

Boosting

(1)

Consider observations (x_i, y_i) , $i=1, \dots, N$, $y_i \in \{-1, 1\}$. The purpose of boosting is to generate a sequence of simple classifiers $G_m(x) \in \{-1, 1\}$, $m=1, \dots, M$, and then combine them through a weighted majority vote to produce the final prediction $G(x) = \text{sign}(\sum_{m=1}^M \alpha_m G_m(x))$.

Algorithm: AdaBoost

1. Initialize $w_i = \frac{1}{N}$, $i=1, \dots, N$
2. For $m=1$ to M
 - (a) Fit a classifier $G_m(x)$ to the training data using weights w_i
 - (b) Compute
$$\text{err}_m = \frac{\sum_{i=1}^N w_i \mathbb{1}_{\{y_i \neq G_m(x_i)\}}}{\sum_{i=1}^N w_i}$$
 - (c) Compute
$$\alpha_m = \log \frac{1 - \text{err}_m}{\text{err}_m}$$
 - (d) Set $w_i \leftarrow w_i e^{\alpha_m \mathbb{1}_{\{y_i \neq G_m(x_i)\}}}$
3. Output $G(x) = \text{sign}[\sum_{m=1}^M \alpha_m G_m(x)]$

Note that we should choose G_m such that $1 - \text{err}_m > \text{err}_m$, i.e.
$$\sum_{i=1}^N w_i \mathbb{1}_{\{y_i = G_m(x_i)\}} > \sum_{i=1}^N w_i \mathbb{1}_{\{y_i \neq G_m(x_i)\}}$$
. In particular, if $1 - \text{err}_m < \text{err}_m$, we should choose $\tilde{G}_m(x) = -G_m(x)$ to get $1 - \tilde{\text{err}}_m > \tilde{\text{err}}_m$.

Then we have $\alpha_m > 0$, and thus observations that are misclassified by $G_m(x)$ have their weights scaled by e^{α_m} , increasing their relative influence for inducing the next classifier $G_{m+1}(x)$ in the sequence.

The reason why AdaBoost works for well-chosen $G_m(x)$ is that it can be considered as a Forward Stagewise Additive Modeling algorithm.

Suppose prediction is made based on a function

$$f(x) = \sum_{m=1}^M \beta_m b(x, \theta_m)$$

(2)

where $\beta_m, m=1, \dots, M$, are coefficients, and $b(x, \gamma) \in \mathbb{R}$ are function of x , characterized by a set of parameters γ . Given a Loss function $L(y, f(x))$, the Forward Stagewise Additive Modeling is

1. Initialize $f_0(x) = 0$

2. For $m=1$ to M

(a) Compute $(\beta_m, \gamma_m) = \arg \min_{\beta, \gamma} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + \beta b(x_i, \gamma))$

(b) Set $f_m(x) = f_{m-1}(x) + \beta_m b(x, \gamma_m)$

This algorithm works because $\sum_{i=1}^N L(y_i, f_m(x_i))$ is decreasing. In this sense, we may not need to compute the minimizers in 2(a).

Consider the loss function $L(y, f(x)) = \exp(-yf(x))$ which is small when $\text{sign}(y) = \text{sign}(f(x))$. Consider the minimization problem

$$(\beta_m, G_m) = \arg \min_{\beta, G} \sum_{i=1}^N \exp[-y_i(f_{m-1}(x_i) + \beta G(x_i))] \quad G(x) \in \{-1, 1\}$$

$$= \arg \min_{\beta, G} \sum_{i=1}^N w_i^{(m)} e^{-y_i \beta G(x_i)} \quad w_i^{(m)} = e^{-y_i f_{m-1}(x_i)}$$

$$= \arg \min_{\beta, G} \left(\sum_{y_i = G(x_i)} w_i^{(m)} e^{-\beta} + \sum_{y_i \neq G(x_i)} w_i^{(m)} e^{\beta} \right)$$

$$= \arg \min_{\beta, G} \left((e^{\beta} - e^{-\beta}) \sum_{i=1}^N w_i^{(m)} \mathbb{1}_{\{y_i \neq G(x_i)\}} + e^{-\beta} \sum_{i=1}^N w_i^{(m)} \right)$$

