studied instability and stated that k -nearest neighbor methods are *stable*, but that neural networks, classification and regression trees, and subset selection in linear regression are *unstable* methods. Breiman [3] found that bagging works well for unstable methods.

Several approaches have been offered to combine different classifiers [20,5,15]. In addition, the following modifications of bagging have been proposed: “nice” bagging [25], sub-bagging or sub-sample aggregating [10], and iterated bagging or de-biasing [7]. We provide a brief overview of the various properties of bagging that have been studied and the application of bagging to available algorithms in the literature.

Friedman and Hall [13] show that bagging reduces variability when applied to highly nonlinear estimators such as decision trees and neural networks, and can also reduce bias for certain types of estimators. Breiman [7] shows that iterated bagging is effective in reducing both bias and variance. Buja and Stuetzle [9] look at bagging statistical functionals and U -statistics and apply bagging to CART [8]. They find that in the case of bagging CART, both variance and bias can be reduced. Bühlmann and Yu [10] define the notion of instability and analyze how bagging reduces variance in hard decision problems. Because hard decisions create instability, bagging is helpful to smooth these out, yielding smaller variance and mean squared error. Bühlmann and Yu also look at the bagging effect on piecewise linear spline functions in multivariate adaptive regression splines (MARS) [12] and find that bagging is unnecessary for MARS. Borra and Ciaccio [2] also apply bagging to MARS, as well as to project pursuit regression (PPR) and local learning based on recursive covering (DART), and note that in most cases bagging reduces the variability of these methods. Bagging has been viewed from its ability to reduce instability [10] and its success with nonlinear features of statistical methods [13,9]. Hall and Samworth [14] address how performance depends on re-sample size.

Skurichina and Duin [25] offer several conclusions about bagging for linear classifiers; these include that bagging is not necessarily a stabilizing technique where stabilization is defined for linear classifiers, that the number of bootstrap replicates should be limited, that the usefulness of bagging can be determined by the instability of the classifier, and that bagging improves the performance of a classifier when the classifier is unstable. For computational considerations, it is helpful to have a sense of how many bootstrap replicates are adequate. Breiman [3] looked at 10–100 replicates and suggested that fewer replicates are required when the outcome is numerical.

3. The cross-validated bagged learner

Suppose that one observes a sample of n i.i.d. observations on a random variable O with data generating distribution P_0 , which is known to be an element of a model \mathcal{M} . Let $\psi_0 = \Psi(P_0)$ be the parameter of interest of the data generating distribution P_0 . We assume that the true parameter (value) ψ_0 can be defined in terms of a *loss function*, $(O, \psi) \rightarrow L(O, \psi)$, as the minimizer of the expected loss, or *risk*. That is, $\psi_0 = \Psi(P_0) = \arg \min_{\psi \in \Psi} \int L(o, \psi) dP_0(o)$, where the minimum is taken over the parameter space $\Psi \equiv \{\Psi(P) : P \in \mathcal{M}\}$.

In regression with a continuous outcome, a common loss function is the squared error loss, $L(O = (Y, W), \psi) = (Y - \psi(W))^2$, corresponding to the conditional mean $\psi_0(W) = E_0[Y | W]$, and if $\psi_0 = dP_0/d\mu$ is the actual density of O or a sub-vector of O , then $L(O, \psi) = -\log \psi(O)$. As in the unified loss-based estimation approach presented in [17], it is allowed that the loss function depends on a nuisance parameter $\Upsilon(P_0)$: that is $L(O, \psi) = L(O, \psi | \Upsilon(P_0))$. By allowing such unknown loss functions, this framework includes most parameters. For the sake of notational convenience, we suppress the possible dependence of the loss function on a nuisance parameter in the notation, but we will point out at the appropriate places how this affects the proposed estimation procedure.

