

Smoothing Effects of Bagging

Andreas Buja¹ Werner Stuetzle²

October 19, 2000

Abstract

Bagging is a device intended for reducing the prediction error of learning algorithms. In its simplest form, bagging draws bootstrap samples from the training sample, applies the learning algorithm to each bootstrap sample, and then averages the resulting prediction rules.

We extend the definition of bagging from statistics to statistical functionals and study the von Mises expansion of bagged statistical functionals. We show that the expansion is related to the Stein-Efron ANOVA expansion of the raw (unbagged) functional. The basic observation is that a bagged functional is always smooth in the sense that the von Mises expansion exists and is finite of length $1 + \text{resample size } M$. This holds even if the raw functional is rough or unstable. The resample size M acts as a smoothing parameter, where a smaller M means more smoothing.

¹AT&T Labs–Research, 180 Park Ave, Florham Park, NJ 07932-0971; andreas@research.att.com

²Department of Statistics, University of Washington, Seattle, WA 98195-4322; wxs@stat.washington.edu. Research partially supported by NSF grant DMS - 9803226. This work was performed while the second author was on sabbatical leave at AT&T Labs.

1 Notations, Definitions and Assumptions for Bagging Statistical Functionals

We need some standard notations and assumptions in order to define bagging for statistics and, more generally, for statistical functionals.

Let θ be a real-valued statistical functional $\theta(F) : \mathcal{P} \rightarrow \mathbb{R}$ defined on a subset \mathcal{P} of the probability measures on a given sample space. By assumption all empirical measures $F_M = \frac{1}{M} \sum_{i=1}^M \delta_{x_i}$ are contained in \mathcal{P} . If θ is evaluated at an empirical measure, it specializes to a statistic which we write as $\theta(F_M) = \theta(x_1, \dots, x_M)$. This is a permutation symmetric function of the M sample points. We will repeatedly need expectations of random variables $\theta(X_1, \dots, X_M)$ where X_1, \dots, X_M are i.i.d. according to some F : $\mathbf{E}_F \theta(X_1, \dots, X_M) = \int \theta(x_1, \dots, x_M) dF(x_1) \cdots dF(x_M)$.

Following Breiman (1996), we define bagging of a statistic $\theta(F_N)$ as the average over bootstrap samples X_1^*, \dots, X_N^* drawn i.i.d. from F_N :

$$\theta^B(F_N) = \mathbf{E}_{F_N} \theta(X_1^*, \dots, X_N^*) .$$

For our purposes we need to generalize the notion of bagging to statistical functionals $\theta(F)$. First, we need to divorce the resample size from the sample size N (compare Friedman and Hall 2000). To this end, we allow the number M of resamples to be drawn from F_N to be arbitrary:

$$\theta_M^B(F_N) = \mathbf{E}_{F_N} \theta(X_1^*, \dots, X_M^*) .$$

This separation of M and N allows us to extend the definition of bagging from empirical measures F_N to arbitrary distributions:

$$\theta_M^B(F) = \mathbf{E}_F \theta(X_1^*, \dots, X_M^*) ,$$

where the random variables X_1^*, \dots, X_M^* are i.i.d. F , and their number M is merely a parameter of the bagging procedure. Unlike for an empirical distribution of an actual sample, for a general probability measure F there is no notion of sample size. The variables X_i^* should still be thought of as bootstrap samples, albeit drawn from an “infinite population”.

Since the resample size M now denotes a parameter of the bagging procedure, we need to distinguish it from the size N of actual data x_1, \dots, x_N . If one models the data as i.i.d. samples from F , one estimates F with the empirical $F_N = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}$. The functional $\theta(F)$ is then estimated by plug-in with the statistic $\theta(F_N)$:

$$\hat{\theta}(F) = \theta(F_N) .$$

The bagged functional $\theta_M^B(F)$ in turn is estimated with the plug-in estimator $\theta_M^B(F_N)$:

$$\hat{\theta}_M^B(F) = \theta_M^B(F_N) = \mathbf{E}_{F_N} \theta(X_1^*, \dots, X_M^*) .$$

