# Measuring the Algorithmic Convergence of Randomized Ensembles: The Regression Setting[*]

Miles E. Lopes[†], Suofei Wu[†], and Thomas C. M. Lee[†]

**Abstract.** When randomized ensemble methods such as bagging and random forests are implemented, a basic question arises: Is the ensemble large enough? In particular, the practitioner desires a rigorous guarantee that a given ensemble will perform nearly as well as an ideal infinite ensemble (trained on the same data). The purpose of the current paper is to develop a bootstrap method for solving this problem in the context of regression—which complements our companion paper in the context of classification [Lopes, *Ann. Statist.*, 47 (2019), 1088–1112]. In contrast to the classification setting, the current paper shows that theoretical guarantees for the proposed bootstrap can be established under much weaker assumptions. In addition, we illustrate the flexibility of the method by showing how it can be adapted to measure algorithmic convergence for variable selection. Lastly, we provide numerical results demonstrating that the method works well in a range of situations.

**1. Introduction.** Ensemble methods are a fundamental approach to prediction, based on the principle that accuracy can be enhanced by aggregating a diverse collection of prediction functions. Two of the most widely used methods in this class are *random forests* and *bagging*, which rely on randomization as a general way to diversify an ensemble [10, 11]. For these types of randomized ensembles, it is generally understood that the predictive accuracy improves and eventually stabilizes as the ensemble size becomes large. Likewise, in the theoretical analysis of randomized ensembles, it is common to focus on the ideal case of an infinite ensemble [14, 28, 6, 50, 5, 55]. However, in practice, the user does not know the true relationship between accuracy and ensemble size. As a result, it is difficult to know if a given ensemble is large enough so that its accuracy will nearly match the ideal level of an infinite ensemble.

Beyond these statistical considerations, the relationship between accuracy and ensemble size is important for computational reasons. Indeed, as an ensemble becomes larger, more resources are needed to train it, to store it in memory, and to make new predictions on unlabeled points—especially when large volumes of data are involved. Consequently, if it were possible for the user to know the true relationship between accuracy and ensemble size, it would be possible to do "just enough" computation to achieve a desired degree of convergence.

[†]Department of Statistics, University of California, Davis, Davis, CA 95616 USA (melopes@ucdavis.edu, swu@ucdavis.edu, tcmlee@ucdavis.edu).

Similarly, this would also make it possible to ensure that the amount of computation is *adaptive* to the unique data the user has at hand.

The purpose of the current paper is develop a solution to the problem of measuring algorithmic convergence for random forests, bagging, and related methods in the context of regression. More specifically, we offer a bootstrap method for estimating how far the prediction error of a finite ensemble is from the ideal prediction error of an infinite ensemble (trained on the same data). To put this into perspective for the setting of regression, it is worth noting that a theoretically justified method for solving this problem has not previously been available. In this way, our work fills a significant gap in the literature by providing users with a more rigorous alternative to informal rules that are used in practice for selecting ensemble size. Furthermore, our approach is of broader conceptual interest, because it indicates new possibilities for applying bootstrap methods to randomized algorithms outside the scope of classical statistical inference (see subsection 1.3 for additional details).

In the remainder of the introduction, we give a precise description of the problem formulation in subsection 1.1, followed by a summary of related work and contributions in subsection 1.2 and subsection 1.3.

**1.1. Background and setup.** To fix some basic notation for the regression setting, let $\mathcal{D} = \{(X_j, Y_j)\}_{j=1}^n$ denote a set of training data in a space $\mathcal{X} \times \mathbb{R}$, where each $Y_j$ is the scalar response variable associated to $X_j$, and the space $\mathcal{X}$ is arbitrary. In addition, for each $i = 1, \ldots, t$, we write $T_i : \mathcal{X} \to \mathbb{R}$ to refer to the $i$th regression function in an ensemble of size $t$ trained on $\mathcal{D}$.

*Randomized regression ensembles.* For the purpose of understanding our setup, it is helpful to quickly review the methods of bagging and random forests. The method of bagging works by generating random sets $\mathcal{D}_1^*, \ldots, \mathcal{D}_t^*$, each of size $n$, by sampling with replacement from $\mathcal{D}$. Next, a standard "base" regression algorithm is used to train a regression function $T_i$ on $\mathcal{D}_i^*$ for each $i = 1, \ldots, t$. For instance, it is especially common to apply a decision tree algorithm like CART [12] to each set $\mathcal{D}_i^*$. In turn, future predictions are made by using the averaged regression function, which is defined for each $x \in \mathcal{X}$ by

$$(1.1) \qquad \bar{T}_t(x) = \frac{1}{t} \sum_{i=1}^t T_i(x).$$

Much like bagging, the method of random forests uses sampling with replacement to generate the same type of random sets $\mathcal{D}_1^*, \ldots, \mathcal{D}_t^*$. However, random forests adds an additional source of randomness. Namely, if the space $\mathcal{X}$ is (say) $p$-dimensional, and CART is the base regression algorithm, then random forests uses randomly chosen subsets of the $p$ features when "split points" are selected for the CART regression trees. Apart from this distinction, random forests also uses the average (1.1) when making final predictions. A more detailed description may be found in [22].

In order to unify the methods of bagging and random forests within a common theoretical framework, our analysis will consider a more general class of randomized ensembles. This class consists of regression functions $T_1, \ldots, T_t$ that can be represented in the abstract form

$$(1.2) \qquad T_i(x) = \varphi(x; \mathcal{D}, \xi_i),$$

where $\xi_1, \ldots, \xi_t$ are i.i.d. "randomizing parameters" generated independently of $\mathcal{D}$, and $\varphi$ is a deterministic function that does not depend on $n$ or $t$. In particular, the representation (1.2) implies that the random functions $T_1, \ldots, T_t$ are conditionally i.i.d., given $\mathcal{D}$. To see why bagging is representable in this form, note that $\xi_i$ can be viewed as a random vector that specifies which points in $\mathcal{D}$ are randomly sampled into $\mathcal{D}_i^*$. Similarly, in the case of random forests, each $\xi_i$ encodes the points in $\mathcal{D}_i^*$, as well as randomly chosen sets of features used for training $T_i$. More generally, the representation (1.2) is relevant to other types of randomized ensembles, such as those based on random rotations [9], random projections [15], and posterior sampling [46, 17].

*Algorithmic convergence.* In our analysis of algorithmic convergence, we will focus on quantifying how the mean-squared error (MSE) of an ensemble behaves as the ensemble size $t$ becomes large. To define this measure of error in more precise terms, let $\boldsymbol{\xi}_t := (\xi_1, \ldots, \xi_t)$ denote the randomizing parameters of the ensemble, and let $\nu = \mathcal{L}(X, Y)$ denote the joint distribution of a test point $(X, Y) \in \mathcal{X} \times \mathbb{R}$, which is drawn independently of $\mathcal{D}$ and $\boldsymbol{\xi}_t$. Accordingly, we define

$$(1.3) \qquad \text{MSE}_t = \int_{\mathcal{X} \times \mathbb{R}} \big(y - \bar{T}_t(x)\big)^2 d\nu(x, y) = \mathbb{E}\Big[(Y - \bar{T}_t(X))^2 \,\Big|\, \boldsymbol{\xi}_t, \mathcal{D}\Big],$$
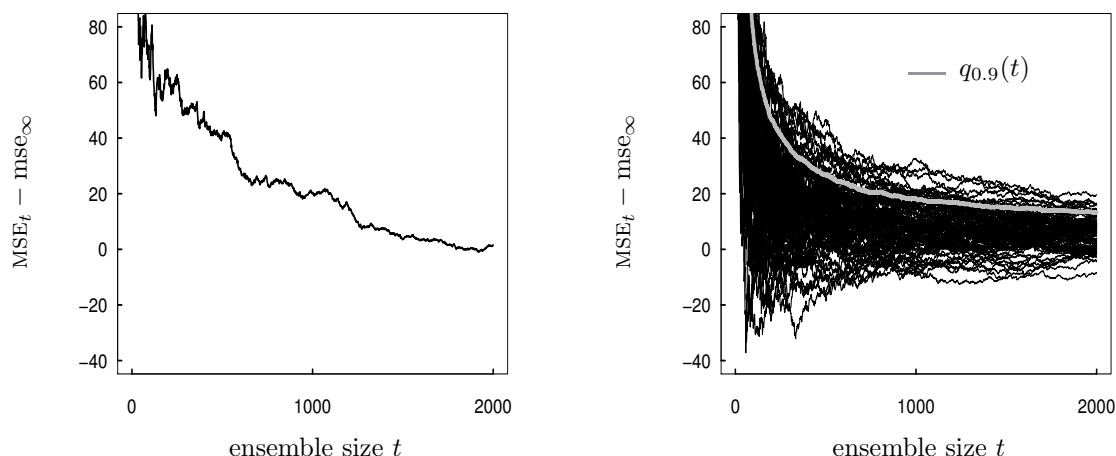
where the expectation on the right is only over the test point $(X, Y)$. In this definition, it is important to notice that $\text{MSE}_t$ is a random variable that depends on both $\boldsymbol{\xi}_t$ and $\mathcal{D}$. However, due to the fact that the *algorithmic* fluctuations of $\text{MSE}_t$ arise only from $\boldsymbol{\xi}_t$, we will view the set $\mathcal{D}$ as a fixed input to the training algorithm, and likewise, our analysis will always be conditional on $\mathcal{D}$. Indeed, the conditioning on $\mathcal{D}$ is motivated by the fact that the user would like to assess convergence for the particular set $\mathcal{D}$ that they actually have, and this viewpoint has been adopted in several other analyses of algorithmic convergence for randomized ensembles [46, 35, 53, 15, 36].

As a conceptual illustration, Figure 1 shows what algorithmic convergence looks like when random forests is applied to a fixed training set $\mathcal{D}$. In detail, the left panel displays values of the convergence gap $\text{MSE}_t - \text{mse}_\infty$ as decision trees are added during a single run of random forests, from $t = 1$ up to $t = 2,000$, where $\text{mse}_\infty$ denotes the limit of $\text{MSE}_t$ as $t \to \infty$. If this entire process is repeated by running random forests many more times on the same set $\mathcal{D}$, then the result is a large collection of overlapping sample paths, as shown in the right panel of Figure 1. (Note also that none of these sample paths are observable in practice, and that the figure is given only for illustration.)