If $\beta > 0$, then $G_m = \arg \min_G \sum_{i=1}^N w_i^{(m)} \mathbb{1}_{\{y_i \neq G(x_i)\}}$. Again, we don't need to choose G_m to be the minimizer of $\sum_{i=1}^N w_i^{(m)} \mathbb{1}_{\{y_i \neq G(x_i)\}}$. What we need is G_m is good enough such that the loss function is decreasing.

In particular, we may choose G_m such that $1 - \text{err}_m > \text{err}_m$, with

$$\text{err}_m = \frac{\sum_{i=1}^N w_i^{(m)} \mathbb{1}_{\{y_i \neq G_m(x_i)\}}}{\sum_{i=1}^N w_i^{(m)}}$$

Now, with G_m is chosen, set the derivative of the object function with respect to β to 0, we get $(e^{\beta} + e^{-\beta}) \sum_{i=1}^N w_i^{(m)} \mathbb{1}_{\{y_i \neq G_m(x_i)\}} - e^{-\beta} \sum_{i=1}^N w_i^{(m)} = 0$

$$\Rightarrow \beta_m = \frac{1}{2} \log \frac{1 - \text{err}_m}{\text{err}_m} = \frac{1}{2} \Delta_m$$

(3)

The approximation is then updated $f_m(x) = f_{m-1}(x) + \beta_m G_m(x)$

$$\Rightarrow w_i^{(m+1)} = e^{-y_i f_m(x_i)} = e^{-y_i f_{m-1}(x_i) - y_i \beta_m G_m(x_i)} = w_i^{(m)} e^{-\beta_m y_i G_m(x_i)}$$

Using the fact that $-y_i G_m(x_i) = 2 \mathbb{1}\{y_i \neq G_m(x_i)\} - 1$ and $\alpha_m = 2\beta_m$

$$\therefore w_i^{(m+1)} = w_i^{(m)} e^{\alpha_m \mathbb{1}\{y_i \neq G_m(x_i)\}} e^{-\beta_m}$$

The factor $e^{-\beta_m}$ multiplies all weights by the same value, so it has no effect.

The idea of boosting (Forward Stagewise Additive Modeling) can be applied for other problems. For regression, we consider $L(y, f(x)) = (y - f(x))^2$.

For K-class classification, we need to first define the prediction function. Let $p_k(x) = \Pr(Y = f_{jk} | x)$ where $Y \in \{f_{j1}, f_{j2}, \dots, f_{jK}\}$

Recall in logistic regression, $Y \in \{0, 1\}$ and we let

$$p_1(x) = \Pr(Y=1|x) = \frac{e^{f_1(x)}}{1 + e^{f_1(x)}} \quad \text{where } f_1(x) = x^T \beta \text{ in logistic regression}$$

$$\Leftrightarrow f_1(x) = \log \frac{p_1(x)}{1 - p_1(x)} = \log \frac{\Pr(Y=1|x)}{\Pr(Y=0|x)}$$

Now suppose $Y \in \{0, 1, 2, \dots, K-1\}$, we extend the idea to define

$$f_2(x) = \log \frac{\Pr(Y=2|x)}{\Pr(Y=0|x)}, \dots, f_{K-1}(x) = \log \frac{\Pr(Y=K-1|x)}{\Pr(Y=0|x)}$$

$$\Rightarrow \Pr(Y=2|x) = e^{f_2(x)} \Pr(Y=0|x), \dots, \Pr(Y=K-1|x) = e^{f_{K-1}(x)} \Pr(Y=0|x)$$