Let P_n denote the empirical probability distribution of the sample O_1, \dots, O_n , which puts mass $1/n$ on each observation. Consider now a collection of candidate estimators $P_n \rightarrow \hat{\Psi}_s(P_n)$ indexed by a fine-tuning parameter s ranging over a set \mathcal{A}_n . For example, if $\psi_0(W) = E(Y | W)$, then $\hat{\Psi}_s(P_n)$ might represent a particular learning algorithm for estimation of $E(Y | W)$ indexed by fine-tuning parameters s which are user supplied, such as a support-vector machine algorithm, a forward step-wise algorithm, logic regression [22], the D/S/A-polynomial regression algorithm [23], and so on. Another class of general examples is obtained by defining $\hat{\Psi}_s(P_n) \equiv \arg \min_{\psi \in \Psi_s} \sum_{i=1}^n L(O_i, \psi)$ as the minimizer of the empirical risk $\sum_i L(O_i, \psi)$ over a sub-parameter space $\Psi_s \subset \Psi$ indexed by s , given a collection of subspaces $\Psi_s, s \in \mathcal{A}_n$.

In the case that the loss function depends on a nuisance parameter $\Upsilon(P_0)$, one would estimate the nuisance parameter with an estimator $\hat{\Upsilon}(P_n)$, and minimize the empirical risk corresponding with the estimated loss function. Most estimators can be considered as approximate minimizers of the empirical risk, indexed by parameters defining the search algorithm, such as the space the algorithm searches over and the depth to which it searches the space. We note that we view the estimators as mappings $\hat{\Psi}_s$ from data, P_n , to the parameter space.

Given the empirical distribution P_n , let $P_n^\#$ denote the empirical distribution of a sample of n i.i.d. observations $O_1^\#, \dots, O_n^\#$ from P_n . Given the s -specific estimator $\hat{\Psi}_s$ we can define a corresponding s -specific bagged estimator as $\tilde{\Psi}_s(P_n) \equiv E(\hat{\Psi}_s(P_n^\#) | P_n)$. To evaluate the conditional expectation, given the data P_n , one needs to draw many bootstrap samples $P_n^\#$ from the empirical probability distribution P_n . For each of these draws of size n , $P_{n1}^\#, \dots, P_{nB}^\#$, the estimators $\hat{\Psi}_s(P_{nb}^\#)$ based on the bootstrap sample $P_{nb}^\#, b = 1, \dots, B$, are calculated and averaged: $\tilde{\Psi}_s(P_n) = \lim_{B \rightarrow \infty} \frac{1}{B} \sum_{b=1}^B \hat{\Psi}_s(P_{nb}^\#)$.

This results in a set of candidate bagged estimators $\tilde{\Psi}_s(P_n)$ indexed by s . Our goal is to data adaptively select the s which minimizes the risk of $\tilde{\Psi}_s(P_n)$ over \mathcal{A}_n . For this purpose, we propose the *cross-validated bagged estimator* defined as: $\tilde{\Psi}(P_n) = \tilde{\Psi}_{\hat{S}(P_n)}(P_n)$ where $\hat{S}(P_n)$ is the cross-validation selector based on the loss function $L(\cdot, \cdot)$ corresponding to a cross-validation scheme defined by a random n dimensional vector $B_n \in \{0, 1\}^n$. A realization of $B_n = (B_{n,1}, \dots, B_{n,n})$ defines a particular split of the learning sample of n observations into a training set, $\{i \in \{1, \dots, n\} : B_{n,i} = 0\}$, and a validation set, $\{i \in \{1, \dots, n\} : B_{n,i} = 1\}$. We will denote the proportion of observations in the validation set with p . The empirical distributions of the training and validation sets are denoted by P_{n,B_n}^0 and P_{n,B_n}^1 , respectively. Formally, the cross-validation selector $\hat{S}(P_n)$ is defined as

$$\begin{aligned} \hat{S}(P_n) &= \arg \min_{s \in \mathcal{A}_n} E_{B_n} P_{n,B_n}^1 L(\cdot, \tilde{\Psi}_s(P_{n,B_n}^0)) \\ &= \arg \min_{s \in \mathcal{A}_n} E_{B_n} \sum_{i, B_n(i)=1} L(O_i, \tilde{\Psi}_s(P_{n,B_n}^0)). \end{aligned} \tag{1}$$