From a practical standpoint, the user would like to know the size of the convergence gap $\text{MSE}_t - \text{mse}_\infty$ as a function of $t$. For this purpose, it is useful to consider the $(1 - \alpha)$-quantile of the random variable $\text{MSE}_t - \text{mse}_\infty$, which is defined for any $\alpha \in (0, 1)$ by

$$q_{1-\alpha}(t) = \inf\Big\{ q \in \mathbb{R} \,\Big|\, \mathbb{P}\big(\text{MSE}_t - \text{mse}_\infty \leq q \,\big|\, \mathcal{D}\big) \geq 1 - \alpha \Big\}.$$

In other words, the value $q_{1-\alpha}(t)$ is the *tightest possible* upper bound on the gap that holds with probability at least $1 - \alpha$, conditionally on the set $\mathcal{D}$. This interpretation of $q_{1-\alpha}(t)$ can also be understood from the right panel of Figure 1, where we have plotted $q_{1-\alpha}(t)$ in gray, with $\alpha = 1/10$.

**Figure 1.** *Left panel: A sample path of $\mathrm{MSE}_t - mse_\infty$ over a single run of random forests on the Housing data described in section 5. Right panel: Many sample paths of $\mathrm{MSE}_t - mse_\infty$, with the 90% quantile $q_{.90}(t)$ overlaid in gray. (The curves in these panels are not observable to the user.)*

   *The problem to be solved.* Although it is clear that the quantile $q_{1-\alpha}(t)$ represents a precise measure of algorithmic convergence, this function is unknown in practice. This leads to the problem of estimating $q_{1-\alpha}(t)$, which we propose to solve.

   Beyond the fact that $q_{1-\alpha}(t)$ is unknown, it is also important to keep in mind that estimating $q_{1-\alpha}(t)$ involves some additional constraints. First, the user would like to be able to assess convergence from the output of a *single* run of the ensemble method. However, at first sight, it is not obvious that the output of a single run provides enough information to successfully estimate $q_{1-\alpha}(t)$. Second, the method for estimating $q_{1-\alpha}(t)$ should be computationally inexpensive, so that the cost of checking convergence is manageable in comparison to the cost of training the ensemble itself. Accordingly, we will show that the proposed method is able to handle both of these constraints in section 2 and section 4, respectively.

**1.2. Related work.** The general problem of measuring the algorithmic convergence of randomized ensembles has attracted sustained interest over the past two decades. For instance, there have been numerous empirical studies of algorithmic convergence for both classification and regression (e.g., [31, 2, 52, 48, 49]).

   With regard to the theoretical analysis of convergence, we will now review the existing results for classification and regression separately. In the setting of classification, much of the literature has studied convergence in terms of the misclassification probability for majority voting, denoted $\mathrm{ERR}_t$ (a counterpart of $\mathrm{MSE}_t$), which is viewed as a random variable that depends on $\boldsymbol{\xi}_t$ and $\mathcal{D}$. For this measure of error, the convergence of $\mathbb{E}[\mathrm{ERR}_t|\mathcal{D}]$ and $\mathrm{var}(\mathrm{ERR}_t|\mathcal{D})$ as $t \to \infty$ has been analyzed in the papers [46, 35, 15], which have developed asymptotic formulas for $\mathbb{E}[\mathrm{ERR}_t|\mathcal{D}]$, as well as bounds for $\mathrm{var}(\mathrm{ERR}_t|\mathcal{D})$. Related results for a different measure of error can also be found in [29]. More recently, our companion paper [36] has developed a bootstrap method for measuring the convergence of $\mathrm{ERR}_t$ which is able to circumvent some of the limitations of formula-based results.

In the setting of regression, algorithmic convergence results on $\mathrm{MSE}_t$ are scarce in comparison to those for $\mathrm{ERR}_t$. Instead, much more attention in the regression literature has focused on how the size of $t$ influences the variance of point predictions $\bar{T}_t(x)$, with $x \in \mathcal{X}$ held fixed, e.g., [56, 1, 62, 43, 53]. To the best of our knowledge, the only paper that has analyzed algorithmic convergence in terms of a prediction error measure is [53], which considers the risk $r_t := \mathbb{E}[(\bar{T}_t(X) - \mu(X))^2]$, where $\mu(x) := \mathbb{E}[Y|X = x]$ is the true regression function, and the expectation in the definition of $r_t$ is over all of the objects $(X, \mathcal{D}, \boldsymbol{\xi}_t)$. In particular, the paper [53] develops an elegant theoretical bound on the gap between $r_t$ and the risk of an infinite ensemble, denoted $r_\infty$. Under the assumption of a Gaussian regression model with $\mathcal{X} = [0,1]^p$, this bound has the form

$$(1.4) \qquad r_t - r_\infty \;\leq\; \tfrac{8}{t}\Big(\|\mu\|_\infty^2 + \sigma^2(1 + 4\log(n))\Big),$$

where $\sigma^2 = \mathrm{var}(Y)$, and $\|\mu\|_\infty = \sup_{x \in \mathcal{X}} |\mu(x)|$. However, due to the fact that the parameters $\sigma$ and $\|\mu\|_\infty$ are unknown, and that $\|\mu\|_\infty$ is inherently conservative, this bound does not lend itself to a practical method for measuring convergence, and is primarily of theoretical interest.

### 1.3. Contributions.

*Methodology.* From a methodological standpoint, the approach taken here differs in several ways from previous works in the regression setting. Most notably, our work looks at algorithmic convergence in terms of an error measure that is conditional on $\mathcal{D}$. (For instance, this differs from the analysis of $r_t$ mentioned above, which averages over $\mathcal{D}$.) In more concrete terms, we will provide a quantile estimate $\hat{q}_{1-\alpha}(t)$, such that the bound

$$\mathrm{MSE}_t - \mathrm{mse}_\infty \;\leq\; \hat{q}_{1-\alpha}(t)$$

holds with a probability that is nearly $1-\alpha$ or larger, conditionally on $\mathcal{D}$. This conditioning is especially important from the viewpoint of the user, who is typically interested in algorithmic convergence with respect to the *actual dataset at hand*. Another distinct feature of our method is that it provides the user with a *direct numerical estimate of convergence*, whereas formula-based results are more likely to depend on specialized models, involve conservative constants, or depend on unknown parameters, such as in the bound (1.4).

In addition, the scope of the proposed method goes beyond $\mathrm{MSE}_t$, and in subsection 2.3 we will show how the bootstrap method is flexible enough that it can also be applied to variable selection. In this context, the ensemble provides a ranking of variables according to an "importance measure," and this ranking typically stabilizes as $t \to \infty$. However, the notion of convergence is somewhat subtle, because the importance measure for some variables may converge more slowly than for others—*which can distort the overall ranking of variables.* As far as we know, this issue has not been addressed in the literature, and the method proposed in subsection 2.3 provides a way to check that convergence has been achieved uniformly across variables, so that they can be compared fairly.

*Theory.* From a theoretical standpoint, the most important aspect of our work is that it establishes consistency guarantees for the proposed methods under very mild assumptions. To place our assumptions into context, it should be emphasized that most analyses of randomized ensembles deal with specialized types of prediction functions $T_1, \ldots, T_t$ that are simpler than

the ones used in practice, e.g., [34, 1, 6, 5, 55, 53, 54, 36]. By contrast, our current results for regression only rely on (1.2) and basic moment assumptions (to be detailed in section 3). In particular, the crucial ingredient that enables us to handle general types of prediction functions in our main result (Theorem 3.1) is a version of Rosenthal's inequality due to Talagrand [61], which is applicable to sums of independent Banach-valued random variables. Moreover, this allows our analysis to be fully *nonasymptotic*.

To make a more direct comparison with the main theoretical result in our previous paper in the classification setting [36], there are three points to highlight. First, the previous analysis requires that the classifier functions, say $Q_1, \ldots, Q_t$, have a particular form, which is not generally satisfied by the decision tree classifiers in random forests—whereas our current theory is applicable to *actual* random forests. Second, if we let $\omega(x) = \mathbb{E}[Q_1(x)|\mathcal{D}]$ for any fixed $x \in \mathcal{X}$, then the previous analysis assumes that the distribution $\mathcal{L}(\omega(X)|\mathcal{D})$ has a continuously differentiable density function, while the current analysis involves no analogue of this condition. Third, the previous result on bootstrap consistency is stated in terms of a distributional limit and does not provide a rate of convergence. Instead, our current result avoids the reliance on such a limit and gives a more quantitative description of coverage probability.

*Links between inference and computation.* Traditionally, bootstrap methods have been viewed by statisticians as a way to use computation in the service of inference. For this reason, it should be emphasized that our work looks at bootstrap methods from a *reciprocal perspective*, since we aim to use inference in the service of computation (viz. using a quantile estimate to measure algorithmic convergence).

More generally, this way of looking at bootstrap methods has the potential to be applied to the convergence analysis of other randomized algorithms. For instance, in the growing field of randomized numerical linear algebra (or "matrix sketching"), it turns out that convergence can often be framed in terms of the quantiles of certain error variables. Some specific examples include randomized algorithms for matrix multiplication, least-squares, and singular value decomposition [39, 40, 37], and we refer the reader to the recent survey [42, pp. 14–18] for a related discussion of the potential of bootstrap methods in this context. In addition, several variants of bootstrap methods have attracted interest as a way to assess the quality of solutions obtained from stochastic gradient descent (SGD) algorithms [21, 32, 60, 20]. Likewise, given the rising use of randomized algorithms in data science, it seems that considerable opportunity remains for developing bootstrap methods along these lines.

*Outline.* The remainder of the paper is organized as follows. The proposed methods are described in section 2, and our theoretical results on bootstrap consistency are presented in section 3. Next, computational cost is assessed in section 4, and numerical experiments are given in section 5. Finally, all proofs are given in the supplementary material (M134330_01.pdf [local/web 586KB]).

**2. Methodology.** Below, we present our core method for measuring algorithmic convergence with respect to $\mathrm{MSE}_t$ in subsection 2.1. Later on, we show how this approach can be extended to measuring convergence with respect to variable importance in subsection 2.3.

**2.1. Measuring convergence with respect to mean-squared error.** The intuition for the proposed method is based on two main considerations. First, the definition of $\mathrm{MSE}_t$ in (1.3)

shows that it can be interpreted as a functional of $\bar{T}_t$. Namely, if we let $f : \mathcal{X} \to \mathbb{R}$ denote a generic function, then we define the functional $\psi$ according to

$$(2.1) \qquad \psi(f) = \int_{\mathcal{X} \times \mathbb{R}} (y - f(x))^2 d\nu(x, y),$$

and it follows that $\text{MSE}_t$ can be written as

$$(2.2) \qquad \text{MSE}_t = \psi(\bar{T}_t).$$

Second, it is a general principle that bootstrap methods are well suited to approximating distributions derived from smooth functionals of sample averages—which is precisely what the representation (2.2) entails.

