

Lecture 5

Ensembles – Part 1

(February 6, 2015)

Mu Zhu
University of Waterloo

An Annus Mirabilis

L. Breiman (1996), “Bagging predictors”, *Machine Learning* **24**, pp. 123–140.

Y. Freund & R. Schapire (1996), “Experiments with a new boosting algorithm”, in *Proceedings of the 13th International Conference on Machine Learning*, pp. 148–156.

- idea (same): instead of one classifier, beneficial to **generate** many different classifiers (from the same training data) and **aggregate** over their classification results
- details (different): how to **generate**; how to **aggregate**

Breiman (1996)

- suppose $y_i = \pm 1$
 - for direct comparison with Freund & Schapire (1996)
- **generation**:
 - use weighted training data $\{(\mathbf{x}_i, y_i, w_i^{(b)})\}$ to learn $f_b(\mathbf{x})$

- **aggregation**:

$$F(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B f_b(\mathbf{x})$$

- **weights**: independently determined by **bootstrap**

called **Bagging** ... for bootstrapped **aggregating**

Freund & Schapire (1996)

- suppose $y_i = \pm 1$
- **generation:**
 - use weighted training data $\{(\mathbf{x}_i, y_i, w_i^{(b)})\}$ to learn $f_b(\mathbf{x})$
 - compute

$$\varepsilon_b = \frac{\sum_{i=1}^n w_i^{(b)} I(y_i \neq f_b(\mathbf{x}_i))}{\sum_{i=1}^n w_i^{(b)}} \quad \text{and} \quad R_b = \frac{1 - \varepsilon_b}{\varepsilon_b}$$

- **aggregation:**

$$F(\mathbf{x}) = \sum_{b=1}^B a_b f_b(\mathbf{x}) \quad \text{where} \quad a_b \propto \log(R_b)$$

- **weights:** sequentially determined adaptively (next slide)

Freund & Schapire (1996)

- initially,

$$w_i^{(1)} = \frac{1}{n}, \quad \text{for all } i = 1, 2, \dots, n$$

- at iteration b , update

$$w_i^{(b+1)} = \begin{cases} w_i^{(b)} \times R_b, & y_i \neq f_b(\mathbf{x}_i); \\ w_i^{(b)}, & y_i = f_b(\mathbf{x}_i) \end{cases}$$

adaptively

called [AdaBoost](#) ... for adaptive boosting
quite a mystery

Dawn of the New Millennium

J. H. Friedman, T. J. Hastie & R. J. Tibshirani (2000), “Additive logistic regression: A statistical view of boosting (with discussion)”, *Ann. Stat.* **28**, pp. 337–407.

L. Breiman (2001), “Random forest”, *Machine Learning* **45**, pp. 5–32.

Friedman, Hastie & Tibshirani (2000)

AdaBoost amounts to constructing the function

$$F(\mathbf{x}) = \sum_{b=1}^B a_b f_b(\mathbf{x})$$

sequentially by choosing a_b and f_b at each step to minimize the so-called **exponential loss**, i.e.,

$$(a_b, f_b) = \arg \min_{\substack{a \in \mathbb{R} \\ f \in \mathcal{F}}} \frac{1}{n} \sum_{i=1}^n e^{-y_i [F_{b-1}(\mathbf{x}_i) + a f(\mathbf{x}_i)]},$$

where

$$F_{b-1}(\mathbf{x}) = \sum_{b'=0}^{b-1} a_{b'} f_{b'}(\mathbf{x}),$$

and \mathcal{F} is a given family of binary classifiers, e.g., **classification trees**.