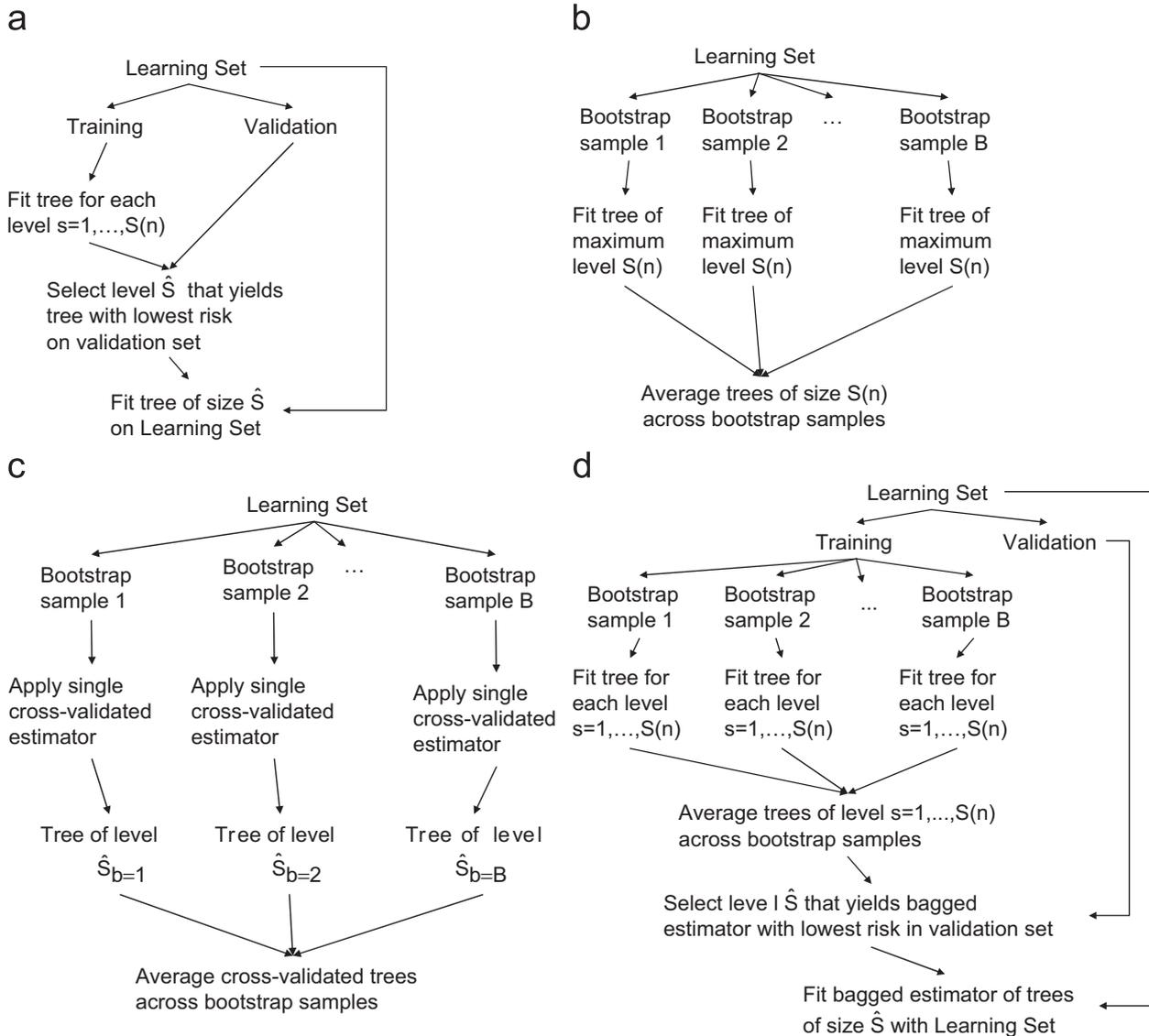


Fig. 1. Schematic contrasting alternative estimators, as applied in Simulations 1–3 using CART; (a) single (non-bagged) cross validated estimator; (b) bagged non-cross-validated estimator; (c) bagged cross-validated estimator; (d) cross-validated bagged estimator.

At the first equality we used the notation $Pf \equiv \int f(o) dP(o)$. If the loss function depends on an unknown nuisance parameter, then one estimates the unknown loss function on the training sample: that is, one replaces the loss function in (1) by $L(\cdot, \tilde{\Psi}_s(P_{n,B_n}^0) | \hat{\Upsilon}(P_{n,B_n}^0))$.