$$\text{With } \Pr(Y=0|x) + \Pr(Y=1|x) + \dots + \Pr(Y=K-1|x) = 1$$

$$\Rightarrow \Pr(Y=0|x) (1 + e^{f_1(x)} + e^{f_2(x)} + \dots + e^{f_{K-1}(x)}) = 1$$

$$\Rightarrow \Pr(Y=0|x) = \frac{1}{1 + \sum_{k=1}^{K-1} e^{f_k(x)}} \quad \Pr(Y=k|x) = \frac{e^{f_k(x)}}{1 + \sum_{k=1}^{K-1} e^{f_k(x)}}$$

$$\text{Let } f_0(x) = 0, \text{ then we have } \Pr(Y=k|x) = \frac{e^{f_k(x)}}{\sum_{k=0}^{K-1} e^{f_k(x)}} \quad \text{for } k=0, 1, \dots, K-1$$

Note that adding an arbitrary $h(x)$ to each $f_k(x)$ leaves the model unchanged, instead of setting $f_0(x) = 0$, we can impose the constraint $\sum_{k=0}^{K-1} f_k(x) = 0$ to retain the symmetry.

For $Y \in \{f_1, \dots, f_K\}$, $p_k(x) = \Pr(Y = f_k | x)$, we define

$$p_k(x) = \frac{e^{f_k(x)}}{\sum_{\ell=1}^K e^{f_\ell(x)}} \quad \text{with} \quad \sum_{\ell=1}^K f_\ell(x) = 0$$

(4)

Consider the K -class multinomial deviance loss function

$$\begin{aligned} L(y, \vec{f}(x)) &= - \sum_{k=1}^K \mathbb{1}_{\{y = f_k\}} \log p_k(x) \\ &= - \sum_{k=1}^K \mathbb{1}_{\{y = f_k\}} f_k(x) + \log \left(\sum_{\ell=1}^K e^{f_\ell(x)} \right) \end{aligned} \quad \vec{f}(x) = (f_1(x), \dots, f_K(x))^T$$

Recall for logistic regression, the likelihood function is $p^y (1-p)^{1-y}$ and we choose $p = \Pr(Y=1)$ to minimize the negative likelihood

$$-\ell(p) = -[y \log p + (1-y) \log(1-p)] = -[\mathbb{1}_{\{y=1\}} \log p_1 + \mathbb{1}_{\{y=0\}} \log p_0]$$

Now, we consider the update $\vec{f}_{m+1}(x) = \vec{f}_m(x) + \vec{\delta}$ for $x \in R$, $\vec{\delta} = (\delta_1, \dots, \delta_K)^T$

Since we only consider $x \in R$, given (x_i, y_i) , $i=1, \dots, N$, the loss function is

$$- \sum_{x_i \in R} \sum_{k=1}^K \mathbb{1}_{\{y_i = f_k\}} f_k(x_i) + \sum_{x_i \in R} \log \left(\sum_{\ell=1}^K e^{f_\ell(x_i)} \right)$$

Let $y_{ik} = \mathbb{1}_{\{y_i = f_k\}}$, then we want to solve

$$\min_{\vec{\delta}} - \sum_{x_i \in R} \sum_{k=1}^K y_{ik} (f_k(x_i) + \delta_k) + \sum_{x_i \in R} \log \left(\sum_{\ell=1}^K e^{f_\ell(x_i) + \delta_\ell} \right) = \min_{\vec{\delta}} L(\vec{\delta})$$

It is costly to find the minimizer. Instead we use Quasi-Newton method to find $\vec{\delta}^*$ such that $L(\vec{\delta}^*) < L(0)$. $\vec{\delta}^*$ will be of the form

$$\vec{\delta}^* = 0 - D^{-1} \nabla L(0) = -D \nabla L(0), \text{ where } D \text{ is a positive definite matrix}$$

By Taylor expansion, $L(\vec{\delta}^*) = L(0) + \nabla L(0)^T \vec{\delta}^* + O(\|\vec{\delta}^*\|^2)$

$$= L(0) - \nabla L(0)^T D^{-1} \nabla L(0) + O(\|\vec{\delta}^*\|^2) < L(0) \quad \text{if } \|\vec{\delta}^*\| \text{ small}$$