The idea of bagging is to smooth θ , with the number M playing the role of a smoothing parameter. It is not a priori clear, though, whether more smoothing occurs for small M or large M . Here is an intuition that proves to be correct: bagging averages over empiricals F_M , hence more smoothing occurs when F_M is allowed to roam further from F , effectively using a larger neighborhood (“bandwidth”) around F ; due to $F_M \rightarrow F$ as $M \rightarrow \infty$, F_M is on the average closer to F for large M , hence the “bandwidth” is larger for *small* M . The calculations below verify that this is so, but curiously the reason has nothing to do with F_M being close to, or far from, F : it turns out that the von Mises expansion of an M -bagged functional is finite of length M ; because the von Mises expansion is essentially a Taylor expansion, the m -bagged functional is smoother if the expansion is shorter, that is, if M is smaller.

The above definition of a bagged statistical functional has a blind spot: It would be interesting to consider both bootstrap sampling with replacement (conventional) and bootstrap sampling without replacement where M is strictly smaller than N (as in Friedman and Hall (2000) and Buhlmann and Yu (2000)). If bootstrap is extended to infinite populations, however, the difference between sampling with and without replacement disappears. Thus, in order to capture both modes of sampling, one is limited to finite populations and correspondingly to statistics as opposed to statistical functionals.

If bagging is smoothing by averaging over nearby empirical distributions, one may wonder whether other types of bagging could be conceived. In fact, one can more generally define a smoothed version θ^S of θ by

$$\theta^S(F) = \mathbf{ave} \{ \theta(G) \mid G \in \mathcal{N}(F) \} ,$$

where $\mathcal{N}(F)$ is some sort of neighborhood of F , and **ave** denotes some way of averaging. This suggests a number of generalizations of bagging, for example by varying the neighborhoods and the meaning of **ave**. In the present note, however, we remain with Breiman’s original version of bagging and pursue some implications of averaging over empirical distributions.

2 Preliminaries 1: The von Mises Expansion of a Statistical Functional

The von Mises expansion of a functional θ around a distribution F is an expansion of the form

$$\begin{aligned} \theta(G) &= \theta(F) + \int \psi_1(x) d(G - F)(x) + \frac{1}{2} \int \psi_2(x_1, x_2) d(G - F)^{\otimes 2} + \dots \\ &= \theta(F) + \sum_{k=1}^{\infty} \frac{1}{k!} \int \psi_k(x_1, \dots, x_k) d(G - F)^{\otimes k} . \end{aligned}$$

It can be interpreted as the Taylor expansion of $\theta((1-s)F + sG) = \theta(F + s(G - F))$ evaluated at $t = 1$. The first term in the sum is a linear functional, the second term is a quadratic functional, etc. There is of course no guarantee that the expansion exists. Reeds (1976) gives a discussion of conditions under which this expansion is meaningful in terms of remainders and convergence. We are not concerned with technical difficulties because the expansions we encounter below are finite and exact.

The functions ψ_k are not uniquely determined. We can choose them such that all the integrals w.r.t. F vanish, i.e.,

$$\begin{aligned} 0 &= \int \psi_1(x) dF \\ 0 &= \int \psi_2(x_1, x_2) dF(x_1) = \int \psi_2(x_1, x_2) dF(x_2) , \end{aligned}$$

and so on. The von Mises expansion then simplifies to

$$\begin{aligned} \theta(G) &= \theta(F) + \mathbf{E}_G \psi_1(X) + \frac{1}{2} \mathbf{E}_G \psi_2(X_1, X_2) + \dots \\ &= \theta(F) + \sum_{k=1}^{\infty} \frac{1}{k!} \mathbf{E}_G \psi_k(X_1, \dots, X_k) . \end{aligned}$$

The function $\psi_1(x)$ is also known as the influence function of θ , but we will similarly call $\psi_k(x_1, \dots, x_k)$ the k 'th order influence function. Influence functions of any order are permutation symmetric in their arguments.