To calculate this selector of s , for each possible realization of the sample split B_n and index $s \in \mathcal{A}_n$, B bootstrap samples $P_{n,B_n,b}^{0\#}$ of size $n(1-p)$ are drawn from the training sample P_{n,B_n}^0 , $b = 1, \dots, B$. For each of these B bootstrap samples we compute the corresponding s -specific estimators $\hat{\Psi}_s(P_{n,B_n,b}^{0\#})$ and average them to obtain: $\tilde{\Psi}_s(P_{n,B_n}^0) = \frac{1}{B} \sum_{b=1}^B \hat{\Psi}_s(P_{n,B_n,b}^{0\#})$. The empirical risk of this estimator over the validation sample can now be computed, and averaged over the different splits B_n , as in (1), which results in the so-called cross-validated risk of the estimator $\tilde{\Psi}_s(P_n)$. The cross-validation selector is defined as the one which minimizes this cross-validated risk over $s \in \mathcal{A}_n$. A schematic illustrating the implementation of the cross-validated bagged learner, applied using classification and regression trees, as employed in the simulations reported in Section 4, is shown in Fig. 1d.

3.1. Cross-validated bagged learning versus bagged cross-validated learning

It is of interest to contrast the proposed cross-validated bagged estimator to the bagged cross-validated estimator as used in [3] and followed by other authors. In the bagged cross-validated approach, the selection of s via cross-validation is performed within each bootstrap sample. Subsequently the B bootstrap-specific estimators are averaged to arrive at the final estimator. Formally, within a bootstrap sample $P_n^\#$ the cross-validation selector of s can be defined as

$$\hat{S}_{br}(P_n^\#) = \arg \min_{s \in \mathcal{A}_n} E_{B_n} P_{n, B_n}^{1, \#} L(\cdot, \hat{\Psi}_s(P_{n, B_n}^{0, \#})).$$

The corresponding estimator based on a single bootstrap sample $P_n^\#$ is thus defined as $\hat{\Psi}_{CV}(P_n^\#) = \hat{\Psi}_{\hat{S}_{br}(P_n^\#)}(P_n^\#)$. Finally, the corresponding bagged cross-validated estimator is the average over a large collection of bootstrap-specific estimators: $\tilde{\Psi}_{br}(P_n) = E(\hat{\Psi}_{CV}(P_n^\#) | P_n)$. Fig. 1 illustrates implementation of the bagged cross-validated estimator, and contrasts it with the cross-validated bagged estimator, as well as with a single cross-validated estimator, and a bagged non-cross-validated estimator.

Using cross-validation within a bootstrap sample provides the right trade-off between bias and variance among the estimators $\hat{\Psi}_s(P_n^\#)$, $s \in \mathcal{A}_n$. However, one would expect it not to perform the right trade-off between bias and variance for the bagged estimators $\tilde{\Psi}_s(P_n)$. The reason for this is that the bagged estimator should be less variable as a result of the averaging, and might be more biased due to the double sampling. An increase in bias is due to two (probably cumulative) sources: first, the bias introduced by applying an estimator to a bootstrap sample relative to the empirical sample; and second, the bias introduced by applying the estimator to the empirical sample relative to the truth. In general, the estimators $\hat{\Psi}_s$, $s \in \mathcal{A}_n$, and $\tilde{\Psi}_s$, $s \in \mathcal{A}_n$, are very different classes of estimators. Thus a good selector among the un-bagged estimators is not necessarily a good selector among the corresponding bagged estimators.

3.2. Performance of the cross-validation selector

Let $d(\psi, \psi_0) = \int L(o, \psi) dP_0(o) - \int L(o, \psi_0) dP_0(o)$ denote the risk dissimilarity between a candidate ψ and the true ψ_0 implied by the loss function $L(\cdot, \cdot)$. Results on the cross-validation selector (see [18,19]) state that if the loss function is uniformly bounded in its arguments, then the difference of the risk dissimilarity of the cross-validated selected bagged estimator and the risk dissimilarity of the oracle selected bagged estimator is of the order $\log K(n)/np$ plus possibly a term due to estimation of the nuisance parameter in the loss function. The oracle selected bagged estimator is defined as $\tilde{\Psi}_{\tilde{S}_{n(1-p)}(P_n)}(P_n)$, where $\tilde{S}_{n(1-p)}(P_n) = \arg \min_s E_{B_n} \int L(o, \tilde{\Psi}_s(P_{n, B_n}^0)) dP_0(o)$. Thus for a given data set the oracle selector $\tilde{S}_{n(1-p)}(P_n)$ selects the bagged estimator (based on $n(1-p)$ observations) closest to the truth with respect to the risk dissimilarity.