To make a more detailed connection between these general ideas and the problem of estimating $q_{1-\alpha}(t)$, recall that we aim to approximate the distribution of the gap $\text{MSE}_t - \text{mse}_\infty$ rather than just $\text{MSE}_t$ itself. Fortunately, the limiting value $\text{mse}_\infty$ can be linked with $\psi$ through the function $\vartheta$ defined by

$$(2.3) \qquad \vartheta(x) = \mathbb{E}[\bar{T}_t(x)|\mathcal{D}],$$

where the expectation is only over the algorithmic randomness in $\bar{T}_t$ (i.e., over the random vector $\boldsymbol{\xi}_t$). More specifically, when the functions $T_1, \ldots, T_t$ satisfy the representation (1.2), the law of large numbers implies $\text{mse}_\infty = \psi(\vartheta)$ under basic integrability assumptions, which leads to the relation

$$(2.4) \qquad \text{MSE}_t - \text{mse}_\infty = \psi(\bar{T}_t) - \psi(\vartheta).$$

This relation is the technical foundation for the proposed method, since it suggests that in order to mimic the fluctuations of $\text{MSE}_t - \text{mse}_\infty$, we can develop a bootstrap method by viewing the functions $T_1, \ldots, T_t$ as "observations" and viewing $\bar{T}_t$ as an estimator of $\vartheta$. In other words, if we sample $t$ functions $T_1^*, \ldots, T_t^*$ with replacement from $T_1, \ldots, T_t$, then we can formally define a bootstrap sample of $\text{MSE}_t - \text{mse}_\infty$ according to

$$(2.5) \qquad \text{MSE}_t^* - \text{MSE}_t = \psi(\bar{T}_t^*) - \psi(\bar{T}_t),$$

where $\bar{T}_t^* := \frac{1}{t} \sum_{i=1}^{t} T_i^*$. In turn, after generating a collection of such bootstrap samples, we can use their empirical $(1 - \alpha)$-quantile as an estimate of $q_{1-\alpha}(t)$. However, as a technical point, it should be noted that (2.5) is a "theoretical" bootstrap sample of $\text{MSE}_t - \text{mse}_\infty$, because the functional $\psi$ depends on the unknown distribution of the test point $\mathcal{L}(X, Y)$. Nevertheless, the same reasoning can still be applied by replacing $\psi$ with an estimate $\hat{\psi}$, which will be explained in detail later in this subsection. Altogether, the method is summarized by Algorithm 2.1 below.

**2.2. Using hold-out or out-of-bag samples.** To complete our discussion of Algorithm 2.1, it remains to clarify how the functional $\psi$ can be estimated from either hold-out samples, or so-called "out-of-bag" (OOB) samples. With regard to the first case, suppose a set of $m$ labeled

---

**Algorithm 2.1.** Bootstrap method for estimating $q_{1-\alpha}(t)$.

**For** $b = 1, \ldots, B$**:**

- Sample $t$ functions $T_1^*, \ldots, T_t^*$ with replacement from $T_1, \ldots, T_t$.
- Compute the bootstrap sample $z_{t,b} := \hat{\psi}(\bar{T}_t^*) - \hat{\psi}(\bar{T}_t)$.

**Return:** the empirical $(1 - \alpha)$-quantile of $z_{t,1}, \ldots, z_{t,B}$ to estimate $q_{1-\alpha}(t)$.

---

samples $\tilde{\mathcal{D}} = \{(\tilde{X}_1, \tilde{Y}_1), \ldots, (\tilde{X}_m, \tilde{Y}_m)\}$ has been held out from the training set $\mathcal{D}$. Using this set, the estimate $\hat{\psi}(\bar{T}_t)$ in Algorithm 2.1 can be easily obtained as

$$(2.6) \qquad \hat{\psi}(\bar{T}_t) \;=\; \frac{1}{m} \sum_{j=1}^{m} (\tilde{Y}_j - \bar{T}_t(\tilde{X}_j))^2.$$

Analogously, we may also obtain $\hat{\psi}(\bar{T}_t^*)$ by using $\bar{T}_t^*$ instead of $\bar{T}_t$ in the formula above.

If the regression functions $T_1, \ldots, T_t$ are trained via bagging or random forests, it is possible to avoid the use of a hold-out set by taking advantage of OOB samples, which are a unique attribute of these methods. To define the notion of an OOB sample, recall that these methods train each function $T_i$ using a random set $\mathcal{D}_i^*$ obtained from $\mathcal{D}$ by sampling with replacement. Due to this sampling mechanism, it follows that each set $\mathcal{D}_i^*$ is likely to exclude approximately $(1 - \frac{1}{n})^n \approx 37\%$ of the training points in $\mathcal{D}$. So, as a matter of terminology, if a particular training point $X_j$ does not appear in $\mathcal{D}_i^*$, we say that $X_j$ is "out-of-bag" for the function $T_i$. Also, we write $\mathrm{OOB}(X_j) \subset \{1, \ldots, t\}$ to denote the index set corresponding to the functions for which $X_j$ is OOB.

From a statistical point of view, OOB samples are important because they serve as "effective" hold-out points. (That is, if $X_j$ is OOB for $T_i$, then the function $T_i$ "never touched" the point $X_j$ during the training process.) Hence, it is natural to consider the following alternative estimate of $\psi$ based on OOB samples:

$$(2.7) \qquad \hat{\psi}_{\mathrm{O}}(\bar{T}_t) \;=\; \frac{1}{n} \sum_{j=1}^{n} (Y_j - \bar{T}_{t,\mathrm{o}}(X_j))^2,$$

where we define $\bar{T}_{t,\mathrm{o}}(X_j)$ to be the average over the functions for which $X_j$ is OOB,

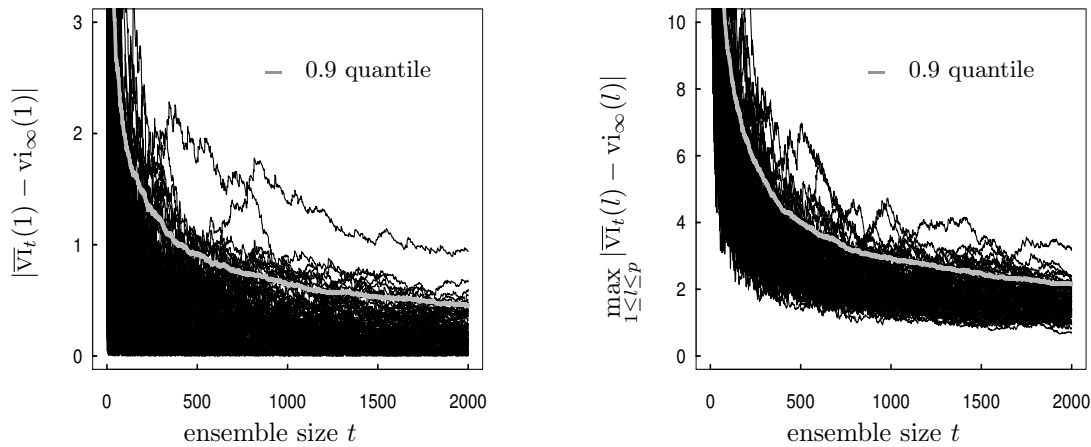$$\bar{T}_{t,\mathrm{o}}(X_j) = \tfrac{1}{|\mathrm{OOB}(X_j)|} \sum_{i \in \mathrm{OOB}(X_j)} T_i(X_j),$$

and $|\cdot|$ refers to the cardinality of a set. Similarly, the quantity $\hat{\psi}_{\mathrm{O}}(\bar{T}_t^*)$ may be defined in terms of a corresponding average with the functions $T_1^*, \ldots, T_t^*$. Lastly, in the case when $\mathrm{OOB}(X_j)$ is empty, we arbitrarily define $\bar{T}_{t,\mathrm{o}}(X_j) = Y_j$, but this occurs very rarely. In fact, it can be checked that for a given point $X_j$, the set $\mathrm{OOB}(X_j)$ is empty with probability approximately equal to $(0.63)^t$.

**2.3. Measuring convergence with respect to variable importance.** In addition to their broad application in prediction problems, randomized ensembles have been very popular for the task of variable selection, e.g., [18, 58, 30, 23, 41, 24, 27]. Although a variety of procedures have been proposed for variable selection in this context, they are generally based on a common approach of ranking the variables according to a measure of averaged variable importance (VI). Under this approach, the averaged VI assigned to each variable typically converges to a limiting value as the ensemble becomes large. However, in practice, the user does not know how this convergence depends on the ensemble size—much like we have seen already for $\text{MSE}_t$.

*Uniform convergence across variables.* Before moving on to the details of our extended method, it is worth mentioning an extra subtlety of measuring algorithmic convergence for VI. Specifically, we must keep in mind that because variable selection is based on ranking, it is important that algorithmic convergence is reached uniformly across variables. In other words, if the VI for some variables converges more slowly than for others, then the ranking of variables will be distorted by purely algorithmic effects. Motivated by this issue, our extended method will provide a way to ensure that algorithmic convergence is achieved in a uniform sense. As an illustration of this point, Figure 2 shows how uniform convergence of VI across several variables can differ considerably from the convergence of VI for a single variable.



**Figure 2.** *Left panel: 1,000 sample paths of $|\overline{\text{VI}}_t(1) - vi_\infty(1)|$, with the true 0.9 quantile curve in gray. Right panel: 1,000 sample paths of the variable $\max_{1 \le l \le p} |\text{VI}_t(l) - vi_\infty(l)|$, with the true 0.9 quantile curve in gray. Both panels were obtained from the Music dataset described in section 5.*

*Setup for variable importance.* To describe algorithmic convergence for VI in detail, let $T_1, \dots, T_t$ be a randomized ensemble that satisfies the representation (1.2), and consider a situation where the space $\mathcal{X}$ is $p$-dimensional. Also, suppose that for each function $T_i$, we have a rule for computing an associated value, say $\text{VI}_i(l)$, to each variable $l \in \{1, \dots, p\}$. (Note that since $T_i$ is a random function, it follows that $\text{VI}_i(l)$ is random as well.) Likewise, the vector of such values associated with $T_i$ is denoted $\text{VI}_i = (\text{VI}_i(1), \dots, \text{VI}_i(p))$, and the

average over $i = 1, \ldots, t$ is denoted as

$$(2.8) \qquad \overline{\mathrm{VI}}_t = \frac{1}{t} \sum_{i=1}^{t} \mathrm{VI}_i.$$

Hence, by comparing the entries of this vector, i.e., $(\overline{\mathrm{VI}}_t(1), \ldots, \overline{\mathrm{VI}}_t(p))$, the user is able to rank the variables, and this is commonly done using a built-in option from the standard random forests software package [33].

Up to this point, we have not specified a particular rule for computing the values $\mathrm{VI}_i(l)$, but several choices are available. For instance, two of the standard choices are based on the notions of "node impurity" (for regression trees) or "random permutations" (for general regression functions). However, from an abstract point of view, our proposed method does not depend on the underlying details of these rules, and so we refer the reader to the book [22, sect. 15.3.2] for additional background. Indeed, our proposed method is applicable to any VI rule, provided that the random vectors $\mathrm{VI}_1, \ldots, \mathrm{VI}_t$ are conditionally i.i.d. given $\mathcal{D}$—and this is satisfied by both of the standard rules when $T_1, \ldots, T_t$ follow the representation (1.2). Also, it should be mentioned that a considerable literature has investigated limitations and improvements of the standard VI rules, e.g., [59, 51, 58, 47, 45], and the study of variable importance in this context continues to be an open direction of research.

