

Lecture 5

Ensembles – Part 2

(February 6, 2015)

Mu Zhu
University of Waterloo

Outline for Part 2

- gradient boosting machine (2001)
- Breiman's theorem (2001)
- some new directions (2005+)
 - variable selection ensembles (2006, 2010, 2011, 2012)
 - machine learning markets (2011)
 - magging (2014)

General Framework for Boosting

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 $F \leftarrow F + f$.

Gradient Boosting Machine

- Friedman (2001; *Ann. Stat.*)
- free to choose $L(y, F)$
- free to choose \mathcal{F} in principle, but mostly **trees** in practice
- choice of optimization method = take one gradient step, i.e.,

$$f = f_0 - \eta \left[\frac{\partial L}{\partial f} \right]_{f=f_0},$$

where $\eta > 0$ is step size and f_0 is initial “estimate” of f

Gradient Boosting Machine

- given current F , initial f_0 for the “next” function = 0
- one gradient step requires setting $f \propto -(\partial L/\partial f)|_{f=0}$
- but must do so within the constraint, $f \in \mathcal{F}$, so

let $r_i \equiv h(\mathbf{x}_i, y_i)$, where $h \equiv -(\partial L/\partial f)|_{f=0}$

$$\text{train } f = \arg \min_{g \in \mathcal{F}} \sum_{i=1}^n [r_i - g(\mathbf{x}_i)]^2$$

$$\text{set } \eta = \arg \min_{\xi} \sum_{i=1}^n L[y_i, F(\mathbf{x}_i) + \xi f(\mathbf{x}_i)]$$

- negative **functional gradient**, $h = -(\partial L/\partial f)|_{f=0}$, next slide

Ex: GBM with L_2 Loss Function

- consider **regression** instead of **classification**, i.e., $y_i \in \mathbb{R}$

- natural to choose

$$L[y, F(\mathbf{x}) + f(\mathbf{x})] = \frac{1}{2} [y - F(\mathbf{x}) - f(\mathbf{x})]^2$$

- then,

$$h = - \left[\frac{\partial L}{\partial f} \right]_{f=0} = (y - F - f) |_{f=0} = y - F$$

- thus, negative gradient = residual of the current model F , so

$$r_i = h(\mathbf{x}_i, y_i) = y_i - F(\mathbf{x}_i)$$

- so, repeatedly fit a model $\in \mathcal{F}$ to residuals of current model

Breiman's Theorem

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

- technically, requires the assumption that $s > 0$
- that is, an “average” classifier in the forest will correctly classify an “average” observation — not a very strong assumption
- next, elucidate the proof in three steps

Step 1

- target is (details later; requires $s > 0$)

$$\varepsilon_{RF} \equiv P_{(\mathbf{x}, y)}(M(\mathbf{x}, y) < 0) \leq \frac{\text{Var}_{(\mathbf{x}, y)}(M(\mathbf{x}, y))}{\mathbb{E}_{(\mathbf{x}, y)}^2(M(\mathbf{x}, y))}$$

- recall definition

$$\begin{aligned} s &= \mathbb{E}_{\Theta} [\mathbb{E}_{(\mathbf{x}, y)}(m(\Theta; \mathbf{x}, y))] = \mathbb{E}_{(\mathbf{x}, y)} [\mathbb{E}_{\Theta}(m(\Theta; \mathbf{x}, y))] \\ &= \mathbb{E}_{(\mathbf{x}, y)}(M(\mathbf{x}, y)), \end{aligned}$$

so **denominator** already s^2

- thus, after this step, suffices to prove **numerator** $\leq \bar{\rho}(1 - s^2)$