These results only rely on the loss function to be uniformly bounded in the support of O and the parameter space. They imply that if the number of candidate estimators is polynomial in sample size (and, in the case that the loss function is unknown, that it can be estimated at a better rate than the convergence rate of the oracle selected estimator), then either the cross-validated selected estimator is asymptotically equivalent (up to the constant) to the oracle selected estimator, or it achieves the essentially parametric rate of convergence $\log n/n$.

3.3. Cross-validation selection of the degree of bagging

In cases where there is a concern that the bagging operation might actually worsen the estimator it is a good idea to let cross-validation select between the original un-bagged estimator and the bagged estimator. In general, the following method might be of interest. Define $\tilde{\Psi}_{s,\alpha} = \alpha\tilde{\Psi}_s + (1-\alpha)\hat{\Psi}_s$ as the weighted average between the bagged and un-bagged estimator, $\alpha \in [0, 1]$, and use cross-validation to select (s, α) . In this manner, the data are used to decide to what degree α the bagging operation should be used, and, by our results establishing asymptotic equivalence with the oracle selector of (α, s) , our cross-validated selected estimator will perform asymptotically at least as well as the non-bagged estimator and bagged estimator.

3.4. Assessing the performance of the cross-validated bagged estimator

One can estimate the risk $\int L(o, \tilde{\Psi}(P_n)) dP_0(o)$ of the cross-validated bagged estimator $\tilde{\Psi}(P_n)$ with the cross-validated risk of the estimator $\tilde{\Psi}(P_n) = \tilde{\Psi}_{\hat{S}(P_n)}(P_n)$ using a cross-validation scheme defined by a random vector $B_n^* \in \{0, 1\}^n$: $E_{B_n^*} \sum_{i, B_n^*(i)=1} L(O_i, \tilde{\Psi}(P_{n, B_n^*}^0))$. This procedure would require carrying out a B_n -specific cross-validation scheme within each learning sample $P_{n, B_n^*}^0$, which is often referred to as double cross-validation.

4. Simulations

In this section we illustrate the proposed cross-validated bagged estimation methodology in comparison to both the bagged cross-validated and bagged non-cross-validated estimators suggested by Breiman, as well as in comparison to the non-bagged estimator. The implementation of each of the estimators is illustrated schematically in Fig. 1. An extended description of these and additional simulation results are provided in an on-line technical report [21].

The models chosen for these simulations are the same as those implemented in both [3,12]. For each simulation, we evaluated four estimators: bagged cross-validated, bagged non-cross-validated, cross-validated bagged, and a single CART tree. For each, the same learning set P_n of size 200 was used for building and choosing the best predictor and the same independent test set of size 1000 was used to evaluate the fit and reflects the empirical risk estimated in the corresponding tables.

For the single tree, P_n was split into a training set P_n^0 and validation set P_n^1 and the training set was used to fit trees $\hat{\Psi}_s(P_n^0)$, where $s = (1, \dots, S(n))$, $minbucket = 7$, $cp = .01$. The best level of the tree was chosen via 10-fold cross-validation based on the validation set, and the learning set was then used to fit a tree of that level, giving a final cross-validated tree estimator $\hat{\Psi}_{CV}(P_n)$.

For the bagged cross-validated estimator, the learning set P_n was used to generate $B = 100$ bootstrap samples and build B bootstrap-specific estimators $\hat{\Psi}_{CV}(P_n^\#)$, where 10-fold cross validation was implemented to select the best tree level $\hat{S}_{br}(P_n^\#)$ for each bootstrap sample. The final estimator $\tilde{\Psi}_{brCV}(P_n)$ was the average of the B bootstrap-specific cross-validated estimators.

For the bagged non-cross-validated estimator, the learning set P_n was used to generate $B = 100$ bootstrap samples and build B bootstrap-specific estimators $\hat{\Psi}_{S(n)}(P_n^\#)$, each of maximum tree size $S(n)$. The final estimator $\tilde{\Psi}_{brFT}(P_n)$ was the average of the B bootstrap-specific estimators.