Exponential Loss

- if y is coded as $+1$ versus -1 (rather than 1 versus 0), standard logistic regression model

$$\log \frac{\mathbb{P}[y = (+1)|\mathbf{x}]}{\mathbb{P}[y = (-1)|\mathbf{x}]} = F(\mathbf{x})$$

means the conditional probability of y given \mathbf{x} can be expressed as

$$\mathbb{P}(y|\mathbf{x}) = \frac{1}{1 + e^{-yF(\mathbf{x})}}, \quad \text{for } y \in \{-1, +1\}$$

- big $\mathbb{P}(y|\mathbf{x}) \Leftrightarrow$ small $e^{-yF(\mathbf{x})}$
- exponential loss, $L[y, F(\mathbf{x})] = e^{-yF(\mathbf{x})}$, “makes sense”

Weights

- at iteration b , objective function can be written as

$$L(a, f) = \sum_{i=1}^n \frac{1}{n} \underbrace{\left[e^{-y_i F_{b-1}(\mathbf{x}_i)} \right]}_{w_i^{(b)}} \left[e^{-ay_i f(\mathbf{x}_i)} \right],$$

where

$$w_i^{(b)} \equiv e^{-y_i F_{b-1}(\mathbf{x}_i)}$$

does not depend on either a or f

- setting $w_i^{(1)} = 1/n$ for all $i \Leftrightarrow$ initializing with $F_0(\mathbf{x}) = 0$

Alternating Minimization

- given f , the function $L(a, f)$ is minimized at

$$a = \frac{1}{2} \log \underbrace{\frac{1 - \varepsilon_b}{\varepsilon_b}}_{R_b}, \quad \varepsilon_b = \frac{\sum_{i=1}^n w_i^{(b)} I(y_i \neq f(\mathbf{x}_i))}{\sum_{i=1}^n w_i^{(b)}},$$

which explains why AdaBoost weights each f_b by a quantity proportional to $\log(R_b)$

- given a , minimizing $L(a, f)$ over $f \in \mathcal{F}$ requires that we fit a classifier f to the training data with each observation (\mathbf{x}_i, y_i) being weighted by $w_i^{(b)}$

Updating Weights

$$F_b(\mathbf{x}) = F_{b-1}(\mathbf{x}) + af(\mathbf{x}) \quad \Rightarrow \quad w_i^{(b+1)} = w_i^{(b)} \times e^{-ay_i f(\mathbf{x}_i)}$$

$$w_i^{(b+1)} = \begin{cases} w_i^{(b)} \times e^a, & \text{if } y_i \neq f(\mathbf{x}_i); \\ w_i^{(b)} \times e^{-a}, & \text{if } y_i = f(\mathbf{x}_i). \end{cases}$$

inflate wrong ones by e^a ; deflate right ones by e^{-a}



inflate wrong ones by $e^{2a} = \frac{1-\varepsilon_b}{\varepsilon_b} = R_b$; leave right ones alone

Details: Minimize $L(\alpha, f)$ Over a

$$f \in \mathcal{F}, \quad f(\mathbf{x}_i) = \pm 1, \quad y_i = \pm 1$$

\Downarrow

either $y_i f(\mathbf{x}_i) = +1$ (right) or $y_i f(\mathbf{x}_i) = -1$ (wrong)

$$L(a, f) = \sum_{y_i f(\mathbf{x}_i) = +1} w_i e^{-a} + \sum_{y_i f(\mathbf{x}_i) = -1} w_i e^a$$

$$\frac{\partial}{\partial a} L(a, f) = 0 \quad \Rightarrow \quad a = \frac{1}{2} \log \frac{1 - \varepsilon}{\varepsilon}, \quad \varepsilon = \frac{\sum_{i=1}^n w_i I(y_i \neq f(\mathbf{x}_i))}{\sum_{i=1}^n w_i}$$

Details: Minimize $L(a, f)$ Over f

$$\begin{aligned}L(a, f) &= \sum_{y_i=f(\mathbf{x}_i)} w_i e^{-a} + \sum_{y_i \neq f(\mathbf{x}_i)} w_i e^a \\&= \sum_{i=1}^n w_i e^{-a} - \sum_{y_i \neq f(\mathbf{x}_i)} w_i e^{-a} + \sum_{y_i \neq f(\mathbf{x}_i)} w_i e^a \\&= \sum_{i=1}^n w_i e^{-a} + (e^a - e^{-a}) \sum_{y_i \neq f(\mathbf{x}_i)} w_i\end{aligned}$$

$$\text{given } a > 0 [e^a > e^{-a}], \quad f = \arg \min_{g \in \mathcal{F}} \sum_{y_i \neq g(\mathbf{x}_i)} w_i$$