Step 2

- target is (again, details later)

$$\text{Var}_{(\mathbf{x}, y)}(M(\mathbf{x}, y)) = \mathbb{E}_{\Theta, \Theta'} (\text{Cov}_{(\mathbf{x}, y)}[m(\Theta; \mathbf{x}, y), m(\Theta'; \mathbf{x}, y)]) ,$$

for $\Theta, \Theta' \stackrel{iid}{\sim} P_{\Theta}$

- definition of $\bar{\rho}$ + a (trivial) identity for $\Theta, \Theta' \stackrel{iid}{\sim} P_{\Theta}$,

$$\mathbb{E}_{\Theta, \Theta'}(g(\Theta)g(\Theta')) = \mathbb{E}_{\Theta}(g(\Theta))\mathbb{E}_{\Theta'}(g(\Theta')) = \mathbb{E}_{\Theta}^2(g(\Theta)),$$

+ Jensen's inequality \Rightarrow

$$\begin{aligned} \text{Var}_{(\mathbf{x}, y)}(M(\mathbf{x}, y)) &= \bar{\rho} \mathbb{E}_{\Theta}^2 (\text{Sd}_{(\mathbf{x}, y)} m(\Theta; \mathbf{x}, y)) \\ &\leq \bar{\rho} \mathbb{E}_{\Theta} (\text{Sd}_{(\mathbf{x}, y)}^2 m(\Theta; \mathbf{x}, y)) = \bar{\rho} \mathbb{E}_{\Theta} (\text{Var}_{(\mathbf{x}, y)} m(\Theta; \mathbf{x}, y)) \end{aligned}$$

Step 3

- after step 2, remains to prove (i.e., target for step 3)

$$\mathbb{E}_{\Theta} (\text{Var}_{(\mathbf{x},y)} (m(\Theta; \mathbf{x}, y))) \leq (1 - s^2)$$

- simply write $m = m(\Theta; \mathbf{x}, y)$, and get

$$\begin{aligned} \mathbb{E}_{\Theta} [\text{Var}_{(\mathbf{x},y)} (m)] &= \mathbb{E}_{\Theta} \left[\mathbb{E}_{(\mathbf{x},y)} (m^2) - \mathbb{E}_{(\mathbf{x},y)}^2 (m) \right] \\ &= \mathbb{E}_{\Theta} \left[\mathbb{E}_{(\mathbf{x},y)} (m^2) \right] - \mathbb{E}_{\Theta} \left[\mathbb{E}_{(\mathbf{x},y)}^2 (m) \right] \equiv \text{(I)} - \text{(II)} \end{aligned}$$

but

$$\text{(I)} = 1, \quad \text{since } m = \pm 1 \Rightarrow m^2 \equiv 1;$$

$$\text{(II)} = \mathbb{E}_{\Theta} \left(\mathbb{E}_{(\mathbf{x},y)}^2 (m) \right) \geq \mathbb{E}_{\Theta}^2 \left(\mathbb{E}_{(\mathbf{x},y)} (m) \right) = s^2,$$

again, by [Jensen's inequality](#)

Step 1: Details

- (Lemma) If X is a random variable with $\mu \equiv \mathbb{E}(X) > 0$, then

$$P(X < 0) \leq \frac{\text{Var}(X)}{\mu^2}.$$

- Lemma above can be proved by simply taking $\epsilon = \mu$ in [Chebyshev's inequality](#),

$$P(|X - \mu| \geq \epsilon) \leq \frac{\text{Var}(X)}{\epsilon^2}, \quad \text{for any } \epsilon > 0.$$

- Step 1 can be completed by simply applying this Lemma to $X = M(\mathbf{x}, y)$, which requires $s = \mathbb{E}_{(\mathbf{x}, y)}(M(\mathbf{x}, y)) > 0$.