For our proposed cross-validated bagging scheme, P_n was split into a training set P_n^0 and validation set P_n^1 , and the training set was used to generate $B = 100$ bootstrap samples. Bootstrap-specific trees $\hat{\Psi}_s(P_n^{0,\#})$ of each level $s = 1, \dots, S(n)$ were averaged across

Table 1
Simulation 1

Estimator	Bootstrap samples	Emp. risk mean	Emp. risk std. dev.	% Improvement in risk
Single cross-valid. tree ($\hat{\Psi}_{CV}$)	0	20.83	3.07	
Bagged cross-valid. ($\tilde{\Psi}_{brCV}$)	100	15.63	1.37	25
Bagged non-cross-valid. ($\tilde{\Psi}_{brFT}$)	100	13.71	0.95	34
Cross-valid. bagged ($\tilde{\Psi}$)	100	13.75	1.01	34

Table 2
Simulation 2

Estimator	Bootstrap samples	Emp. risk mean	Emp. risk std. dev.	% Improvement in risk
Single cross-valid. tree ($\hat{\Psi}_{CV}$)	0	15055.41	4709.78	
Bagged cross-valid. ($\tilde{\Psi}_{brCV}$)	100	6283.21	1177.67	58
Bagged non-cross-valid. ($\tilde{\Psi}_{brFT}$)	100	4883.19	978.73	67
Cross-valid. bagged ($\tilde{\Psi}$)	100	4871.83	964.94	68

bootstrap samples, yielding a vector of bagged estimators based on the training set and indexed by s . Ten-fold cross-validation was used to select the best tree level for the bagged estimator ($\hat{S}(P_n)$) based on the independent validation sample. The resulting cross-validated bagged estimator was then applied to the learning sample $\tilde{\Psi}_{\hat{S}(P_n)}(P_n) = \tilde{\Psi}(P_n)$.

This entire procedure was repeated 100 times and the empirical risk averaged over the repetitions. All simulations were performed in the statistical package R [16] using the recursive partitioning algorithm `rpart` [26] for CART and the `mlbench` package for data generation.

4.1. Simulation 1

The first simulated data set, as described in Breiman [3], has 10 independent predictor variables x_1, \dots, x_{10} each of which is uniformly distributed over $(0, 1)$. The response is given by $y = 10 \sin(\pi x_1 x_2) + 20(x_3 - .5)^2 + 10x_4 + 5x_5 + error$, where the $error \sim N(0, 1)$. The results are shown in Table 1.

4.2. Simulation 2

In the second simulated data set, as described in Breiman [3], there are four independent predictor variables x_1, \dots, x_4 each of which is uniformly distributed over different ranges: $0 \leq x_1 \leq 100, 20 \leq (x_2/2\pi) \leq 280, 0 \leq x_3 \leq 1, 1 \leq x_4 \leq 11$. The response is given by $y = (x_1^2 + (x_2 x_3 - (1/x_2 x_4))^2)^{1/2} + error$, where the $error \sim N(0, .62)$. The results are shown in Table 2.

Table 3
Simulation 3

Estimator	Bootstrap samples	Emp. risk mean	Emp. risk std. dev.	% Improvement in risk
Single cross-valid. tree ($\hat{\Psi}_{CV}$)	0	0.8377	0.0365	
Bagged cross-valid. ($\tilde{\Psi}_{brCV}$)	100	0.8005	0.0354	4
Bagged non-cross-valid. ($\tilde{\Psi}_{brFF}$)	100	0.8217	0.0391	2
Cross-valid. bagged ($\tilde{\Psi}$)	100	0.7951	0.0357	5

4.3. Simulation 3

In the third simulated data set, as described in Breiman [3], there are four independent predictor variables x_1, \dots, x_4 each of which is uniformly distributed over different ranges: $0 \leq x_1 \leq 100$, $20 \leq (x_2/2\pi) \leq 280$, $0 \leq x_3 \leq 1$, $1 \leq x_4 \leq 11$. The response is given by $y = \arctan(\frac{x_2x_3 - (1/x_2x_4)}{x_1}) + error$, where the $error \sim N(0, .86)$. The results are shown in Table 3.

4.4. Summary

The risks of the estimators generated by each bagging approach (bagged cross-validated, cross-validated bagged, and bagged non-cross-validated) were compared with the risk of the single cross-validated tree. In all simulations, all three bagged estimators resulted in a decreased risk as compared to the single tree. The magnitude of this decrease provides a measure of performance on which the three bagged estimators can be compared. Interestingly, in every case, the cross-validated bagged approach proposed in this paper resulted in a greater decrease in risk than did the bagged cross-validated approach. This result agrees with the theoretical argument, presented in Section 3, that the cross-validated bagged approach appropriately trades off bias and variance for the bagged estimator.