When the vectors $\mathrm{VI}_1, \ldots, \mathrm{VI}_t$ are conditionally i.i.d. given $\mathcal{D}$, the vector $\overline{\mathrm{VI}}_t$ will typically converge to a limit, say $\mathrm{vi}_\infty \in \mathbb{R}^p$, as $t \to \infty$ and $\mathcal{D}$ is held fixed. In order to measure this convergence uniformly across $l \in \{1, \ldots, p\}$, we will focus on the (unobserved) random variable defined by

$$(2.9) \qquad \varepsilon_t := \max_{1 \le l \le p} |\overline{\mathrm{VI}}_t(l) - \mathrm{vi}_\infty(l)|,$$

and our goal will be to estimate its $(1 - \alpha)$-quantile, denoted as

$$(2.10) \qquad \mathsf{q}_{1-\alpha}(t) := \inf \left\{ q \in [0, \infty) \,\middle|\, \mathbb{P}\big(\varepsilon_t \le q \,|\, \mathcal{D}\big) \ge 1 - \alpha \right\}.$$

*The bootstrap method for variable importance.* By analogy with our method for estimating the quantiles of $\mathrm{MSE}_t - \mathrm{mse}_\infty$, we propose to construct bootstrap samples of $\varepsilon_t$ by resampling the vectors $\mathrm{VI}_1, \ldots, \mathrm{VI}_t$, and then estimating $\mathsf{q}_{1-\alpha}(t)$ with the empirical $(1 - \alpha)$-quantile. In algorithmic form, the procedure is described in Algorithm 2.2 below.

Numerical results illustrating the performance of this algorithm, as well as Algorithm 2.1, are given section 5. Also, in Appendix F of the supplementary material, we show how Algorithm 2.2 can be adapted to the situation where convergence is measured in terms of the relative error variable $\max_{1 \le l \le p} |\overline{\mathrm{VI}}_t(l) - \mathrm{vi}_\infty(l)|/|\mathrm{vi}_\infty|$. Numerical results for the case of relative error are provided there as well.

---

**Algorithm 2.2.** Bootstrap method for estimating $\mathsf{q}_{1-\alpha}(t)$.

---

**For** $b = 1, \ldots, B$ **:**

- Sample $t$ vectors $(\mathrm{VI}_1^*, \ldots, \mathrm{VI}_t^*)$ with replacement from $(\mathrm{VI}_1, \ldots, \mathrm{VI}_t)$, and let $\overline{\mathrm{VI}}_t^* = \frac{1}{t} \sum_{i=1}^{t} \mathrm{VI}_i^*$.

- Compute the bootstrap sample

$$(2.11) \qquad \varepsilon_{t,b}^* := \max_{1 \le l \le p} |\overline{\mathrm{VI}}_t^*(l) - \overline{\mathrm{VI}}_t(l)|.$$

**Return:** the empirical $(1 - \alpha)$-quantile of $\varepsilon_{t,1}^*, \ldots, \varepsilon_{t,B}^*$ to estimate $\mathsf{q}_{1-\alpha}(t)$.

---

**3. Main result.** In this section, we develop the main theoretical result of the paper (Theorem 3.1), which quantifies the coverage probability of the bootstrap estimate $\hat{q}_{1-\alpha}(t)$ for $q_{1-\alpha}(t)$. Namely, we will show that for a fixed set $\mathcal{D}$, the inequality

$$(3.1) \qquad \mathrm{MSE}_t - \mathrm{mse}_\infty \ \le \ \hat{q}_{1-\alpha}(t)$$

holds with a probability that is not much less than $1 - \alpha$. Later on, we will also show that a corresponding result holds for estimating the quantile $\mathsf{q}_{1-\alpha}(t)$ in the context of variable importance (cf. subsection 3.1).

To establish the main result, we will rely on a common type of simplification, which is to exclude sources of error beyond the resampling process itself. More specifically, we will focus on bootstrap samples of the form $\mathrm{MSE}_t^* - \mathrm{MSE}_t$, defined in (2.5), since these are not affected by the extraneous error from estimating the functional $\psi$. (In other words, these samples are different from those of the form $\hat{\psi}(\bar{T}_t^*) - \hat{\psi}(\bar{T}_t)$ and $\hat{\psi}_\mathrm{o}(\bar{T}_t^*) - \hat{\psi}_\mathrm{o}(\bar{T}_t)$ described in subsection 2.2.) Meanwhile, even with such a simplification, the proof of the result is still quite involved. Also, this same choice was used in our previous analysis of the classification setting for the same reasons [36], but apart from this detail, the analysis in the current paper is entirely different.

With regard to the ensemble, it will only be assumed to satisfy the representation (1.2) and a basic moment condition in Theorem 3.1. From the standpoint of existing theory for randomized ensembles, these assumptions are very mild, because the representation (1.2) is always satisfied by bagging and random forests. By contrast, it is much more common in the theoretical literature to work with ensembles that are simpler than the ones used in practice, and similarly, our previous work in the classification setting relied on a specialized type of ensemble. Finally, it is notable that our result is fully *nonasymptotic*, whereas much existing work on the convergence of randomized ensembles has taken an asymptotic approach that does not always provide explicit rates of convergence.

*Notation.* If $g$ and $h$ are real-valued functions on $\mathcal{X} \times \mathbb{R}$, we denote their inner product with respect to the test point distribution $\nu = \mathcal{L}(X, Y)$ as

$$\langle g, h \rangle = \int_{\mathcal{X} \times \mathbb{R}} g(x, y) \, h(x, y) \, d\nu(x, y),$$

and accordingly, we write $\|g\|_{L_2} = \sqrt{\langle g, g \rangle}$. In addition, recall the function $\vartheta(x) = \mathbb{E}[T_1(x)|\mathcal{D}]$ from (2.3), and define the random variable

$$\zeta = 2 \langle \vartheta - y, T_1 - \vartheta \rangle,$$

where $\vartheta - y$ is understood as the function that sends $(x, y)$ to $\vartheta(x) - y$. When the random variable $\zeta$ is conditioned on $\mathcal{D}$, we denote its standard deviation by

$$\sigma(\mathcal{D}) = \sqrt{\mathrm{var}(\zeta|\mathcal{D})},$$

and the finiteness of this quantity will follow from assumption **A2** below. Also, all expressions involving $1/\sigma(\mathcal{D})$ will be understood as $\infty$ in the exceptional case when $\sigma(\mathcal{D}) = 0$. Lastly, for each positive integer $k$, we define the moment parameter

$$(3.2) \qquad\qquad \gamma_k(\mathcal{D}) = \left( \mathbb{E} \left[ \|T_1 - y\|_{L_2}^{2k} \big| \mathcal{D} \right] \right)^{1/k}.$$

To interpret the role of this parameter, note that the random variable $\mathrm{MSE}_t$ can be written as $\| \frac{1}{t} \sum_{i=1}^{t} (T_i - y) \|_{L_2}^2$. Hence, the fluctuations of $\mathrm{MSE}_t$ are determined by the tail behavior of the summands $T_i - y$, and the parameter $\gamma_k(\mathcal{D})$ describes the tails of the summands through their moments.

   *Assumptions.* With the above notation in place, the two assumptions for our main result may be stated as follows.

**A1.** The ensemble $T_1, \ldots, T_t$ can be represented in the form (1.2).

**A2.** There is at least one integer $k \geq 2$ such that $\gamma_{3k}(\mathcal{D}) < \infty$.

Regarding the finiteness of $\gamma_{3k}(\mathcal{D})$ in **A2**, it is noteworthy that this condition is satisfied for any $k$ whenever the regression functions $T_1, \ldots, T_t$ are trained by the standard method of CART and the test label distribution has moments of all orders. This is because the regression trees trained by CART have a range that is determined by the training labels $Y_1, \ldots, Y_n$. In particular, if we define $M(\mathcal{D}) = \max_{1 \leq i \leq n} |Y_i|$, then every tree $T_i$ satisfies $\sup_{x \in \mathcal{X}} |T_i(x)| \leq M(\mathcal{D})$. The same reasoning also applies beyond CART to any other method whose predictions are obtained as local averages of training labels.

   We now state the main result of the paper.

   **Theorem 3.1.** *Suppose that* **A1** *and* **A2** *hold. In addition, fix any small constant* $\alpha \in (0, 1)$, *and let* $k \geq 2$ *be as in* **A2**. *Lastly, let* $\hat{q}_{1-\alpha}(t)$ *denote the empirical* $(1 - \alpha)$-*quantile of* $B$ *bootstrap samples of the form* (2.5), *and define the quantity*

$$(3.3) \qquad\qquad \delta(\mathcal{D}) \;=\; \frac{k^2}{\sqrt{t}} \left( \frac{\gamma_{3k}(\mathcal{D})}{\sigma(\mathcal{D})} \right)^3 \;+\; e^{-k/2} \;+\; \sqrt{\frac{\log(B)}{B}}.$$

*Then, there is an absolute constant* $c_0 > 0$ *such that* $\hat{q}_{1-\alpha}(t)$ *satisfies*

$$(3.4) \qquad\qquad \mathbb{P}\Big( \mathrm{MSE}_t - \mathrm{mse}_\infty \leq \hat{q}_{1-\alpha}(t) \,\Big|\, \mathcal{D} \Big) \;\geq\; 1 - \alpha - c_0\, \delta(\mathcal{D}).$$

*Remarks.* In essence, the result shows that $\hat{q}_{1-\alpha}(t)$ bounds the unknown convergence gap $\text{MSE}_t - \text{mse}_\infty$ with a probability that is not much less than the ideal value of $1-\alpha$. To comment on some further aspects of the result, note that the inequality (3.4) has the desirable property of being *scale-invariant* with respect to the labels $Y_1, \dots, Y_n$ and the functions $T_1, \dots, T_t$. More precisely, if we were to change the units of the labels and functions by a common scale factor, it can be checked that both sides of (3.4) would remain unchanged.