Assuming sufficient smoothness of the functional, ψ_k can be obtained by differentiation:

$$\psi_k(x_1, \dots, x_k) = \left. \frac{d}{ds_1} \right|_{s_1=0} \dots \left. \frac{d}{ds_k} \right|_{s_k=0} \theta((1 - \sum s_i)F + \sum s_i \delta_{x_i}) .$$

3 Preliminaries 2: The ANOVA Expansion of a Statistic

Efron and Stein (1981) introduced an ANOVA-type expansion for statistics that are functions of independent random variables X_1, \dots, X_M . Because we are only interested in symmetric functions of i.i.d. data as they arise from evaluating statistical functionals on empirical distributions, we introduce a simplified version of the expansion as follows. Define partial expectations by

$$\begin{aligned} \mu_0 &= \mathbf{E}_F \theta(X_1, \dots, X_M) \\ \mu_1(x_1) &= \mathbf{E}_F \theta(x_1, X_2, \dots, X_M) \\ \mu_2(x_1, x_2) &= \mathbf{E}_F \theta(x_1, x_2, X_3, \dots, X_M) \end{aligned}$$

$$\begin{aligned}
& \dots \\
\mu_k(x_1, \dots, x_k) &= E_F \theta(x_1, \dots, x_k, X_{k+1}, \dots, X_M) \\
& \dots \\
\mu_M(x_1, \dots, x_M) &= \theta(x_1, \dots, x_M) .
\end{aligned}$$

Permutation symmetry of $\theta(x_1, \dots, x_M)$ implies that the free arguments x_j could be in any position, a fact that will be used extensively below.

Define ANOVA terms by

$$\begin{aligned}
\alpha_0 &= \mu_0 \\
\alpha_1(x_1) &= \mu_1(x_1) - \mu_0 \\
\alpha_2(x_1, x_2) &= \mu_2(x_1, x_2) - \mu_1(x_1) - \mu_1(x_2) + \mu_0 \\
& \dots \\
\alpha_k(x_1, \dots, x_k) &= \sum_{\nu=0}^k (-1)^{k-\nu} \sum_{1 \leq i_1 < \dots < i_\nu \leq k} \mu_\nu(x_{i_1}, \dots, x_{i_\nu}) \\
& \dots \\
\alpha_M(x_1, \dots, x_M) &= \sum_{\nu=0}^M (-1)^{M-\nu} \sum_{1 \leq i_1 < \dots < i_\nu \leq M} \mu_\nu(x_{i_1}, \dots, x_{i_\nu}) .
\end{aligned}$$

Then the ANOVA expansion of $\theta(x_1, \dots, x_M)$ is:

$$\begin{aligned}
\theta(x_1, \dots, x_M) &= \alpha_0 + \sum_{j=1}^M \alpha_1(x_j) + \sum_{1 \leq j_1 < j_2 \leq M} \alpha_2(x_{j_1}, x_{j_2}) + \dots \\
&= \sum_{k=0}^M \sum_{1 \leq j_1 < \dots < j_k \leq M} \alpha_k(x_{j_1}, \dots, x_{j_k}) .
\end{aligned}$$

This expansion is tautological and holds without assumptions other than permutation symmetry of $\theta(x_1, \dots, x_M)$ in its arguments. The proof is by showing that the partial expectations implicit in the ANOVA terms cancel each other except for $\mu_M = \theta(x_1, \dots, x_M)$.

If one assumes that the variables X_1, \dots, X_M are i.i.d., then the terms α_k have vanishing marginals in all arguments:

$$\mathbf{E}_F \alpha_k(x_1, \dots, x_{j-1}, X_j, x_{j+1}, \dots, x_k) = 0 .$$

As a consequence, all terms in the ANOVA expansion are pairwise uncorrelated.

Note that all functions μ_k and α_k are implicitly dependent on M because they derive from a statistic of M arguments, $\theta(x_1, \dots, x_M)$. If necessary we make the dependence explicit by writing μ_k^M and α_k^M . By contrast, the influence functions ψ_k in the von Mises expansion are independent of any sample size because this expansion is centered at F as opposed to F_M .

The zero'th term $\alpha_0^M = \mu_0^M$ is also called the “grand mean”, and the first term $\alpha_1^M(x)$ the “main effect” function. Correspondingly we call $\alpha_k^M(x_1, \dots, x_k)$ the k 'th order “interaction” function of $\theta(x_1, \dots, x_M)$.