⇓

fit classifier to weighted training data, which guarantees $a > 0$

General Framework

1. Pick a **functional class**, \mathcal{F} . (WHERE)
2. Pick a **loss function**, $L[y, F(\mathbf{x})]$. (WHAT)
3. Pick an **optimization method**. (HOW)
4. Start with $F(\mathbf{x}) = 0$, and sequentially build up the function,

$$F(\mathbf{x}) = f_1(\mathbf{x}) + \dots + f_B(\mathbf{x}),$$

by repeating the following steps:

- (a) apply the chosen optimization method from (3) to find

$$f = \arg \min_{g \in \mathcal{F}} L[y, F(\mathbf{x}) + g(\mathbf{x})];$$

- (b) update $F \leftarrow F + f$.

Breiman (2001)

- a **random forest** (RF) refers to

$$\{f_b(\mathbf{x}) \equiv f(\mathbf{x}; \Theta_b) : \Theta_b \stackrel{iid}{\sim} \mathcal{P}_\Theta, b = 1, 2, \dots\}$$

- each $f(\mathbf{x}; \Theta_b)$ is a classifier completely parameterized by Θ_b
- meaning of “ $\Theta_b \stackrel{iid}{\sim} \mathcal{P}_\Theta$ ”:
 - each $f(\cdot; \Theta_b)$ is generated independently with an identical set of stochastic mechanisms, denoted collectively by \mathcal{P}_Θ

“The” RF Algorithm by Breiman

Algorithm

- (i) (*Bootstrap*) First, draw a bootstrap sample, $\mathbb{D}_*^{(b)}$, of the training data \mathbb{D} . Then, use $\mathbb{D}_*^{(b)}$ — not \mathbb{D} — to train a *classification tree*, f_b .
- (ii) (*Random Subset*) When building f_b , randomly select a subset of predictors before making *each* split — say, \mathcal{S} , $|\mathcal{S}| < d$ — and make the best split over the set \mathcal{S} rather than over all possible predictors.

Stochastic Mechanism \mathcal{P}_Θ

- (i) iid sampling from the empirical distribution
- (ii) iid sampling from the set $\{1, 2, \dots, d\}$

BIG Question

- clearly, **bagging** is **random forest** by mechanism (i) alone
- in five years (1996 to 2001), Breiman added mechanism (ii)



?????

WHY

?????

Breiman's Theorem

$$\varepsilon_{RF} \leq \bar{\rho} \left(\frac{1 - s^2}{s^2} \right)$$

- ε_{RF} = prediction error of random forest
- $\bar{\rho}$ = average pairwise correlation between individual classifiers in the forest
- s = average strength of individual classifiers in the forest

will explain for the case of binary classification

Breiman's Theorem

- building block: for an individual classifier, $f(\cdot; \Theta)$, define its margin at (\mathbf{x}, y) as [for binary case]

$$m(\Theta; \mathbf{x}, y) = I(f(\mathbf{x}; \Theta) = y) - I(f(\mathbf{x}; \Theta) \neq y),$$

so that

$$m(\Theta; \mathbf{x}, y) > 0, \quad \text{if } f(\cdot; \Theta) \text{ right about } (\mathbf{x}, y);$$

$$m(\Theta; \mathbf{x}, y) < 0, \quad \text{if } f(\cdot; \Theta) \text{ wrong about } (\mathbf{x}, y)$$

- ε_{RF} , $\bar{\rho}$ and s all defined in terms of probability- and/or moment-transformations of $m(\Theta; \mathbf{x}, y)$, w.r.t. $\mathcal{P}_{(\mathbf{x}, y)}$ or \mathcal{P}_{Θ}
- relationship can then be established by well-known inequalities that link $\mathbb{P}(Z > c)$ with $\mathbb{E}(Z)$, $\text{Var}(Z)$, etc