Another important aspect of Theorem 3.1 deals with the dependence of $\delta(\mathcal{D})$ on the value of $k$, and it is of interest to develop a bound on $\delta(\mathcal{D})$ that simplifies this dependence. To do this, we can look at a basic situation where the regression functions are trained by CART and the test label variable is bounded. In addition, we may consider the particular choice

$$(3.5) \qquad k = \lceil \log(t) - 4\log\log(t) \rceil,$$

which leads to the following bounds:

$$e^{-k/2} \ \leq \ \frac{\log(t)^2}{\sqrt{t}} \qquad \text{and} \qquad \frac{k^2}{\sqrt{t}} \ \leq \ \frac{c_1 \log(t)^2}{\sqrt{t}}$$

for some absolute constant $c_1 > 0$ and all $t \geq 2$. In turn, it follows that there is a number $c(\mathcal{D}) > 0$ not depending on $t$, $k$, or $B$, such that

$$(3.6) \qquad \delta(\mathcal{D}) \ \leq \ \frac{c(\mathcal{D}) \log(t)^2}{\sqrt{t}} \ + \ \sqrt{\frac{\log(B)}{B}},$$

which provides a considerable simplification. Hence, under the conditions just mentioned, and with $\mathcal{D}$ held fixed, the quantity $\delta(\mathcal{D})$ converges to 0 at *nearly parametric rates* with respect to $t$ and $B$.

**3.1. Bootstrap consistency in the context of variable importance.** Having developed our main result as a consistency guarantee for Algorithm 2.1 in the context of mean-squared error, we now aim to establish a corresponding result for Algorithm 2.2 in the context of variable importance, which is given as Theorem 3.2 below.

*Setting and assumptions.* In order to formulate this result, we will proceed along the lines of the setup described in subsection 2.3. Recall that for each random function $T_i$ with $i \in \{1, \dots, t\}$, there is an associated random vector $\text{VI}_i = (\text{VI}_i(1), \dots, \text{VI}_i(p)) \in \mathbb{R}^p$, where $\text{VI}_i(l)$ refers to the importance assigned to the variable $l$ by the function $T_i$, and the sample average is denoted $\overline{\text{VI}}_t = \frac{1}{t}\sum_{i=1}^{t} \text{VI}_i$. The only two conditions required of $\text{VI}_1, \dots, \text{VI}_t$ are as follows.

**A3.** The random vectors $\text{VI}_1, \dots, \text{VI}_t \in \mathbb{R}^p$ are conditionally i.i.d. given $\mathcal{D}$.

**A4.** There are positive numbers $b(\mathcal{D})$ and $b'(\mathcal{D})$ such that the following inequalities hold almost surely for all $l \in \{1, \dots, p\}$:

$$(3.7) \qquad b(\mathcal{D}) \leq \sqrt{\text{var}(\text{VI}_1(l)|\mathcal{D})} \qquad \text{and} \qquad \text{VI}_1(l) \leq b'(\mathcal{D}).$$

Perhaps the most important point to emphasize about **A3** is that it is automatically satisfied by two of the standard variable importance measures used within random forests,

namely the "node impurity" measure and the "random permutations" measure [33]. More generally, as long as each vector $\mathrm{VI}_i$ can be computed as a function of $T_i$, and as long as $T_i$ can be represented in the abstract form (1.2), then **A3** will hold. With regard to the first inequality in **A4**, this is simply a nondegeneracy condition, which rules out situations where $\mathrm{VI}_i(l)$ has no algorithmic fluctuations. Meanwhile, the second inequality in **A4** is always satisfied by the two standard variable importance measures in random forests when each $T_i$ is trained via CART. Lastly, the condition **A4** ensures that $\overline{\mathrm{VI}}_t$ has a limit as $t \to \infty$ with $\mathcal{D}$ held fixed, which is given by $\mathrm{vi}_\infty = \mathbb{E}[\mathrm{VI}_1|\mathcal{D}]$.

The gist of Theorem 3.2 below is that the output $\hat{\mathsf{q}}_{1-\alpha}(t)$ of Algorithm 2.2 has reliable coverage probability when it is used as an upper bound on $\max_{1 \le l \le p} |\overline{\mathrm{VI}}_t(l) - \mathrm{vi}_\infty(l)|$.

**Theorem 3.2.** *Suppose that **A3** and **A4** hold, and fix any small constants $\alpha, \eta \in (0, 1)$. In addition, let $\hat{\mathsf{q}}_{1-\alpha}(t)$ denote the empirical $(1-\alpha)$-quantile of $B$ bootstrap samples of the form (2.11), and define the quantity*

$$(3.8) \qquad \tilde{\delta} = \sqrt{\frac{\log(2pt)^3}{t}} \; + \; \sqrt{\frac{\log(B)}{B}}.$$

*Then, there is a number $\tilde{c}(\mathcal{D}) > 0$ depending only on the triple $(\eta, b(\mathcal{D}), b'(\mathcal{D}))$ such that $\hat{\mathsf{q}}_{1-\alpha}(t)$ satisfies*

$$(3.9) \qquad \mathbb{P}\left( \max_{1 \le l \le p} |\overline{\mathrm{VI}}_t(l) - \mathrm{vi}_\infty(l)| \le \hat{\mathsf{q}}_{1-\alpha}(t) + \eta \;\middle|\; \mathcal{D} \right) \; \ge \; 1 - \alpha - \tilde{c}(\mathcal{D})\,\tilde{\delta}.$$

*Remarks.* Just like Theorem 3.1 given earlier, this result quantifies coverage probability in a nonasymptotic manner. On the other hand, one small point of contrast with Theorem 3.1 is the constant $\eta \in (0, 1)$ in the present result, which serves only as a theoretical expedient and can be fixed at an *arbitrarily small* value. Concerning the proof, it leverages recent advances on bootstrap methods for "max statistics" [16]. Furthermore, under some extra structural assumptions on the covariance matrix of $\mathrm{VI}_1$, it is possible to replace the error term $\log(2pt)^{3/2}t^{-1/2}$ in (3.8) with a dimension-free term of the form $t^{-1/2+\epsilon_0}$ for an arbitrarily small constant $\epsilon_0 > 0$ [38].

**4. Computation and speedups.** In order for the proposed method to be a practical tool for checking algorithmic convergence, its computational cost should be manageable in comparison to training the ensemble itself. Below, in subsection 4.1, we offer a quantitative comparison, showing that under simple conditions, Algorithm 2.1 and Algorithm 2.2 are not a bottleneck in relation to training $t$ regression functions with CART. Additionally, we show in subsection 4.2 how an extrapolation technique from our previous work on classification can be improved in our current setting with a *bias correction rule*.

**4.1. Cost comparison.** Since the CART method is based on a greedy iterative algorithm, the exact computational cost of training a regression tree is difficult to describe analytically. Due to this difficulty, the authors of CART studied its cost in the simplified situation where each node of a regression tree is split into exactly 2 child nodes (except for the leaves). To be more precise, suppose $\mathcal{X} \subset \mathbb{R}^p$, and let $d \ge 2$ denote the "depth" of the tree, so that there are $2^d$ leaves. In addition, suppose that when the algorithm splits a given node, it searches over

$\lceil p/3 \rceil$ candidate variables that are randomly chosen from $\{1, \ldots, p\}$, which is the default rule when CART is used by random forests [33]. Based on these assumptions, the analysis in the book [12, p. 166] shows that the number of operations involved in training $t$ such trees is at least of order $\Omega(t \cdot p \cdot d \cdot n)$.[1]

*The cost of Algorithm* 2.1. To determine the cost of Algorithm 2.1, it is important to clarify that when bagging and random forests are used in practice, the prediction error of the ensemble is typically estimated automatically using either hold-out or OOB samples. As a result, the predicted values of each tree on these samples can be regarded as being precomputed by the ensemble method. Once these values are available, the subsequent cost of Algorithm 2.1 is simple to measure. Specifically, in the case of hold-out samples, (2.6) shows that the cost to obtain $\hat{\psi}(\bar{T}_t) - \hat{\psi}(\bar{T}_t^*)$ for each bootstrap sample is $\mathcal{O}(t \cdot m)$, which leads to an overall cost that is $\mathcal{O}(B \cdot t \cdot m)$. Similarly, for the case of OOB samples, the overall cost is $\mathcal{O}(B \cdot t \cdot n)$. Altogether, this leads to the conclusion that the cost of Algorithm 2.1 does not exceed that of training the ensemble if the number of bootstrap samples satisfies the very mild condition

$$(4.1) \qquad\qquad B = \mathcal{O}(p \cdot d),$$

and this applies to either the hold-out or OOB cases, provided $m = \mathcal{O}(n)$. Moreover, our discussion in subsection 4.2 will show that the condition (4.1) can be relaxed even further via extrapolation.

Beyond the fact that Algorithm 2.1 compares well with the cost of training an ensemble, there are several other favorable aspects to mention. First, the algorithm only relies on predicted labels for its input, and it never needs to access any points in the space $\mathcal{X}$. In particular, this means that the cost of the algorithm is independent of the dimension of $\mathcal{X}$. Second, the bootstrap samples in Algorithm 2.1 are simple to compute in parallel, which means that the runtime of the algorithm can essentially be reduced by a factor of $B$.

*The cost of Algorithm* 2.2. Many of the previous considerations for Algorithm 2.1 also apply to Algorithm 2.2, but it turns out that the cost of Algorithm 2.2 can be much less when $n$ is large. Because each bootstrap sample in Algorithm 2.2 requires forming an average of $t$ vectors in $\mathbb{R}^p$, it is straightforward to check that the overall cost is $\mathcal{O}(B \cdot t \cdot p)$, where we view the vectors $\text{VI}_1, \ldots, \text{VI}_t$ as being precomputed by the ensemble method. Consequently, the cost is independent of $n$, and the algorithm is thus highly scalable. Furthermore, under the setup of our earlier cost comparison with CART, the cost of Algorithm 2.2 does not exceed the cost of training the ensemble if

$$B = \mathcal{O}(n \cdot d),$$

which allows for plenty of bootstrap samples in practice. Better still, our numerical experiments show that just a few dozen bootstrap samples can be sufficient when $n$ is on the order of $10^4$, indicating that Algorithm 2.2 is quite inexpensive in comparison to training.

**4.2. Further reduction of cost by extrapolation.** The basic idea of extrapolation is to check algorithmic convergence for a small "initial" ensemble, say of size $t_0$, and then use this information to "look ahead" and predict convergence for a larger ensemble of size $t > t_0$. This

---

[1]We use $\Omega(\cdot)$ and $\mathcal{O}(\cdot)$ in the conventional way, so that they respectively refer to lower and upper bounds that hold up to constants [26, sect. 9.2].

general technique has a long history in the development of resampling methods and numerical algorithms, and further background can be found in [8, 3, 4, 7, 13, 57] among others. In the remainder of this section, we first summarize how extrapolation was previously developed in our companion paper [36] and then explain how that approach can be improved in the present context with a bias correction rule for OOB samples.

*A basic version of extrapolation.* At a technical level, our use of extrapolation is based on the central limit theorem, which suggests that the fluctuations of $\text{MSE}_t - \text{mse}_\infty$ should scale like $1/\sqrt{t}$ as a function of $t$. As a result, we expect that the quantile $q_{1-\alpha}(t)$ should behave like

$$q_{1-\alpha}(t) \approx \frac{\kappa}{\sqrt{t}}$$

for some quantity $\kappa$ that may depend on all problem parameters except $t$.