4 A Warm-Up Exercise: The First Order Influence Function of a Bagged Functional

Before deriving a general formula for the terms of the von Mises expansion of θ_M^B , we calculate the linear term to illustrate the idea. The influence function will be denoted $\psi_1^B(x)$ as a reminder that it belongs to the bagged functional:

$$\begin{aligned}\psi_1^B(x) &= \left. \frac{d}{ds} \right|_{s=0} \theta_M^B((1-s)F + s\delta_x) \\ &= \left. \frac{d}{ds} \right|_{s=0} \mathbf{E}_{(1-s)F + s\delta_x} \theta(X_1, \dots, X_M) .\end{aligned}$$

The expectation $\mathbf{E}_{(1-s)F + s\delta_x} \theta(X_1, \dots, X_M)$ is effectively a polynomial of degree M in s and hence arbitrarily differentiable. We expand it by applying the mixture $(1-s)F + s\delta_x$ to each argument X_i , resulting in 2^M terms. They in turn can be bundled according to the number of times δ_x occurs:

$$\begin{aligned}& \mathbf{E}_{(1-s)F + s\delta_x} \theta(X_1, \dots, X_M) \\ &= (1-s)^M \mathbf{E}_F \theta(X_1, \dots, X_M) \\ &+ (1-s)^{M-1} s M \mathbf{E}_F \theta(x, X_2, \dots, X_M) \\ &+ (1-s)^{M-2} s^2 \frac{M(M-1)}{2} \mathbf{E}_F \theta(x, x, X_3, \dots, X_M) \\ &+ O(s^3) .\end{aligned}$$

Also used was permutation symmetry which implies, for example,

$$\mathbf{E}_F \theta(\dots, X_{j-1}, x, X_{j+1}, \dots) = \mathbf{E}_F \theta(x, X_2, \dots, X_M) .$$

As we differentiate w.r.t. s at $s = 0$, only the first two terms make a contribution:

$$\psi_1^B(x) = M [-\mathbf{E}_F \theta(X_1, \dots, X_M) + \mathbf{E}_F \theta(X_1, \dots, X_{M-1}, x)] = M \alpha_1^M(x) ,$$

where as above α_1^M is the main effects function in the ANOVA expansion of $\theta(F_M)$, which is the raw, not the bagged, statistic.

Suppose we have an i.i.d. sample x_1, \dots, x_N of size N from F with empirical distribution $F_N = \frac{1}{N} \sum \delta_{x_i}$. The first order von Mises approximation to $\theta_M^B(F) = \theta_M^B(F_N)$ is

$$\theta_M^B(F_N) \approx \theta_M^B(F) + \frac{1}{N} \sum_{i=1}^N \psi_1^B(x_i) = \mu_0^M + \frac{M}{N} \sum_{i=1}^N \alpha_1^M(x_i) .$$

For the special case $M = N$ this is exactly the grand mean and the main effects in the ANOVA expansion of $\theta(F_N)$.

5 The von Mises Expansion of a Bagged Functional

We now calculate the k -th order influence function. To this end let

$$\tilde{F}_k = (1 - \sum_1^k s_i)F + \sum_1^k s_i \delta_{x_i} .$$

By definition,

$$\psi_k^B(x_1, \dots, x_k) = \frac{\partial^k}{\partial s_1 \dots \partial s_k} \Big|_{s_1, \dots, s_k=0} \theta_M^B(\tilde{F}_k) .$$

Again we note that $\theta_M^B(\tilde{F}_k) = \mathbf{E}_{\tilde{F}_k} \theta(X_1, \dots, X_M)$ is effectively a polynomial of degree M in s . Expanding it into $(k+1)^M$ summands, bundling the summands according to the number of δ_{x_i} 's they contain, and using permutation symmetry, we get:

$$\begin{aligned} \theta_M^B(\tilde{F}_k) &= \mathbf{E}_{\tilde{F}_k} \theta(X_1, \dots, X_M) \\ &= (1 - \sum_{i=1}^k s_i)^M \mathbf{E}_F \theta(X_1, \dots, X_M) \\ &\quad + \sum_{j=1}^k (1 - \sum_{i=1}^k s_i)^{M-1} s_j M \mathbf{E}_F \theta(x_j, X_2, \dots, X_M) \\ &\quad + \sum_{1 \leq j_1 < j_2 \leq k} (1 - \sum_{i=1}^k s_i)^{M-2} s_{j_1} s_{j_2} M(M-1) \mathbf{E}_F \theta(x_{j_1}, x_{j_2}, X_3, \dots, X_M) \\ &\quad + \dots \\ &\quad + O(s_1^2, \dots, s_k^2) \end{aligned}$$

