

Random forests and averaging classifiers

Gábor Lugosi

ICREA and Pompeu Fabra University

Barcelona

joint work with

Gérard Biau (Paris 6)

Luc Devroye (McGill, Montreal)

Leo Breiman



Binary classification

$(\mathbf{X}, \mathbf{Y}) \in \mathbb{R}^d \times \{-1, 1\}$ observation/label pair

$\mathbf{g}_n(\mathbf{X}) = \mathbf{g}_n(\mathbf{X}, \mathbf{D}_n) \in \{0, 1\}$ classifier, based on

$\mathbf{D}_n = (\mathbf{X}_1, \mathbf{Y}_1), \dots, (\mathbf{X}_n, \mathbf{Y}_n)$, i.i.d. training data, distributed as (\mathbf{X}, \mathbf{Y}) .

$L(\mathbf{g}_n) = \mathbb{P}\{\mathbf{g}_n(\mathbf{X}) \neq \mathbf{Y} | \mathbf{D}_n\}$ loss of \mathbf{g}_n .

a posteriori probability $\eta(\mathbf{x}) = \mathbb{P}\{\mathbf{Y} = 1 | \mathbf{X} = \mathbf{x}\}$.

Bayes classifier: $\mathbf{g}^*(\mathbf{x}) = \mathbb{1}_{\eta(\mathbf{x}) \geq 1/2}$.

Bayes risk: $L^* = L(\mathbf{g}^*)$.

$\{\mathbf{g}_n\}$ is *consistent* if $L(\mathbf{g}_n) \rightarrow L^*$ in probability.

Local averaging

Historically the first non-parametric classification rules.

Histogram, k -nearest neighbor, kernel classifiers.

Fix and Hodges (1951-52),

Cover and Hart (1967),

Glick (1973),

Devroye and Wagner (1976),

Stone (1977),

Gordon and Olshen (1978),

Devroye and Györfi (1983).

Stone's 1977 theorem

Local averaging classifiers:

$$g_n(\mathbf{x}) = 1 \quad \text{iff} \quad \sum_{i=1}^n Y_i W_{ni}(\mathbf{x}) \geq 0$$

where $W_{ni}(\mathbf{x}) = W_{ni}(\mathbf{x}, \mathbf{X}_1, \dots, \mathbf{X}_n) \geq 0$

and $\sum_{i=1}^n W_{ni}(\mathbf{x}) = 1$.

Stone's 1977 theorem

Consistency holds if

(i) $\lim_{n \rightarrow \infty} \mathbb{E} \left\{ \max_{1 \leq i \leq n} W_{ni}(\mathbf{X}) \right\} = 0.$

(ii) For all $a > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{E} \left\{ \sum_{i=1}^n W_{ni}(\mathbf{X}) \mathbb{1}_{\|\mathbf{X}_i - \mathbf{X}\| > a} \right\} = 0.$$

(iii) There is a $c > 0$ such that, for every $f \geq 0$,

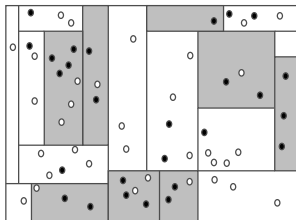
$$\mathbb{E} \left\{ \sum_{i=1}^n W_{ni}(\mathbf{X}) f(\mathbf{X}_i) \right\} \leq c \mathbb{E} f(\mathbf{X}).$$



Tree classifiers

Histograms based on data-dependent partitions.

Partition is constructed by recursive splitting.



See Breiman, Freedman, Olshen, and Stone (1984), Devroye, Györfi, and Lugosi (1996) for surveys.

Consistency of tree classifiers

Many versions suggested in the literature are inconsistent.

General consistency theorems:

Assume the partition depends on $\mathbf{X}_1, \dots, \mathbf{X}_n$ only. Let $\mathbf{A}(\mathbf{X})$ denote the cell containing \mathbf{X} and $\mathbf{N}(\mathbf{X}) = \sum_{i=1}^n \mathbb{1}_{\mathbf{X}_i \in \mathbf{A}(\mathbf{X})}$.

If $\text{diam}(\mathbf{A}(\mathbf{X})) \rightarrow \mathbf{0}$ and $\mathbf{N}(\mathbf{X}) \rightarrow \infty$ in probability then the classifier is consistent (Devroye, Györfi, and Lugosi, 1996).

For general partitions consistency holds under an additional combinatorial condition (Lugosi and Nobel, 1993).

Random forests

Tree classifiers are unstable.