Breiman's Theorem: ϵ_{RF}

- first, generalize the notion of **margin** to the entire forest
- margin of forest at (\mathbf{x}, y) is simply

$$M(\mathbf{x}, y) = \mathbb{E}_{\Theta}(m(\Theta; \mathbf{x}, y))$$

or the average margin of all individual classifiers in the forest

- the **prediction error** of the forest is simply the probability of the forest making incorrect classifications, i.e.,

$$\epsilon_{RF} = \mathbb{P}_{(\mathbf{x}, y)}(M(\mathbf{x}, y) < 0)$$

Breiman's Theorem: s

- recall **positive** margin $m(\Theta; \mathbf{x}, y) > 0$ is **good** for classifier $f(\cdot; \Theta)$ and **negative** margin $m(\Theta; \mathbf{x}, y) < 0$ is **bad**
- natural to measure **average strength** of a typical classifier in the forest by the overall average of individual margins, i.e.,

$$s = \mathbb{E}_{\Theta} \left[\mathbb{E}_{(\mathbf{x}, y)} (m(\Theta; \mathbf{x}, y)) \right]$$

Breiman's Theorem: $\bar{\rho}$

- first, a correlation coefficient measuring if $f(\cdot; \Theta)$ and $f(\cdot; \Theta')$ are more or less likely to be correct together:

$$\rho(\Theta, \Theta') = \frac{\text{Cov}_{(\mathbf{x}, y)}[m(\Theta; \mathbf{x}, y), m(\Theta'; \mathbf{x}, y)]}{\text{Sd}_{(\mathbf{x}, y)}[m(\Theta; \mathbf{x}, y)] \times \text{Sd}_{(\mathbf{x}, y)}[m(\Theta'; \mathbf{x}, y)]}$$

- let

$$w(\Theta, \Theta') \equiv \text{Sd}_{(\mathbf{x}, y)}[m(\Theta; \mathbf{x}, y)] \times \text{Sd}_{(\mathbf{x}, y)}[m(\Theta'; \mathbf{x}, y)],$$

the **average pairwise correlation** is defined as

$$\bar{\rho} = \frac{\mathbb{E}_{\Theta, \Theta'} [\rho(\Theta, \Theta') \times w(\Theta, \Theta')]}{\mathbb{E}_{\Theta, \Theta'} [w(\Theta, \Theta')]},$$

for $\Theta, \Theta' \stackrel{iid}{\sim} \mathcal{P}_{\Theta}$... [in fact, a **weighted** average]

Significance

Example Say $s = 0.4$, $\bar{\rho} = 0.5$, then theory says

$$\varepsilon_{RF} \leq \bar{\rho} \left(\frac{1 - s^2}{s^2} \right) = (0.5) \times \left[\frac{1 - (0.4)^2}{(0.4)^2} \right] = 2.625.$$

But of course $\varepsilon_{RF} \leq 1$. A useless bound?

Implication A good random forest should consist of individual classifiers that are

- as strong as possible (large s), and
- as uncorrelated as possible (small $\bar{\rho}$).

Adding mechanism (ii) further decorrelates the classifiers and diversifies the forest.

Strength-Diversity Trade-off

Unfortunately, there is a trade-off between the two, s and $\bar{\rho}$.

Discussion

- (a) What's a tuning parameter in “the” random forest algorithm?
- (b) What happens to s and to $\bar{\rho}$, respectively, as this parameter increases?