Terms containing a second or higher power of any s_j have vanishing derivatives at zero and hence will disappear in what follows. This is why the summation on the fourth line can run over index pairs $j_1 \neq j_2$ only, the omitted summands being summarily lumped into $O(s_1^2, \dots, s_k^2)$. Thus, with the abbreviated notation for partial expectations:

$$\begin{aligned} \theta_M^B(\tilde{F}_k) &= \sum_{\nu=0}^{\min(k,M)} \sum_{1 \leq j_1 < \dots < j_\nu \leq k} (1 - \sum_{i=1}^k s_i)^{M-\nu} s_{j_1} \dots s_{j_\nu} \frac{M!}{(M-\nu)!} \mu_\nu^M(x_{j_1}, \dots, x_{j_\nu}) \\ &\quad + O(s_1^2, \dots, s_k^2) . \end{aligned}$$

Note that the outer sum extends to $\min(k, M)$ only. As the derivatives can be pulled inside the double sum, we have to calculate

$$\frac{\partial^k}{\partial s_1 \cdots \partial s_k} \Big|_{s_1, \dots, s_k=0} \left[\left(1 - \sum_{i=1}^k s_i\right)^{M-\nu} s_{j_1} \cdots s_{j_\nu} \right] .$$

We first take partial derivatives w.r.t. $s_{j_1}, \dots, s_{j_\nu}$ in turn:

$$\begin{aligned} & \frac{\partial}{\partial s_{j_1}} \Big|_{s_{j_1}=0} \left[\left(1 - \sum s_i\right)^{M-\nu} s_{j_1} \cdots s_{j_\nu} \right] \\ &= \left[(M-\nu)(1 - \sum s_i)^{M-\nu-1} (-1) s_{j_1} \cdots s_{j_\nu} + (1 - \sum s_i)^{M-\nu} s_{j_2} \cdots s_{j_\nu} \right] \Big|_{s_{j_1}=0} \\ &= (1 - \sum s_i)^{M-\nu} s_{j_2} \cdots s_{j_\nu} . \end{aligned}$$

Repeating this process we obtain:

$$\frac{\partial^\nu}{\partial s_{j_1} \cdots \partial s_{j_\nu}} \Big|_{s_1, \dots, s_k=0} \left[\left(1 - \sum s_i\right)^{M-\nu} s_{j_1} \cdots s_{j_\nu} \right] = (1 - \sum s_i)^{M-\nu} .$$

We still have to take the derivatives w.r.t. indices not among j_1, \dots, j_ν . Pick one such index l :

$$\frac{\partial}{\partial s_l} \Big|_{s_l=0} \left[(1 - \sum s_i)^{M-\nu} \right] = (M-\nu)(1 - \sum s_i)^{M-\nu-1} (-1)$$

Repeating for all such l we get:

$$\begin{aligned} & \frac{\partial^k}{\partial s_1 \cdots \partial s_k} \Big|_{s_1, \dots, s_k=0} \left[\left(1 - \sum s_i\right)^{M-\nu} s_{j_1} \cdots s_{j_\nu} \right] \\ &= \begin{cases} (M-\nu)(M-\nu-1) \cdots (M-k+1) (-1)^{k-\nu} & = \frac{(M-\nu)!}{(M-k)!} (-1)^{k-\nu} \quad \text{for } k \leq M , \\ 0 & \text{for } k > M . \end{cases} \end{aligned}$$

Putting everything together, we get first of all

$$\psi_k^B(x_1, \dots, x_k) = 0 \quad \text{for } k > M .$$

For $k \leq M$ we get

$$\begin{aligned} \psi_k^B(x_1, \dots, x_k) &= \frac{M!}{(M-k)!} \sum_{\nu=0}^k (-1)^{k-\nu} \sum_{1 \leq j_1 < \dots < j_\nu \leq k} \mu_\nu^M(x_{j_1}, \dots, x_{j_\nu}) \\ &= \frac{M!}{(M-k)!} \alpha_k^M(x_1, \dots, x_k) \end{aligned}$$