To take advantage of this heuristic scaling property, suppose that we train an initial ensemble of size $t_0$ and run Algorithm 2.1 to obtain an estimate $\hat{q}_{1-\alpha}(t_0)$. We can then extract an estimate of $\kappa$ by defining

$$\hat{\kappa} = \sqrt{t_0}\, \hat{q}_{1-\alpha}(t_0).$$

Next, we can rapidly estimate $q_{1-\alpha}(t)$ for all subsequent $t \geq t_0$ by defining the extrapolated estimate

$$(4.2) \qquad \hat{q}_{1-\alpha}^{\text{ext}}(t) = \frac{\hat{\kappa}}{\sqrt{t}} = \frac{\sqrt{t_0}\hat{q}_{1-\alpha}(t_0)}{\sqrt{t}}.$$

In particular, there are two crucial benefits of this estimate: (1) It is much faster to apply Algorithm 2.1 to a small initial ensemble of size $t_0$ than to a large one of size $t$. (2) If we would like $\text{MSE}_t$ to be within some tolerance $\epsilon > 0$ of the limit $\text{mse}_\infty$, then we can use the condition

$$\hat{q}_{1-\alpha}^{\text{ext}}(t) \leq \epsilon$$

to *dynamically predict* how large $t$ must be chosen to reach that tolerance, namely $t \geq (\sqrt{t_0}\hat{q}_{1-\alpha}(t_0)/\epsilon)^2$.

*Bias-corrected extrapolation.* If the initial estimate $\hat{q}_{1-\alpha}(t_0)$ is obtained by implementing Algorithm 2.1 with OOB samples, it turns out to be a biased estimate of $q_{1-\alpha}(t_0)$. Fortunately, however, it is possible to correct for this bias in a simple way, as we now explain.

To understand the source of the bias, consider a particular training point $X_j$ and note that for an initial ensemble of size $t_0$, the expected number of functions for which $X_j$ is OOB is given by

$$(4.3) \qquad \tau_n(t_0) = (1 - 1/n)^n \cdot t_0.$$

In other words, this means that when an ensemble of size $t_0$ makes a prediction on an OOB point, the "effective" size of the ensemble is $\tau_n(t_0)$ rather than $t_0$. As a result, if we implement Algorithm 2.1 using OOB samples with an initial ensemble of size $t_0$, then the output $\hat{q}_{1-\alpha}(t_0)$ should really be viewed as an estimate of $q_{1-\alpha}(\tau_n(t_0))$ rather than $q_{1-\alpha}(t_0)$.

Based on this reasoning, we can adjust our previous definition of the estimate $\hat{q}_{1-\alpha}^{\text{ext}}(t)$ in (4.2) by using

$$(4.4) \qquad \hat{q}_{1-\alpha}^{\text{ext,o}}(t) = \frac{\sqrt{\tau_n(t_0)}\hat{q}_{1-\alpha}(t_0)}{\sqrt{t}} \qquad \text{for} \qquad t \geq \tau_n(t_0).$$

Later on, in section 5 we will demonstrate that this simple adjustment works well in practice.

*Remark.* As a clarification, it should be noted that the definition (4.4) is only to be used when Algorithm 2.1 is implemented with OOB samples, and the basic rule (4.2) should be used in the case of hold-out samples. Also, the basic rule (4.2) can be easily adapted to extrapolate the estimate produced by Algorithm 2.2, and so we omit the details in the interest of brevity.

**5. Numerical results.** We now demonstrate the bootstrap's numerical accuracy in the tasks of measuring algorithmic convergence with respect to both mean-squared error and variable importance. Overall, our results show that the extrapolated OOB estimate is accurate at predicting the effect of increasing $t$. In fact, the results show that extrapolation succeeds at predicting what will happen when $t$ is increased by a factor of 4 beyond $t_0$, and possibly much farther.

### 5.1. Organization of experiments.
*Data preparation.* Our experiments were based on several natural datasets that were each randomly partitioned in the following way. Letting $\mathcal{F}$ denote the full set of observation pairs $(X_1, Y_1), (X_2, Y_2), \ldots$ for a given dataset, we evenly split $\mathcal{F}$ into a disjoint union $\mathcal{F} = \mathcal{D} \sqcup \mathcal{T}$, where $\mathcal{D}$ was used as a training set, and $\mathcal{T}$ was used as a "ground truth set" to approximate the true quantile curves $q_{1-\alpha}(t)$ and $\mathsf{q}_{1-\alpha}(t)$.

Since Algorithm 2.1 relies on a hold-out set, we also used a relatively small subset $\mathcal{H} \subset \mathcal{T}$ for that purpose. Specifically, the hold-out set $\mathcal{H}$ was chosen so that its cardinality satisfied $|\mathcal{H}|/(|\mathcal{H}| + |\mathcal{D}|) = 1/6$, up to rounding error. This reflects a practical situation where the user can only afford to allocate 1/6 of the available data for the hold-out set. In other words, the idea is to think of the user as only having access to $\mathcal{D} \sqcup \mathcal{H}$, while the set $\mathcal{T}$ is used externally to determine $q_{1-\alpha}(t)$ and $\mathsf{q}_{1-\alpha}(t)$.

Each of the full datasets are briefly summarized below.

- *Housing.* This dataset originates from the 1990 California census and is available as part of the online supplement to the book [25]. The observations correspond to different housing districts, and for each one, there are nine features for predicting the median home price in that district. ($|\mathcal{F}| = 20{,}640$, $|\mathcal{D}| = |\mathcal{T}| = 10{,}320$, $|\mathcal{H}| = 4{,}128$.)

- *Protein.* This dataset was collected from the fifth through the ninth series of CASP experiments [44] and is available at the UCI repository [19] under the title *Physicochemical Properties of Protein Tertiary Structure Data Set*. The observations correspond to artificially generated conformations of proteins (known as decoys) that are described by nine biophysical features. Each decoy can be thought of as a perturbation of an associated "target" protein, and the features are used to predict how far the decoy is from its target. ($|\mathcal{F}| = 45{,}730$, $|\mathcal{D}| = |\mathcal{T}| = 22{,}865$, $|\mathcal{H}| = 4{,}573$.)

- *Music.* This dataset consists of audio recordings (observations) described by 68 features that are used to predict the geographic latitude of the recording, as described in [64]. The dataset is available at the UCI repository [19] under the title *Geographical Origin of Music Data Set*. ($|\mathcal{F}| = 1{,}059$, $|\mathcal{D}| = |\mathcal{T}| = 530$, $|\mathcal{H}| = 106$.)

- *Diamond.* This dataset arises from a collection of diamonds, each described by nine features that are used to predict the diamond's price. The dataset was obtained as

a downsampled version of `diamonds` in the package `ggplot2` [63]. ($|\mathcal{F}| = 10{,}000$, $|\mathcal{D}| = |\mathcal{T}| = 5{,}000$, $|\mathcal{H}| = 1{,}000$.)

*Computing the true quantile curves $q_{1-\alpha}(t)$ and $\mathsf{q}_{1-\alpha}(t)$.* Once a full dataset $\mathcal{F}$ was partitioned as above, we ran the random forests algorithm 1,000 times on the associated set $\mathcal{D}$, using the R package `randomForest` [33]. The overall process was a serious computational undertaking, because 2,000 regression trees were trained during every run, and hence a total of $2 \times 10^6$ trees were trained on each dataset.

During each run, as the ensemble size increased from $t = 1$ up to $t = 2{,}000$, the corresponding true values of $\mathrm{MSE}_t$ were approximated with the ensemble's error rate on $\mathcal{T}$. Also, the true value of $\mathrm{mse}_\infty$ was approximated with the average of the 1,000 approximate values of $\mathrm{MSE}_{2{,}000}$. In this way, the collection of runs produced 1,000 approximate sample paths of $\mathrm{MSE}_t - \mathrm{mse}_\infty$, similar to those illustrated in the right panel of Figure 1. Finally, the quantile curve $q_{.90}(t)$ was approximated by using the empirical 90% quantile of the sample paths at each $t \in \{1, \ldots, 2{,}000\}$.

To handle the setting of variable importance, essentially the same steps were used. Specifically, we measured variable importance in terms of node impurity to compute $\overline{\mathrm{VI}}_t \in \mathbb{R}^p$ at every value $t \in \{1, \ldots, 2{,}000\}$, for each of the 1,000 runs mentioned above. In addition, we approximated $\mathrm{vi}_\infty \in \mathbb{R}^p$ with the average of the 1,000 realizations of $\overline{\mathrm{VI}}_{2{,}000}$. Altogether, these computations provided us with 1,000 approximate sample paths of $\varepsilon_t = \max_{1 \le l \le p} |\overline{\mathrm{VI}}_t(l) - \mathrm{vi}_\infty(l)|$, and then we used the empirical 90% quantile at each $t \in \{1, \ldots, 2{,}000\}$ as an approximation to $\mathsf{q}_{.90}(t)$.

*Applying the bootstrap algorithms with extrapolation.* For each of the described 1,000 runs of random forests, we applied the extrapolated versions of Algorithm 2.1 and Algorithm 2.2 at the initial ensemble size of $t_0 = 500$, using a choice of $B \in \{25, 50, 100\}$ bootstrap samples. (The extrapolation was carried out to a final ensemble size of $t = 2{,}000$.) Also, for Algorithm 2.1, we implemented both the hold-out and OOB versions, including the bias correction for the OOB samples described in (4.4). Hence, this provided us with 1,000 realizations of each type of estimate, allowing for an assessment of their variability.

### 5.2. Numerical results for mean-squared error.

*Organization of the plots.* The hold-out and OOB estimates for $q_{.90}(t)$ are illustrated in Figures 3–6. For each choice of $B \in \{25, 50, 100\}$, the colored curves represent the averages of the estimates over the 1,000 runs described previously, and the error bars display the fluctuations of the estimates over repeated runs—corresponding to the 10th and 90th percentiles of the estimates. (For the values of $t$ between the endpoints, we omit the error bars for clarity. Also, the error bars should *not* be interpreted as confidence intervals for $q_{.90}(t)$; they are only intended to illustrate the variability of the estimates.)

With regard to computation, another point to mention is that the estimates were only computed for the initial ensemble size $t_0 = 500$, and the rest of the estimated curves were obtained essentially *for free* by extrapolation. Lastly, as a clarification, it should be noted that the OOB curves are shifted to the left of the hold-out curves because of the bias correction rule (4.4) for OOB samples.

*Remarks on performance.* The main point to take away from the plots is that the OOB estimate performs well overall and can be noticeably more accurate than the hold-out estimate

(cf. Figures 3, 5, and 6). Furthermore, the OOB estimate has an extra advantage because it does not require the user to hold out any data. For these reasons, we recommend the OOB estimate in practice.