Step 2: Details

- shorthand: $m = m(\Theta; \mathbf{x}, y)$; $m' = m(\Theta'; \mathbf{x}, y)$
- if $\Theta, \Theta' \stackrel{iid}{\sim} \mathcal{P}_\Theta$, $\mathbb{E}_\Theta(m) = \mathbb{E}_{\Theta'}(m')$ [trivial]
- from the definition of M , get

$$\begin{aligned} & \text{Var}_{(\mathbf{x}, y)}(M(\mathbf{x}, y)) \\ &= \mathbb{E}_{(\mathbf{x}, y)}(M^2(\mathbf{x}, y)) - \mathbb{E}_{(\mathbf{x}, y)}^2(M(\mathbf{x}, y)) \\ &= \mathbb{E}_{(\mathbf{x}, y)}(\mathbb{E}_\Theta^2[m]) - \mathbb{E}_{(\mathbf{x}, y)}^2(\mathbb{E}_\Theta[m]) \\ &= \mathbb{E}_{(\mathbf{x}, y)}(\mathbb{E}_\Theta[m]\mathbb{E}_{\Theta'}[m']) - \mathbb{E}_{(\mathbf{x}, y)}(\mathbb{E}_\Theta[m])\mathbb{E}_{(\mathbf{x}, y)}(\mathbb{E}_{\Theta'}[m']) \\ &= \mathbb{E}_{\Theta, \Theta'}(\mathbb{E}_{(\mathbf{x}, y)}[mm']) - \mathbb{E}_{\Theta, \Theta'}(\mathbb{E}_{(\mathbf{x}, y)}[m]\mathbb{E}_{(\mathbf{x}, y)}[m']) \\ &= \mathbb{E}_{\Theta, \Theta'}(\text{Cov}_{(\mathbf{x}, y)}[m, m']) \end{aligned}$$

Diversity and Ensembles

- importance of [diversity](#)
- power of [ensembles](#)
- testimony
 - high-profile, million-dollar Netflix contest ([2006-09](#))
 - prediction contests on www.kaggle.com

Some New Directions

1. variable selection ensembles

- use ensembles for a problem other than prediction

2. machine learning markets

- how to aggregate other than taking averages

3. magging

- how to aggregate, but for a different problem

disclaimer

selection of directions purely dictated by *personal taste*

Variable Selection Ensembles

Idea

- standard variable selection algorithms all **deterministic**
- do multiple runs
 - use a **stochastic** mechanism to “force” standard algorithms to give **different** answers for each run
- aggregate the results

Challenges

- delicate balance between **false positives** and **false negatives**
- characterizing/finding the “right” stochastic mechanism

Stochastic Mechanisms

	Stochastic Mechanism	Type
Z&C	Darwinian evolution in parallel universes	SO
M&B	Bootstrap + random scaling	Bt/Sub
W+	Bootstrap + random subset	Bt/Sub
X&Z	Stochastic stepwise search	SO