We summarize:

Theorem: *The k 'th order influence function ψ_k^B of an M -bagged functional $\theta_M^B(F)$ is proportional to the k 'th order interaction function α_k^M of the statistic $\theta(F_M)$:*

$$\psi_k^B(x_1, \dots, x_k) = \begin{cases} \frac{M!}{(M-k)!} \alpha_k^M(x_1, \dots, x_k) & \text{for } k \leq M, \\ 0 & \text{for } k > M. \end{cases}$$

It is now a simple matter to write down the full von Mises expansion of an M -bagged functional:

$$\begin{aligned} \theta_M^B(G) &= \theta_M^B(F) + \sum_{k \geq 1} \frac{1}{k!} \mathbf{E}_G \psi_k(X_1, \dots, X_k) \\ &= \alpha_0^M + \sum_{k=1}^M \binom{M}{k} \mathbf{E}_G \alpha_k^M(X_1, \dots, X_k). \end{aligned}$$

Again we summarize:

Theorem: *Bagged functionals are smooth in the sense that the von Mises expansion exists and is of finite length M :*

$$\theta_M^B(G) = \sum_{k=0}^M \binom{M}{k} \mathbf{E}_G \alpha_k^M(X_1, \dots, X_k).$$

Since the von Mises expansion is effectively a Taylor expansion, it would seem natural for exact finite expansions to consider the length as an inverse measure of smoothness: the shorter the expansion the smoother the functional. With this interpretation and in light of the theorem, *bagging performs more smoothing for smaller M .*

Suppose now we have an i.i.d. sample y_1, \dots, y_N of size N from the distribution F . The von Mises expansion of θ_M^B at $F_N = \frac{1}{N} \sum_{j=1}^N \delta_{y_j}$ centered at F is:

$$\theta_M^B(F_N) = \sum_{k=0}^M \binom{M}{k} \frac{1}{N^k} \sum_{1 \leq j_1, \dots, j_k \leq N} \alpha_k^M(y_{j_1}, \dots, y_{j_k}).$$

Note that the inner sum is unconstrained. The bagging parameter M is unconstrained w.r.t. the sample size N : M can be chosen to be smaller or larger than N , which raises the question of criteria for choosing among values for M . This is then just another form of the problem of smoothing parameter selection.

For the conventional choice $M = N$ one obtains an interesting comparison with the ANOVA expansion of $\theta(F_N)$:

Theorem: *The terms in the von Mises expansion of the conventional N -bagged statistic $\theta_N^B(F_N)$ form a superset of the terms in the ANOVA expansion of $\theta(F_N)$.*

$$\begin{aligned}\theta_N^B(F_N) &= \sum_{k=0}^N \binom{N}{k} \frac{1}{N^k} \sum_{1 \leq j_1, \dots, j_k \leq N} \alpha_k^N(y_{j_1}, \dots, y_{j_k}) , \\ \theta(F_N) &= \sum_{k=0}^N \sum_{1 \leq j_1 < \dots < j_k \leq N} \alpha_k^N(y_{j_1}, \dots, y_{j_k}) .\end{aligned}$$

The inner sums of the first and the second line have N^k and $\binom{N}{k}$ terms, respectively, the difference being that the first inner sum runs over unconstrained indices, the second over strictly ordered indices. The ratio $\binom{N}{k}/N^k$ downweights the inner sum in the first line to match the smaller number of terms in the second line.

The difference between the raw and the N -bagged statistic is that the latter includes “diagonal” terms such as $\alpha_2^N(y_1, y_1)$, arising from sampling with replacement in the bootstrap procedure. By comparison sampling without replacement amounts to a mere permutation of the data and hence leaves the value of a permutation symmetric statistic unchanged.

References

- [1] L. Breiman. Bagging predictors. *Machine Learning*, 26:123–140, 1996.
- [2] Peter Buhlmann and Bin Yu. Explaining bagging. Can be downloaded from <http://stat.ethz.ch/~buhlmann/bibliog.html>, September 2000.
- [3] J.H. Friedman and O. Hall. On bagging and nonlinear estimation. Can be downloaded from <http://www-stat.stanford.edu/~jhf/#reports>, May 2000.