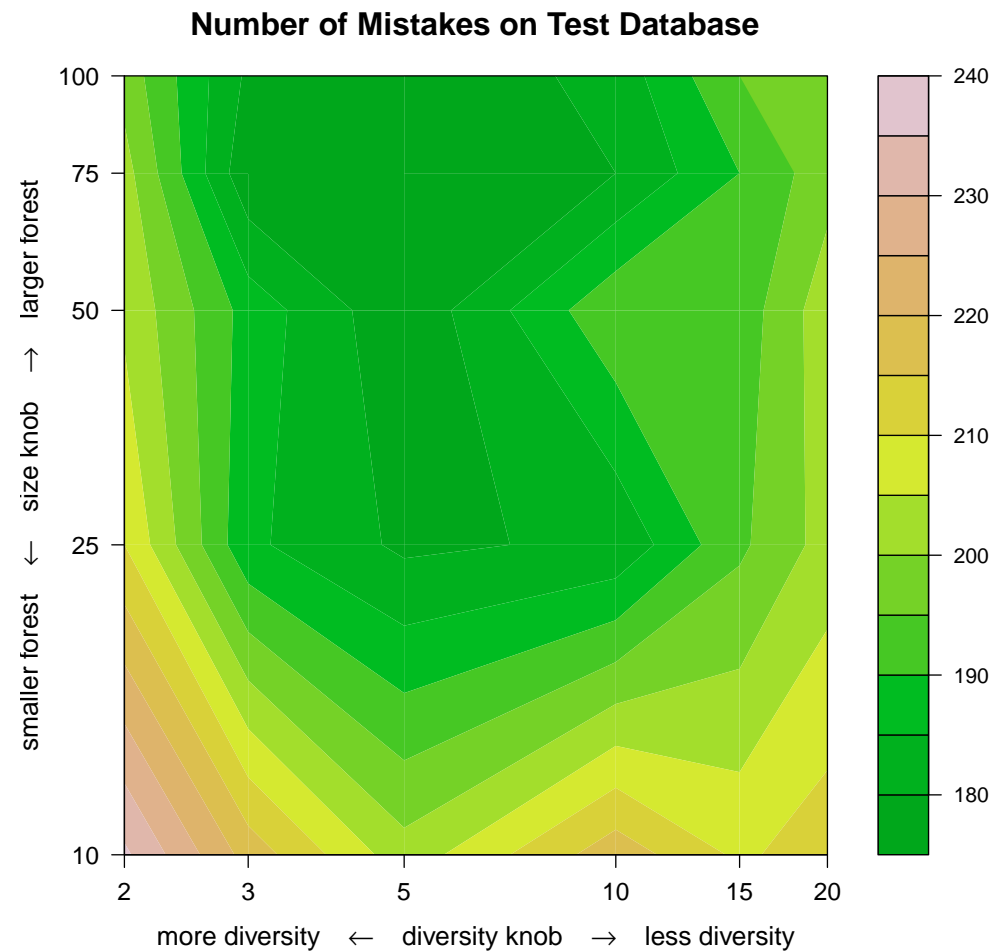


Illustration using the “spam” data set from <http://statweb.stanford.edu/~tibs/ElemStatLearn/datasets/spam.data>.

Decoupling in AdaBoost

$$\varepsilon_b = \frac{\sum_{i=1}^n w_i I(y_i \neq f_b(\mathbf{x}_i))}{\sum_{i=1}^n w_i} \Rightarrow \varepsilon_b \sum_{\text{all}} w_i = \sum_{\text{wrong}} w_i$$
$$(1 - \varepsilon_b) \sum_{\text{all}} w_i = \sum_{\text{right}} w_i$$

$$w_i^{\text{new}} = \begin{cases} w_i \times \frac{1 - \varepsilon_b}{\varepsilon_b}, & i \in \text{wrong}; \\ w_i, & i \in \text{right} \end{cases}$$

$$\sum_{\text{wrong}} w_i^{\text{new}} = \frac{1 - \varepsilon_b}{\varepsilon_b} \sum_{\text{wrong}} w_i = \sum_{\text{right}} w_i = \sum_{\text{right}} w_i^{\text{new}}$$

Decoupling in AdaBoost

- error of f_b under w_i^{new} = exactly 50%, worst possible
- f_{b+1} trained on a (weighted) data set f_b cannot classify
- AdaBoost actively de-correlates f_{b+1} from f_b

Summary

- Part 1
 - bagging (1996)
 - AdaBoost (1996)
 - statistical view of boosting (2000)
 - random forest (2001)
- Part 2 (next, after short break)
 - gradient boosting machine (2001)
 - Breiman's theorem (2001)
 - some new directions (2005+)