SO = stochastic optimization; Bt/Sub = bootstrap and/or sub-sampling

References

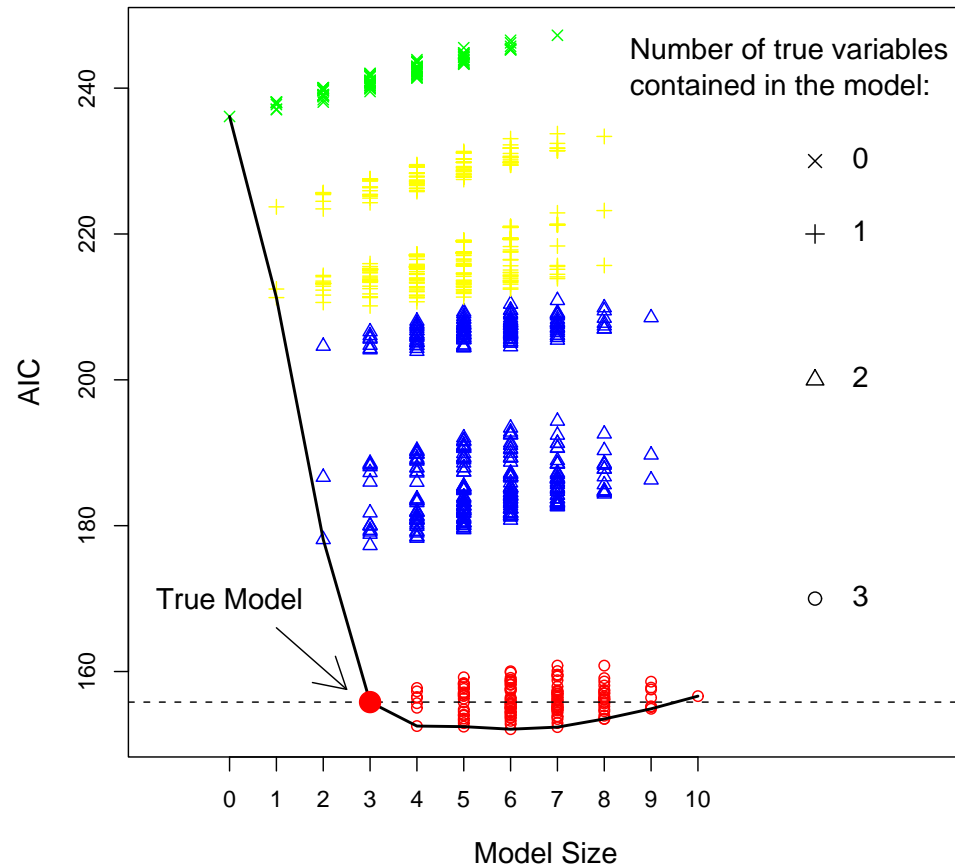
Z&C: Zhu & Chipman (2006; *Technometrics*)

M&B: Meinshausen & Bühlmann (2010; *JRSSB*)

W+: Wang *et al.* (2011; *Ann. Appl. Stat.*)

X&Z: Xin & Zhu (2012; *J. Comp. Graph. Stat.*)

AIC of All 2^{10} Models



Total of ten predictors x_1, x_2, \dots, x_{10} . True model:
 $y = x_2 + x_5 + x_8 + \varepsilon$. [M. Zhu (2008; *Am Stat*)]

Machine Learning Markets

- classification **outcomes**: c_1, c_2, \dots, c_K
- **price** vector (p_k is the price for betting on the outcome c_k):

$$\mathbf{p} = (p_1, p_2, \dots, p_K)^T$$

- **reward** vectors:

$$\mathbf{r}_1 = (1, 0, 0, \dots, 0, 0)^T$$

$$\mathbf{r}_2 = (0, 1, 0, \dots, 0, 0)^T$$

$$\vdots \qquad \qquad \qquad \vdots$$

$$\mathbf{r}_K = (0, 0, 0, \dots, 0, 1)^T$$

- **players** $b = 1, 2, \dots, B$

Each Player

- **wealth**: w_b
- **utility function**: $u_b(\cdot)$
- **belief** for outcome c_k : $f_b(k)$
- **holding** vector (**fractional** and **negative** positions allowed):

$$\mathbf{z}_b = [z_b(1), z_b(2), \dots, z_b(K)]^T$$

- each player maximizes his/her **expected utility**,

$$\max_{\mathbf{z}_b} \sum_{k=1}^K u_b (w_b - \mathbf{p}^T \mathbf{z}_b + \mathbf{r}_k^T \mathbf{z}_b) f_b(k),$$

under various market constraints

Equilibrium Prices

A. Storkey (2011), “Machine learning markets”, in *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics*, pp. 716-724.