The improvement in risk resulting from cross-validating the bagged estimator, as compared to simply using the bagged estimator with the maximum tree level, is less dramatic; however, cross-validation does appear to provide some small benefit in terms of decreased risk. In comparing the cross-validated and non-cross-validated bagged estimators, it is interesting to note that the fine-tuning parameter (tree level) selected by cross-validation is based on bagged estimators fit using the training sample. Because the training sample is a subset of the learning sample, one expects that the estimator based on the training sample will be more biased than an equivalent estimator based on the entire learning sample (for which cross-validation was performed with a completely independent data set). In contrast, the non-cross-validated estimator fit on the learning sample is generally expected to be over-fit. Thus, the truth is expected to lie somewhere between the two estimators. It is interesting to note that, even under the current simulation settings, where lack of cross-validation seems to have little cost in terms of prediction error, the cross-validated bagged estimator performed no worse than the non-cross-validated bagged estimator.

Tables 1–3 also present estimates of the variability of the cross-validated risk. While estimating this variability is straightforward in a simulation setting, estimation of the variability in risk estimate based on a single sample poses a challenge. While beyond the scope of the present article, we point out that Dudoit and van der Laan [11] derive the influence curve for the cross-validated risk, making possible first-order inference for the risk estimate.

5. Discussion

In this article we developed a general class of statistical learning algorithms for parameters that can be defined as a minimizer over the parameter space of the population mean of a loss function. The proposed class combines (1) a given statistical learning algorithm indexed by fine-tuning parameters, (2) bootstrap aggregation, and (3) cross-validation. The result is a high dimensional fit of the parameter of interest that corresponds with a sensible trade-off between bias and variance. In particular, we implemented such an algorithm for prediction based on classification and regression trees. The results support our claim that when cross-validation is combined with bagging, cross-validation should be performed at the level of the bagged estimator itself (cross-validated bagged estimation), rather than prior to bagging at the level of the bootstrap sample (bagged cross-validated estimation).

As discussed by previous studies, the real benefit of bagging occurs if the original statistical learning algorithms are local learning algorithms. The learning algorithm employed in our simulations, CART, is vary variable. In contrast, other learning algorithms rely more on extrapolation; for example, the Deletion/Substitution/Addition algorithm [23] fits the true regression using a linear combination of polynomial basis functions. The optimal trade-off between local learning and global smoothing via extrapolation will depend on the true underlying data-generating distribution. In a recent paper, Sinisi et al. introduce the concept of super learning [24], in which the data are used to decide (via cross-validation) between candidate learning algorithms or convex combinations of such learners. This concept can readily be applied to bagging, by creating convex combinations of global and local learning algorithms, bagging the resulting estimators, and selecting between the resulting candidate estimators (both bagged and un-bagged) using cross-validation.

Acknowledgments

Maya Petersen was supported by a Predoctoral Fellowship from the Howard Hughes Medical Institute. Annette Molinaro was supported by an NCI-funded career award (K22 CA123146). This work was partially supported by NIH RO1 Grant GM071397.

References

- [1] M.D. Birkner, S.E. Sinisi, M.J. van der Laan, Multiple testing and data adaptive regression: an application to HIV-1 sequence data, *Statist. Appl. Genetics Molecular Biol.* 4 (1), article 8, 2005. URL: <http://www.bepress.com/sagmb/vol4/iss1/art8>.
- [2] S. Borra, A.D. Ciaccio, Improving nonparametric regression methods by bagging and boosting, *Comput. Statist. Data Anal.* 38 (2002) 407–420.
- [3] L. Breiman, Bagging predictors, *Mach. Learning* 24 (2) (1996) 123–140.
- [4] L. Breiman, Heuristics of instability and stabilization in model selection, *Ann. Statist.* 24 (1996) 2350–2383.
- [5] L. Breiman, Stacked regressions, *Mach. Learning* 24 (1996) 49–64.
- [6] L. Breiman, Random forests, *Mach. Learning* 45 (2001) 5–32.
- [7] L. Breiman, Using iterated bagging to debias regressions, *Mach. Learning* 45 (2001) 261–277.
- [8] L. Breiman, J.H. Friedman, R.A. Olshen, C.J. Stone, *Classification and Regression Trees*, The Wadsworth Statistics/Probability Series, Wadsworth International Group, Belmont, CA, 1984.
- [9] A. Buja, W. Stuetzle, Observations on bagging, 2002. URL: <http://www-stat.wharton.upenn.edu/~buja/PAPERS/paper-bag-wxs.pdf>.
- [10] P. Bühlmann, B. Yu, Analyzing bagging, *Ann. Statist.* 30 (2002) 927–961.
- [11] S. Dudoit, M. van der Laan, Asymptotics of cross-validated risk estimation in estimator selection and performance assessment, *Statist. Methodology* 2 (2) (2005) 131–154.