We choose D to be the diagonal entries of $(\nabla^2 L(0))$ for Quasi-Newton method. To get rid of the constraint $\sum_{k=1}^K \delta_k = 0$, we first assume $\delta_K = 0$ and assume no constraints on $\delta_1, \dots, \delta_{K-1}$. Consider

$$\frac{\partial}{\partial \delta_k} L(\vec{\delta}) = - \sum_{x_i \in R} y_{ik} + \sum_{x_i \in R} \frac{e^{f_k(x_i) + \delta_k}}{\sum_{\ell=1}^K e^{f_\ell(x_i) + \delta_\ell}} \quad 1 \leq k \leq K-1$$

$$\frac{\partial^2 L(\delta)}{\partial \delta_k^2} = \sum_{i \in R} \left(- \frac{e^{f_k(x_i) + \delta_k}}{\left(\sum_{l=1}^K e^{f_l(x_i) + \delta_l} \right)^2} + \frac{e^{f_k(x_i) + \delta_k}}{\sum_{l=1}^K e^{f_l(x_i) + \delta_l}} \right) \quad (5)$$

$$\therefore \delta^* = -D\delta L(0) \Rightarrow \delta_k^* = - \left(\frac{\partial^2 L(0)}{\partial \delta_k^2} \right)^{-1} \frac{\partial}{\partial \delta_k} L(0)$$

$$p_{ik} = \frac{e^{f_k(x_i)}}{\sum_{l=1}^K e^{f_l(x_i)}}$$

$$\begin{aligned} &= - \left(\sum_{i \in R} (-p_{ik}^2 + p_{ik}) \right)^{-1} \left(- \sum_{i \in R} (y_{ik} - p_{ik}) \right) \\ &= \frac{\sum_{i \in R} (y_{ik} - p_{ik})}{\sum_{i \in R} p_{ik} (1 - p_{ik})} \quad \text{for } 1 \leq k \leq K-1 \end{aligned}$$

Note that δ_k^* does not depend on K (recall $\delta_K = 0$). Therefore, if we assume $\delta_1 = 0$, instead of $\delta_K = 0$, we will get the same δ_k^* for $k = 2, \dots, K-1$.

Instead of setting $\delta_K = 0$ (and thus $f_K(x) = 0$ by the updates $\vec{f}_{m+1}(x) = \vec{f}_m(x) + \vec{\delta}$ and $\vec{f}_0(x) = 0$), one would prefer $\sum_{l=1}^K f_l(x) = 0$. It means that $\sum_{l=1}^K \delta_l = 0$ is required. From $p_k(x) = \frac{e^{f_k(x)}}{\sum_{l=1}^K e^{f_l(x)}} = \frac{e^{f_k(x) - a}}{\sum_{l=1}^K e^{f_l(x) - a}}$ for any a independent of index l ,

we can consider $\hat{\delta}_k = \delta_k^* - \frac{1}{K} \sum_{l=1}^K \delta_l^*$ so that we have $\sum_{l=1}^K \hat{\delta}_l = 0$

However, it is strange to have the constraint $\delta_K = 0$ while it is not required. To retain the symmetry (i.e. to make the role of the K^{th} variable the same as the others), we can consider the average of the solution under each constraint $\delta_i = 0$. That is, let $\delta_k^* = \frac{\sum_{i \in R} (y_{ik} - p_{ik})}{\sum_{i \in R} p_{ik} (1 - p_{ik})}$ for $1 \leq k \leq K$

for $\delta_1 = 0$, we have $(0, \delta_2^*, \delta_3^*, \dots, \delta_K^*)$

for $\delta_2 = 0$, we have $(\delta_1^*, 0, \delta_3^*, \dots, \delta_K^*)$

\vdots

for $\delta_K = 0$, we have $(\delta_1^*, \delta_2^*, \delta_3^*, \dots, 0)$

Note that they are solutions for the same minimization problem with different constraints $\delta_i = 0$. We take the average to get $\frac{K-1}{K} \delta^*$ and then we consider