- computed various equilibrium prices for different utility functions
- e.g., if $u_b(\cdot) = \log(\cdot)$, then

$$p_k^* = \frac{\sum_b w_b f_b(k)}{\sum_b w_b}$$

- implication: with different u_1, u_2, \dots, u_B , p_k^* would represent different (perhaps more complex) combinations of $f_b(k)$

Magging

- Meinshausen & Bühlmann (2014; arXiv:1406.0596) +
Bühlmann & Meinshausen (2014; arXiv:1409.2638)

- suppose \exists heterogeneous groups, G_1, G_2, \dots, G_B , and

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta}_i + \varepsilon_i, \quad \text{where } \boldsymbol{\beta}_i = \boldsymbol{\beta}_b \quad \forall \quad i \in G_b$$

- don't know which group \mathbf{x}_{new} belongs to @ prediction time
- use a regression estimate $\boldsymbol{\beta}^*$ that maximizes the “explained variance in the most adversarial scenario” instead of, say, averaging individual estimates from each group

Explained Variance

- suppose true coefficient = β and estimate = γ
- then, explained variance is

$$\begin{aligned}\Omega(\beta, \gamma) &= \text{Var}(y) - \text{Var}(y - \mathbf{x}^T \gamma) \\ &= \text{Var}(y) - [\text{Var}(y) - 2\text{Cov}(y, \mathbf{x}^T \gamma) + \text{Var}(\mathbf{x}^T \gamma)] \\ &= 2\mathbf{c}^T \gamma - \gamma^T \Sigma \gamma\end{aligned}$$

where

$$\mathbf{c} \equiv \text{Cov}(y, \mathbf{x}) = \text{Cov}(\mathbf{x}^T \beta + \varepsilon, \mathbf{x}) \quad \text{and} \quad \Sigma \equiv \text{Var}(\mathbf{x})$$

- under (common) assumption of $\text{Cov}(\varepsilon, \mathbf{x}) = 0$, get

$$\mathbf{c} = \Sigma \beta \quad \Rightarrow \quad \Omega(\beta, \gamma) = 2\beta^T \Sigma \gamma - \gamma^T \Sigma \gamma$$

The Maximin Effect

- assume $\beta_i \sim \mathcal{P}_\beta$ [e.g., for heterogeneous groups, $\mathcal{P}_\beta =$ finite mixture of point masses]
- go after the quantity,

$$\beta^* = \max_{\gamma} \min_{\beta \in \text{Supp}(\mathcal{P}_\beta)} \Omega(\beta, \gamma),$$

which they called the “**maximin effect**”

- the term “**most adversarial scenario**” refers to the case in which our estimate is least capable of explaining the variance of y

Estimating the Maximin Effect

- first, **generate** different individual estimates $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_B$ using random sub-samples of the training data
- then, **aggregate** these individual estimates by taking the point in their **convex hull** $\mathcal{H}(\hat{\beta}_1, \dots, \hat{\beta}_B)$ that is closest to the origin (in a particular norm), i.e.,

$$\hat{\beta}^* = \arg \min_{\mathbf{u} \in \mathcal{H}(\hat{\beta}_1, \dots, \hat{\beta}_B)} \|\mathbf{u}\|_{\hat{\Sigma}}^2$$

where $\|\mathbf{u}\|_{\hat{\Sigma}}^2 = \mathbf{u}^T \hat{\Sigma} \mathbf{u}$

called **Magg**ing ... for **maximin aggregating**

Our Journey: Bagging → Magging

- many successes, but successes also perpetuate myths,
 - e.g., “if a collection of classifiers are ‘nearly independent’ of each other, then taking majority votes over their results will always achieve greater accuracy”
- for a counter example, see
 - S. B. Vardeman & M. D. Morris (2013), “Majority voting by independent classifiers can increase error rates”, *The American Statistician* **67**, pp. 94–96.
- for today’s lecture + more on such myths, watch out for
 - M. Zhu (forthcoming), “Use of majority votes in statistical learning”, *Wiley Interdisciplinary Reviews: Comp. Stat.*

Summary

- key ideas:
 - ensembles
 - strength-diversity trade-off
- specific methods:
 - bagging; AdaBoost
 - random forest; gradient boosting machine
 - variable selection ensemble
 - machine learning market
 - magging