Concerning the number of bootstrap samples, we see the expected pattern that larger values of $B$ reduce the fluctuations of the estimates. Nevertheless, even at $B = 25$, the fluctuations are well behaved. So, for practical purposes, this indicates that the speedup from a small choice of $B$ may outweigh a relatively minor reduction in variance.

Another conclusion to draw from the plots is that the bias correction plays a significant role in the extrapolation of the OOB estimate. To see this, note that if the bias correction were not used, this would be equivalent to shifting the blue curves so that they start at the same point as the green curves, which would clearly lead to a loss in accuracy. Also, it is remarkable that the extrapolated OOB estimate continues to be accurate at a final ensemble size of $t = 2,000$ that is 4 times larger than the initial ensemble size $t_0 = 500$. Hence, this provides the user with a very inexpensive way to predict how quickly the ensemble will converge.

To explain the inferior performance of the hold-out estimate, recall that it uses the small set $\mathcal{H}$ in order to estimate $\mathrm{MSE}_t$. As a result of the small size of $\mathcal{H}$, the estimate of $\mathrm{MSE}_t$ has high variability, which inflates the upper extremes and ultimately leads to a larger estimate of $q_{.90}(t)$. On the other hand, the OOB estimate is able to take advantage of the OOB samples in the much larger set $\mathcal{D}$, which reduces this detrimental effect.
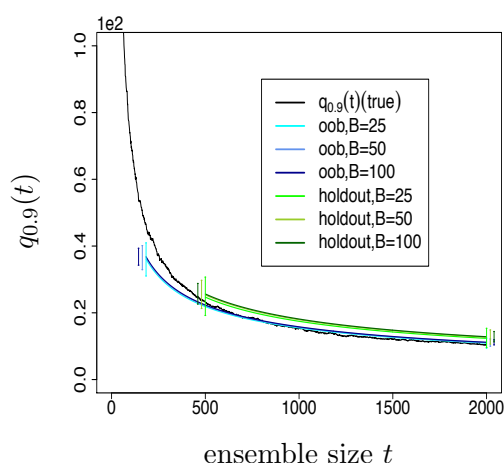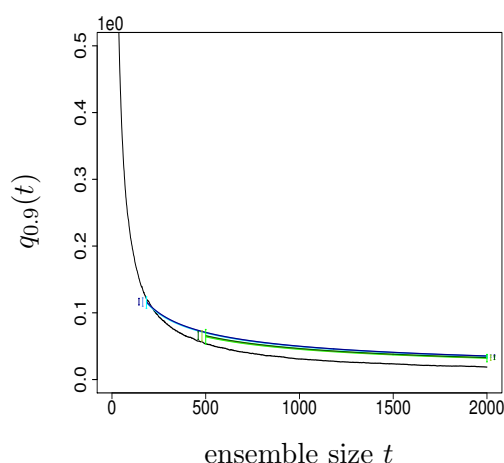


**Figure 3.** *Housing data.*



**Figure 4.** *Protein data.*

**5.3. Numerical results for variable importance.** The results in the setting of variable importance are simpler to describe, since there is only one type of estimate for $q_{.90}(t)$. Figures 7–10 display the average of the 1,000 realizations of the estimate using a blue curve (corresponding to $B = 50$), and as before, the error bars at the endpoints represent the 10th and 90th percentiles. Also, the extrapolation procedure was performed using an initial ensemble size of $t_0 = 500$, as in the previous subsection. From the four plots, it is clear that the extrapolated estimate displays excellent overall performance, with its bias and variance both
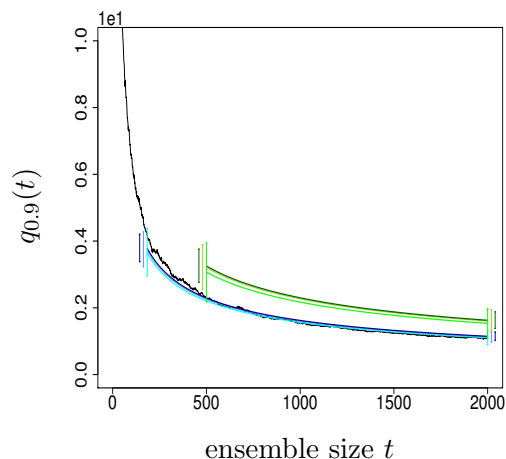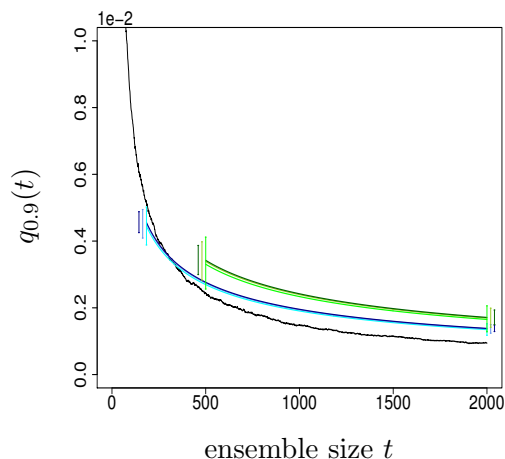
**Figure 5.** *Music data.*
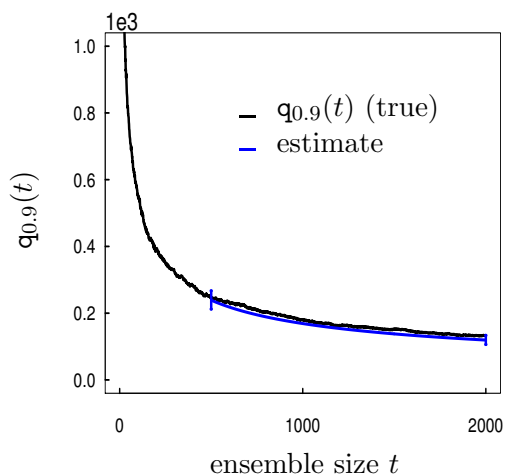


**Figure 6.** *Diamond data.*



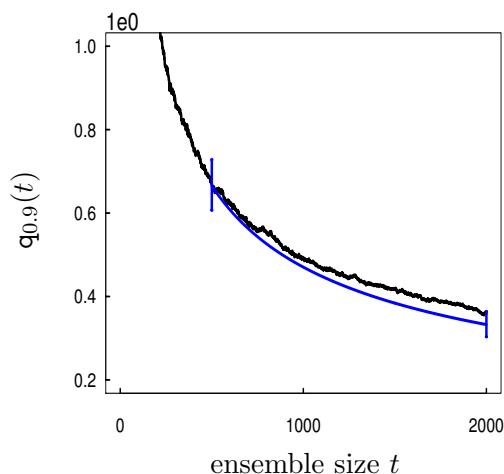**Figure 7.** *Housing data.*



**Figure 8.** *Protein data.*

being very small.

**6. Conclusion.** In this paper, we have developed a bootstrap method that allows users to measure the algorithmic convergence of regression ensembles with a level of precision that has not previously been available. In particular, the method provides users with a systematic way to determine when the ensemble is large enough so that it will perform nearly as well as an ideal infinite ensemble—with respect to either mean-squared error or variable importance. With regard to theory, our approach is supported by guarantees in Theorems 3.1 and 3.2 that quantify the coverage probabilities of the quantile estimates produced by Algorithms 2.1 and 2.2. Computationally, the method incurs only modest cost in comparison to training the ensemble itself, and furthermore, the method naturally lends itself to speedups via parallel computing and extrapolation. Lastly, we have shown empirically that the method has
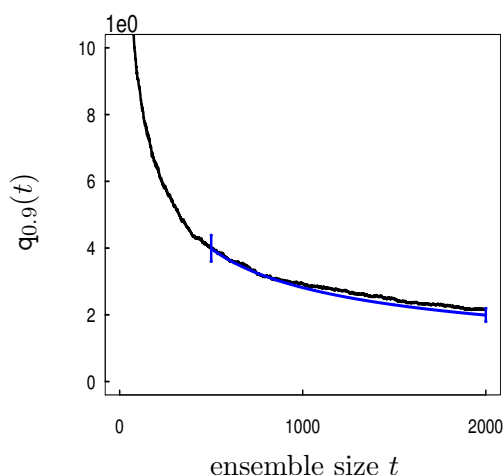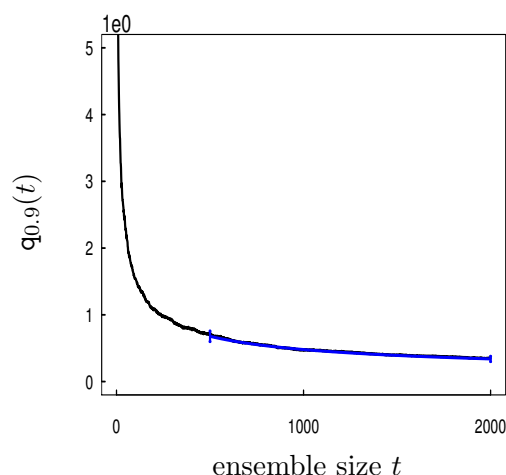
**Figure 9.** *Music data.*



**Figure 10.** *Diamond data.*

encouraging finite-sample performance in a range of situations.

## REFERENCES

[1] S. ARLOT AND R. GENUER, *Analysis of Purely Random Forests Bias*, preprint, https://arxiv.org/abs/1407.3939, 2014.

[2] J. BASILICO, M. MUNSON, T. KOLDA, K. DIXON, AND W. KEGELMEYER, *Comet: A recipe for learning and using large ensembles on massive data*, in Proceedings of the 11th IEEE International Conference on Data Mining (ICDM), IEEE, Washington, DC, 2011, pp. 41–50.

[3] P. BERTAIL, *Second-order properties of an extrapolated bootstrap without replacement under weak assumptions*, Bernoulli, 3 (1997), pp. 149–179.

[4] P. BERTAIL AND D. N. POLITIS, *Extrapolation of subsampling distribution estimators: The i.i.d. and strong mixing cases*, Canad. J. Statist., 29 (2001), pp. 667–680.

[5] G. BIAU, *Analysis of a random forests model*, J. Mach. Learn. Res., 13 (2012), pp. 1063–1095.

[6] G. BIAU, L. DEVROYE, AND G. LUGOSI, *Consistency of random forests and other averaging classifiers*, J. Mach. Learn. Res., 9 (2008), pp. 2015–2033.

[7] P. J. BICKEL AND A. SAKOV, *Extrapolation and the bootstrap*, Sankhyā Ser. A, 64 (2002), pp. 640–652.

[8] P. J. BICKEL AND J. A. YAHAV, *Richardson extrapolation and the bootstrap*, J. Amer. Statist. Assoc., 83 (1988), pp. 387–393.