$$\hat{\delta}_k = \frac{K-1}{K} \left(\delta_k^* - \frac{1}{K} \sum_{l=1}^K \delta_l^* \right)$$

to fulfill the constraint $\sum_{l=1}^K \hat{\delta}_l = 0$

Note that in each update of the k^{th} entry of $\vec{f}(x)$, we assign (6) a constant γ_k for $x \in R$, i.e. $(\vec{f}_{m+1} - \vec{f}_m)_k(x) = \gamma_k$. Suppose the domain of X is divided into R_1, R_2, \dots, R_J , and we apply the update for each R_j , then $(\vec{f}_{m+1} - \vec{f}_m)_k(x) = \gamma_{kj}$ for $x \in R_j$

$$= \sum_{j=1}^J \gamma_{kj} \mathbb{1}_{\{x \in R_j\}}$$

which can be considered as a regression tree with the regions $R_j, j=1, \dots, J$ represented by the terminal nodes of the tree. Let $\Theta = \{R_j, \gamma_j\}$ (we drop the index k as the arguments are the same for each entry of $\vec{f}(x)$) and $T(x; \Theta_m) = (f_{m+1} - f_m)(x) = \sum_{j=1}^J \gamma_j \mathbb{1}_{\{x \in R_j\}}, \quad f_m(x) \in \mathbb{R}$

We have seen how to compute γ_j 's when R_j 's are given. Next we will talk about how to choose R_j 's.

Gradient Boosting

At each step in the forward stagewise procedure, we would like to solve

$$\hat{\Theta}_m = \arg \min_{\Theta_m} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + T(x_i; \Theta_m)), \quad (1)$$

which is difficult to solve. Actually what we really want to find is a function $\hat{f}(x)$ such that $\hat{f} = \arg \min_f \sum_{i=1}^N L(y_i, f(x_i))$.

While this problem is difficult to solve, it is much easier if we consider finding $\vec{f} = (f(x_1), f(x_2), \dots, f(x_N)) \in \mathbb{R}^N$ such that $\sum_{i=1}^N L(y_i, f(x_i))$ is minimum.

For example, we can apply steepest descent method with the update

$$\vec{f}_{m+1} = \vec{f}_m - p_m \vec{g}_m, \quad p_m \in \mathbb{R}, \quad \vec{g}_m \in \mathbb{R}^N$$

The i th component of \vec{g}_m is

$$(\vec{g}_m)_i = \left(\nabla \sum_{i=1}^N L(y_i, f(x_i)) \right)_i \Big|_{\vec{f} = \vec{f}_m} = \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right] f(x_i) = f_m(x_i)$$

The step length p_m is the solution of $\min_p \sum_{i=1}^N L(y_i, f_m(x_i) - p(\vec{g}_m)_i)$

If we stop at $m=M$, the output is a vector

(7)

$$\vec{f}_M = -\rho_M \vec{g}_M - \rho_{M-1} \vec{g}_{M-1} - \dots - \rho_0 \vec{g}_0 = (f_M(x_1), f_M(x_2), \dots, f_M(x_N))^T$$

but what we want is a function $f(x)$ for a new $x \notin \{x_1, \dots, x_N\}$.