- [12] J.H. Friedman, Multivariate adaptive regression splines, *Ann. Statist.* 19 (1) (1991) 1–141 discussion by A. R. Barron and X. Xiao.
- [13] J.H. Friedman, P. Hall, On bagging and nonlinear estimation, 2000. URL: (<http://www-stat.stanford.edu/~jhf/ftp/bag.ps>).
- [14] P. Hall, R.J. Samworth, Properties of bagged nearest-neighbour classifiers, *J. Roy. Statist. Soc. Ser. B* 67 (3) (2005) 363–379.
- [15] T. Hothorn, B. Lausen, Bundling classifiers by bagging trees, *Comput. Statist. Data Anal.* 49 (4) (2005) 1068–1078.
- [16] R. Ihaka, R. Gentleman, R: a language for data analysis and graphics, *J. Comput. Graphical Statist.* 5 (1996) 299–314.
- [17] M.J. van der Laan, S. Dudoit, Unified cross-validation methodology for selection among estimators and a general cross-validated adaptive epsilon-net estimator: finite sample oracle inequalities and examples, Technical Report 130, Division of Biostatistics, University of California, Berkeley, November 2003, URL: (www.bepress.com/ucbbiostat/paper130/).
- [18] M.J. van der Laan, S. Dudoit, S. Keleş, Asymptotic optimality of likelihood based cross-validation, Technical Report 125, Division of Biostatistics, University of California, Berkeley, February 2003. URL: (www.bepress.com/ucbbiostat/paper125/).
- [19] M.J. van der Laan, S. Dudoit, A.W. van der Vaart, The cross-validated adaptive epsilon-net estimator, Technical Report 142, Division of Biostatistics, University of California, Berkeley, February 2004. URL: (www.bepress.com/ucbbiostat/paper142/).
- [20] M. LeBlanc, R. Tibshirani, Combining estimates in regression and classification, *J. Amer. Statist. Assoc.* 91 (1996) 1641–1650.
- [21] A.M. Molinaro, M.J. van der Laan, Cross-validating and bagging partitioning algorithms with variable importance, Technical Report 185, Division of Biostatistics, University of California, Berkeley, November 2005. URL: (www.bepress.com/ucbbiostat/paper185/).
- [22] I. Ruczinski, C. Kooperberg, M. LeBlanc, Logic regression, *J. Comput. Graphical Statist.* 12 (3) (2003) 475–511 URL: (<http://www.biostat.jhsph.edu/~iruczins/publications/publications.html>).
- [23] S.E. Sinisi, M.J. van der Laan, Deletion/Substitution/Addition algorithm in learning with applications in genomics, *Statist. Appl. Genetics and Molecular Biology* 3 (1), article 18, 2004. URL: (<http://www.bepress.com/sagmb/vol3/iss1/art18>).
- [24] S.E. Sinisi, E.C. Polley, M.L. Petersen, S.-Y. Rhee, M.J. van der Laan, Super learning: an application to the prediction of HIV-1 drug resistance, *Statist. Appl. Genetics* 6 (1), article 7. URL: (<http://www.bepress.com/sagmb/vol6/iss1/art7>).
- [25] M. Skurichina, R.P.W. Duin, Bagging for linear classifiers, *Pattern Recognition* 31 (1998) 909–930.
- [26] T. Therneau, E. Atkinson, An introduction to recursive partitioning using the rpart routines, Technical Report 61, Division of Biostatistics, Mayo Clinic, Rochester, November 1997. URL: (www.mayo.edu/hsr/techrpt/61.pdf).