Breiman (2001) suggests “bootstrap” randomization in building trees:

- (1) choose a cell at random
- (2) choose $m < d$ coordinates at random
- (3) cut at a point (and direction) giving the largest decrease in empirical error.

Repeat until every cell is pure.

Repeat the random tree classifier a zillion times and take majority vote.

Additional randomization is achieved by bootstrap sampling.

Simpler version

Breiman asked if this classifier was consistent. Performs well in practice.

Simpler version:

- (1) choose a cell at random
- (2) choose a coordinate at random
- (3) cut at a random point.

Repeat k times.

Repeat the random tree classifier a zillion times and take majority vote.

Local averaging rule with weights

$$\mathbf{W}_{ni}(\mathbf{X}) \sim \mathbb{P}_{\mathbf{Z}}\{\mathbf{X}, \mathbf{X}_i \text{ are in same cell}\}.$$

Averaged classifiers

Let $\mathbf{g}_n(\mathbf{X}, \mathbf{Z}, \mathbf{D}_n) = \mathbf{g}_n(\mathbf{X}, \mathbf{Z})$ be a randomized classifier.

Probability of error:

$$L(\mathbf{g}_n) = \mathbb{P}_{(\mathbf{X}, \mathbf{Y}), \mathbf{Z}} \{ \mathbf{g}_n(\mathbf{X}, \mathbf{Z}, \mathbf{D}_n) \neq \mathbf{Y} \} .$$

Averaged classifier: $\bar{\mathbf{g}}_n(\mathbf{x}) = \mathbb{1}_{\mathbb{E}_{\mathbf{Z}} \mathbf{g}_n(\mathbf{x}, \mathbf{Z}) \geq 1/2}$

Main lemma:

If \mathbf{g}_n is consistent then $\bar{\mathbf{g}}_n$ is also consistent.

Averaging “stabilizes.”

Consistency of simple version

We obtain consistency without computing the weights $\mathbf{W}_{ni}(\mathbf{X})$.

Assume \mathbf{X} is supported in $[0, 1]^d$.

Then $\bar{\mathbf{g}}_n$ is consistent whenever $\mathbf{k} \rightarrow \infty$ and $\mathbf{k}/n \rightarrow 0$ as $\mathbf{k} \rightarrow \infty$.

It suffices to prove consistency of the randomized “base” classifier.

It is enough to show $\text{diam}(\mathbf{A}(\mathbf{X}, \mathbf{Z})) \rightarrow 0$ and $\mathbf{N}(\mathbf{X}, \mathbf{Z}) \rightarrow \infty$ in probability.

Consistency of simple version

$\mathbf{N}(\mathbf{X}, \mathbf{Z}) \rightarrow \infty$ and $\text{diam}(\mathbf{A}(\mathbf{X}, \mathbf{Z})) \rightarrow \mathbf{0}$, in probability, are both easy to show.

Interestingly, for $\mathbf{d} > \mathbf{1}$, $\sup_{\mathbf{x}} \text{diam}(\mathbf{A}(\mathbf{x}, \mathbf{Z})) \not\rightarrow \mathbf{0}$.

If $\mathbf{d} \geq \mathbf{3}$, the number of cells with diameter $\mathbf{1}$ (in sup norm) is a supercritical branching process.

A scale invariant version

- (1) choose a cell at random
- (2) choose a coordinate at random
- (3) cut at a random data point.

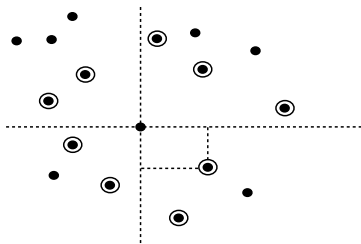
Repeat k times.

Repeat the random tree classifier a zillion times and take majority vote.

If the distribution of \mathbf{X} has non-atomic marginals in \mathbb{R}^d , then $\bar{\mathbf{g}}_n$ is consistent whenever $k \rightarrow \infty$ and $k/n \rightarrow 0$ as $k \rightarrow \infty$.

Breiman's original random forest

Lin and Jeon (2006) point out that any random forest classifier that cuts down to pure cells is a weighted *layered nearest neighbor rule*.



No such rule can be consistent if the distribution of \mathbf{X} is concentrated on a diagonal.

Randomizing inconsistent classifiers

Averaging consistent randomized classifiers preserves consistency.

The converse is not true: averaging inconsistent classifiers may lead to consistency.

This may be the case with Breiman's original random forest if \mathbf{X} has a density.

We work out a stylized example.