Consider the two updates $\vec{f}_{m+1} = \vec{f}_m - \rho_m \vec{g}_m$ and $f_{m+1}(x) = f_m(x) + T(x; \Theta_m)$

The idea of gradient boosting is to choose $T(x; \Theta_m)$ such that

$\| -\rho_m \vec{g}_m - (T(x_1, \Theta_m), \dots, T(x_N, \Theta_m))^T \|_2^2$ is small. That is to choose

$$\begin{aligned} \hat{\Theta}_m &= \arg \min_{\Theta} \sum_{i=1}^N (-\rho_m (\vec{g}_m)_i - T(x_i; \Theta))^2 \\ &= \arg \min_{\{R_j, \gamma_j\}} \sum_{i=1}^N (-\rho_m (\vec{g}_m)_i - \sum_{j=1}^J \gamma_j \mathbb{1}_{\{x_i \in R_j\}})^2 \\ &= \arg \min_{\{R_j, \gamma_j\}} \sum_{i=1}^N (-(\vec{g}_m)_i - \sum_{j=1}^J \frac{\gamma_j}{\rho_m} \mathbb{1}_{\{x_i \in R_j\}})^2 \quad \text{--- (2)} \end{aligned}$$

Note that for $L(y_i, \vec{f}(x_i)) = -\sum_{k=1}^K \mathbb{1}_{\{y_i = f_k\}} f_k(x_i) + \log(\sum_{k=1}^K e^{f_k(x_i)})$

$$\frac{\partial L}{\partial f_k(x_i)} = -\mathbb{1}_{\{y_i = f_k\}} + \frac{e^{f_k(x_i)}}{\sum_{k=1}^K e^{f_k(x_i)}} = -y_{ik} + p_k(x_i)$$

Therefore, for the update of $f_k(x)$, we can first fit a regression trees on the data $(x_i, -(\vec{g}_m)_i) = (x_i, y_{ik} - p_k(x_i)), i=1, \dots, N$ to get a partition \tilde{R}_{jm} . Although the solution regions \tilde{R}_{jm} to (2) will not be identical to the regions R_{jm} that solve (1), it is generally similar enough to serve the same purpose.

Now, with given $\tilde{R}_{1m}, \tilde{R}_{2m}, \dots, \tilde{R}_{Jm}$, we can find γ_{jm} by

$$\gamma_{jm} = \arg \min_{\gamma} - \sum_{x_i \in \tilde{R}_{jm}} \sum_{k=1}^K y_{ik} (f_k(x_i) + \gamma_k) + \sum_{x_i \in \tilde{R}_{jm}} \log(\sum_{k=1}^K e^{f_k(x_i) + \gamma_k})$$

as before.

Gradient Boosting for K-class Classification

(8)

1. Initialize $f_{k0}(x) = 0$, $k = 1, 2, \dots, K$

2. For $m = 1$ to M

(a) Set $p_k(x) = \frac{e^{f_k(x)}}{\sum_{k=1}^K e^{f_k(x)}}$, $k = 1, 2, \dots, K$

(b) For $k = 1$ to K :

i. Compute $r_{ikm} = y_{ik} - p_k(X_i)$, $i = 1, \dots, N$

ii. Fit a regression tree to the targets r_{ikm} , $i = 1, 2, \dots, N$, giving terminal regions R_{jkm} , $j = 1, 2, \dots, J_m$

iii. Compute

$$\gamma_{jkm} = \frac{K-1}{K} \frac{\sum_{i \in R_{jkm}} r_{ikm}}{\sum_{i \in R_{jkm}} |r_{ikm}| (1 - |r_{ikm}|)}, \quad j = 1, 2, \dots, J_m$$

iv. Update $f_{km}(x) = f_{k,m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jkm} \mathbb{1}_{\{x \in R_{jkm}\}}$

3. Output $\hat{f}_k(x) = f_{kM}(x)$, $k = 1, 2, \dots, K$

Gradient boosting for linear regression

Suppose the given observations are (\vec{X}_i, y_i) , $\vec{X}_i = (X_{i1}, \dots, X_{ip})^T$, $i = 1, \dots, N$.