[9] R. BLASER AND P. FRYZLEWICZ, *Random rotation ensembles*, J. Mach. Learn. Res., 17 (2016), pp. 126–151.

[10] L. BREIMAN, *Bagging predictors*, Mach. Learn., 24 (1996), pp. 123–140.

[11] L. BREIMAN, *Random forests*, Mach. Learn., 45 (2001), pp. 5–32.

[12] L. BREIMAN, J. FRIEDMAN, C. J. STONE, AND R. A. OLSHEN, *Classification and Regression Trees*, CRC, Boca Raton, FL, 1984.

[13] C. BREZINSKI AND M. R. ZAGLIA, *Extrapolation Methods: Theory and Practice*, Elsevier, New York, 2013.

[14] P. BÜHLMANN AND B. YU, *Analyzing bagging*, Ann. Statist., 30 (2002), pp. 927–961.

[15] T. I. CANNINGS AND R. J. SAMWORTH, *Random-projection ensemble classification (with discussion)*, J. R. Stat. Soc. Ser. B Stat. Methodol., 79 (2017), pp. 959–1035.

[16] V. CHERNOZHUKOV, D. CHETVERIKOV, K. KATO, AND Y. KOIKE, *Improved Central Limit Theorem and Bootstrap Approximations in High Dimensions*, preprint, https://arxiv.org/abs/1912.10529, 2019.

[17] H. A. CHIPMAN, E. I. GEORGE, AND R. E. MCCULLOCH, *BART: Bayesian additive regression trees*, Ann. Appl. Stat., 4 (2010), pp. 266–298.

[18] R. DÍAZ-URIARTE AND S. A. DE ANDRES, *Gene selection and classification of microarray data using random forest*, BMC Bioinformatics, 7 (2006), 3.

[19] D. DUA AND C. GRAFF, *UCI Machine Learning Repository*, http://archive.ics.uci.edu/ml, 2017.

[20] Y. FANG, *Scalable statistical inference for averaged implicit stochastic gradient descent*, Scand. J. Stat., 46 (2019), pp. 987–1002.

[21] Y. FANG, J. XU, AND L. YANG, *Online bootstrap confidence intervals for the stochastic gradient descent estimator*, J. Mach. Learn. Res., 19 (2018), pp. 3053–3073.

[22] J. FRIEDMAN, T. HASTIE, AND R. TIBSHIRANI, *The Elements of Statistical Learning*, Springer, New York, 2001.

[23] R. GENUER, J.-M. POGGI, AND C. TULEAU-MALOT, *Variable selection using random forests*, Pattern Recognition Lett., 31 (2010), pp. 2225–2236.

[24] R. GENUER, J.-M. POGGI, AND C. TULEAU-MALOT, *VSURF: An R package for variable selection using random forests*, The R Journal, 7 (2015), pp. 19–33.

[25] A. GÉRON, *Hands-on Machine Learning with Scikit-Learn and TensorFlow*, O'Reilly Media, https://github.com/ageron/handson-ml/tree/master/datasets, 2017.

[26] R. L. GRAHAM, D. E. KNUTH, O. PATASHNIK, AND S. LIU, *Concrete Mathematics: A Foundation for Computer Science*, Addison & Wesley, Reading, MA, 1994.

[27] B. GREGORUTTI, B. MICHEL, AND P. SAINT-PIERRE, *Correlation and variable importance in random forests*, Stat. Comput., 27 (2017), pp. 659–678.

[28] P. HALL AND R. J. SAMWORTH, *Properties of bagged nearest neighbour classifiers*, J. R. Stat. Soc. Ser. B Stat. Methodol., 67 (2005), pp. 363–379.

[29] D. HERNÁNDEZ-LOBATO, G. MARTÍNEZ-MUÑOZ, AND A. SUÁREZ, *How large should ensembles of classifiers be?*, Pattern Recognition, 46 (2013), pp. 1323–1336.

[30] H. ISHWARAN, *Variable importance in binary regression trees and forests*, Electron. J. Stat., 1 (2007), pp. 519–537.

[31] P. LATINNE, O. DEBEIR, AND C. DECAESTECKER, *Limiting the number of trees in random forests*, in Multiple Classifier Systems, Springer, Berlin, 2001, pp. 178–187.

[32] T. LI, L. LIU, A. KYRILLIDIS, AND C. CARAMANIS, *Statistical inference using SGD*, in Thirty-Second AAAI Conference on Artificial Intelligence, AAAI, 2018, 16619.

[33] A. LIAW AND M. WIENER, *Classification and regression by randomForest*, R News, 2 (2002), pp. 18–22, https://cran.r-project.org/web/packages/randomForest/randomForest.pdf.

[34] Y. LIN AND Y. JEON, *Random forests and adaptive nearest neighbors*, J. Amer. Statist. Assoc., 101 (2006), pp. 578–590.

[35] M. E. LOPES, *Estimating a Sharp Convergence Bound for Randomized Ensembles*, preprint, https://arxiv.org/abs/1303.0727, 2016.

[36] M. E. LOPES, *Estimating the algorithmic variance of randomized ensembles via the bootstrap*, Ann. Statist., 47 (2019), pp. 1088–1112.

[37] M. E. LOPES, N. B. ERICHSON, AND M. W. MAHONEY, *Error estimation for sketched SVD via the bootstrap*, in Proceedings of the 37th International Conference on Machine Learning, Vienna, Austria, 2020, 2978.

[38] M. E. LOPES, Z. LIN, AND H.-G. MUELLER, *Bootstrapping max statistics in high dimensions: Near-parametric rates under weak variance decay and application to functional and multinomial data*, Ann. Statist., 48 (2020), pp. 1214–1229.

[39] M. E. LOPES, S. WANG, AND M. W. MAHONEY, *Error estimation for randomized least-squares algorithms via the bootstrap*, in Proceedings of the International Conference on Machine Learning, Stockholm, Sweden, 2018, pp. 3223–3232.

[40] M. E. LOPES, S. WANG, AND M. W. MAHONEY, *A bootstrap method for error estimation in randomized matrix multiplication*, J. Mach. Learn. Res., 20 (2019), pp. 1–40.

[41] G. LOUPPE, L. WEHENKEL, A. SUTERA, AND P. GEURTS, *Understanding variable importances in forests of randomized trees*, in Advances in Neural Information Processing Systems, NeurIPS, San Diego, CA, 2013, pp. 431–439.

[42] P.-G. MARTINSSON AND J. TROPP, *Randomized numerical linear algebra: Foundations & algorithms*, Acta Numer., to appear; preprint, https://arxiv.org/abs/2002.01387, 2020.

[43] L. MENTCH AND G. HOOKER, *Quantifying uncertainty in random forests via confidence intervals and hypothesis tests*, J. Mach. Learn. Res., 17 (2016), pp. 1–41.

[44] J. MOULT, K. FIDELIS, A. KRYSHTAFOVYCH, AND A. TRAMONTANO, *Critical assessment of methods of protein structure prediction (CASP)—round* IX, Proteins Struct. Funct. Bioinform., 79 (2011), pp. 1–5.

[45] S. NEMBRINI, I. R. KÖNIG, AND M. N. WRIGHT, *The revival of the Gini importance?*, Bioinformatics, 34 (2018), pp. 3711–3718.

[46] A. Y. NG AND M. I. JORDAN, *Convergence rates of the voting Gibbs classifier, with application to Bayesian feature selection*, in Proceedings of the International Conference on Machine Learning, Williamstown, MA, 2001, pp. 377–384.

[47] K. K. NICODEMUS, J. D. MALLEY, C. STROBL, AND A. ZIEGLER, *The behaviour of random forest permutation-based variable importance measures under predictor correlation*, BMC Bioinformatics, 11 (2010), 110.

[48] T. M. OSHIRO, P. S. PEREZ, AND J. A. BARANAUSKAS, *How many trees in a random forest?*, in Machine Learning and Data Mining in Pattern Recognition, Springer, New York, 2012, pp. 154–168.

[49] P. PROBST AND A.-L. BOULESTEIX, *To tune or not to tune the number of trees in random forest*, J. Mach. Learn. Res., 18 (2018), pp. 1–18.

[50] R. J. SAMWORTH, *Optimal weighted nearest neighbour classifiers*, Ann. Statist., 40 (2012), pp. 2733–2763.

[51] M. SANDRI AND P. ZUCCOLOTTO, *A bias correction algorithm for the Gini variable importance measure in classification trees*, J. Comput. Graph. Statist., 17 (2008), pp. 611–628.

[52] A. SCHWING, C. ZACH, Y. ZHENG, AND M. POLLEFEYS, *Adaptive random forest—How many "experts" to ask before making a decision?*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, Washington, DC, 2011, pp. 1377–1384.

[53] E. SCORNET, *On the asymptotics of random forests*, J. Multivariate Anal., 146 (2016), pp. 72–83.

[54] E. SCORNET, *Random forests and kernel methods*, IEEE Trans. Inform. Theory, 62 (2016), pp. 1485–1500.

[55] E. SCORNET, G. BIAU, AND J.-P. VERT, *Consistency of random forests*, Ann. Statist., 43 (2015), pp. 1716–1741.

[56] J. SEXTON AND P. LAAKE, *Standard errors for bagged and random forest estimators*, Comput. Statist. Data Anal., 53 (2009), pp. 801–811.

[57] A. SIDI, *Practical Extrapolation Methods: Theory and Applications*, Cambridge University Press, Cambridge, UK, 2003.

[58] C. STROBL, A.-L. BOULESTEIX, T. KNEIB, T. AUGUSTIN, AND A. ZEILEIS, *Conditional variable importance for random forests*, BMC Bioinformatics, 9 (2008), 307.

[59] C. STROBL, A.-L. BOULESTEIX, A. ZEILEIS, AND T. HOTHORN, *Bias in random forest variable importance measures: Illustrations, sources and a solution*, BMC Bioinformatics, 8 (2007), 25.

[60] W. J. SU AND Y. ZHU, *Uncertainty Quantification for Online Learning and Stochastic Approximation via Hierarchical Incremental Gradient Descent*, preprint, https://arxiv.org/abs/1802.04876, 2018.

[61] M. TALAGRAND, *Isoperimetry and integrability of the sum of independent Banach-space valued random variables*, Ann. Probab., 17 (1989), pp. 1546–1570.

[62] S. WAGER, T. HASTIE, AND B. EFRON, *Confidence intervals for random forests: The jackknife and the infinitesimal jackknife*, J. Mach. Learn. Res., 15 (2014), pp. 1625–1651.

[63] H. WICKHAM, *ggplot2: Elegant Graphics for Data Analysis*, 2nd ed., Springer, New York, 2016.

[64] F. ZHOU, Q. CLAIRE, AND R. D. KING, *Predicting the geographical origin of music*, in Proceedings of the 2014 IEEE International Conference on Data Mining, IEEE, Washington, DC, 2014, pp. 1115–1120.