A randomized nearest neighbor rule

For $\mathbf{x} \in \mathbb{R}$, let $\mathbf{X}_{(1)}(\mathbf{x}), \mathbf{X}_{(2)}(\mathbf{x}), \dots, \mathbf{X}_{(n)}(\mathbf{x})$ be $\mathbf{X}_1, \dots, \mathbf{X}_n$ ordered according to distances to \mathbf{x} .

Let $\mathbf{U}_1, \dots, \mathbf{U}_n$ be i.i.d. uniform $[0, \mathbf{m}]$.

Let $\mathbf{g}_n(\mathbf{x}, \mathbf{Z}) = \mathbf{Y}_{(i)}(\mathbf{x})$ if and only if

$$\max(i, \mathbf{U}_i) \leq \max(j, \mathbf{U}_j) \quad \text{for } j = 1, \dots, n$$

$\mathbf{X}_{(i)}(\mathbf{x})$ is the *perturbed nearest neighbor* of \mathbf{x} .

$\bar{\mathbf{g}}_n(\mathbf{x}) = \mathbb{1}_{\mathbb{E}_{\mathbf{Z}} \mathbf{g}_n(\mathbf{x}, \mathbf{Z}) \geq 1/2}$ is the *averaged perturbed nearest neighbor classifier*.

Consistency

The averaged perturbed nearest neighbor classifier is consistent if $m \rightarrow \infty$ and $m/n \rightarrow 0$.

Proof: \bar{g}_n is a local averaging classifier with

$$\begin{aligned} W_{ni}(\mathbf{x}) &= \mathbb{P}_Z\{\mathbf{X}_{(i)}(\mathbf{x}) \text{ is the perturbed nearest neighbor of } \mathbf{x}\} \\ &= \dots \text{ can be written explicitly} \end{aligned}$$

Stone's theorem may be used.

Bagging

In bagging, suggested by Breiman (1996), bootstrap samples are generated from the original data set.

Let $\mathbf{q}_n \in [0, 1]$. In a bootstrap sample $\mathbf{D}_n(\mathbf{Z})$ each $(\mathbf{X}_i, \mathbf{Y}_i)$ is present with probability \mathbf{q}_n .

Given a classifiers $\{\mathbf{g}_n\}$, let

$$\mathbf{g}_n(\mathbf{X}, \mathbf{Z}, \mathbf{D}_n) = \mathbf{g}_n(\mathbf{X}, \mathbf{D}_n(\mathbf{Z})) ,$$

By drawing many bootstrap samples, one obtains the averaged classifier $\bar{\mathbf{g}}_n(\mathbf{x}, \mathbf{D}_n) = \mathbb{1}_{\mathbb{E}_{\mathbf{Z}} \mathbf{g}_n(\mathbf{x}, \mathbf{D}_n(\mathbf{Z})) \geq 1/2}$.

If $n\mathbf{q}_n \rightarrow \infty$ as $n \rightarrow \infty$ then the bagging classifier is consistent.

Bagging the 1-NN classifier

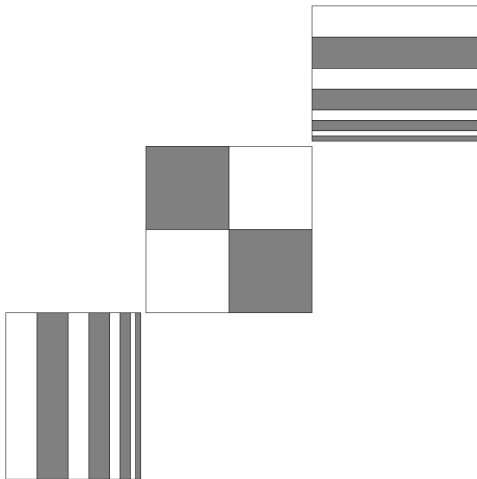
It may help to choose much smaller values of \mathbf{q}_n .

The **1**-nearest neighbor rule is not consistent unless either $\mathbf{L}^* = \mathbf{0}$ or $\mathbf{L}^* = \mathbf{1}/2$.

However, the bagging averaged **1**-nearest neighbor classifier is consistent for all distributions of (\mathbf{X}, \mathbf{Y}) if and only if $\mathbf{q}_n \rightarrow \mathbf{0}$ and $n\mathbf{q}_n \rightarrow \infty$.

Greedy trees

Greedy trees like Breiman's may be inconsistent for another reason:



Questions

Is Breiman's original random forest consistent if \mathbf{X} has a density?

In what situations does randomizing and averaging help?

Random forests