We assume a linear model $y = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$, with $E(\varepsilon) = 0$ and ε is independent of X_1, \dots, X_p . If $p > n$, then we cannot use least squares estimator to estimate β_1, \dots, β_p . However, if we knew that most of the β_j 's are zeros, we can estimate β_1, \dots, β_p by some variable selection methods. One of them is boosting

1. Initialize $f_0(x) = 0$

2. For $m = 1$ to M

(a) Compute $(\alpha_m, j_m) = \arg \min_{\alpha, j} \sum_{i=1}^N L(y_i, f_{m-1}(X_i) + \alpha X_j)$ $\alpha \in \mathbb{R}$, $j \in \{1, \dots, p\}$
 $= \arg \min_{\alpha, j} \sum_{i=1}^N \frac{1}{2} (y_i - f_{m-1}(X_i) - \alpha X_j)^2$

(b) Set $f_m(x) = f_{m-1}(x) + \alpha_m X_{j_m}$

The final outcome is $f_M(x) = \sum_{m=0}^M \alpha_m X_{j_m}$

(9)

By gradient boosting, the step 2(a) can be approximately done by first choosing index j_{m+1} that minimize $\min_{\alpha} \|\vec{g}_m - \alpha X_j\|^2$

where $(\vec{g}_m)_i = \frac{\partial}{\partial f(x_i)} \sum_{i=1}^N \frac{1}{2} (y_i - f(x_i))^2 \Big|_{f=f_m} = -(y_i - f_m(x_i))$ and

$$\|\vec{g}_m - \alpha X_j\|^2 = \sum_{i=1}^N (y_i - f_m(x_i) - \alpha X_{ij})^2 = \sum_{i=1}^N (y_i - f_m(x_i))^2 - 2\alpha \sum_{i=1}^N X_{ij} (y_i - f_m(x_i)) + \alpha^2 \sum_{i=1}^N X_{ij}^2$$

$$\Rightarrow \tilde{\alpha} = \arg \min_{\alpha} \|\vec{g}_m - \alpha X_j\|^2 = \frac{\sum_{i=1}^N X_{ij} (y_i - f_m(x_i))}{\|X_j\|^2}$$

$$\text{and } \|\vec{g}_m - \tilde{\alpha} X_j\|^2 = \sum_{i=1}^N (y_i - f_m(x_i))^2 - 2 \frac{(\sum_{i=1}^N X_{ij} (y_i - f_m(x_i)))^2}{\|X_j\|^2} + \frac{(\sum_{i=1}^N X_{ij} (y_i - f_m(x_i)))^2}{\|X_j\|^2}$$

Therefore, we first choose $j_{m+1} = \arg \max_j \frac{|\sum_{i=1}^N X_{ij} (y_i - f_m(x_i))|}{\|X_j\|}$

And then we estimate $\alpha_{m+1} = \arg \min_{\alpha} \sum_{i=1}^N \frac{1}{2} (y_i - f_m(x_i) - \alpha X_{i,j_{m+1}})^2$
 $= \frac{X_{j_{m+1}}^T U_m}{\|X_{j_{m+1}}\|^2}$ where $(U_m)_i = y_i - f_m(x_i)$

This is called L_2 -Boosting in Bühlmann and Yu (2003)

1. Initialize $f_0(x) = 0$, $U_0 = Y = (y_1, \dots, y_N)^T$

2. For $m=1$ to M

(a) Choose $j_m = \arg \max_j \frac{|X_j^T U_{m-1}|}{\|X_j\|^2}$

(b) Compute $\alpha_m = \frac{X_{j_m}^T U_{m-1}}{\|X_{j_m}\|^2}$

(c) Update $U_m = U_{m-1} - \alpha_m X_{j_m}$
 $f_m = f_{m-1} + \alpha_m X_{j_m}$

Ing and Lai (2011) point out that it is unclear how to choose the upper bound M on the number of iterations as the same predictor variable can be entered at several iterations. For variable selection procedure, a common stopping criterion is stopping the algorithm when no variables are selected. For L_2 -Boosting, since the same predictor can enter again, it may

take much more iterations to select m_0 predictors, or the error $\sum_{i=1}^n \frac{1}{2} (y_i - f_m(x_i))^2$ is still large after m_0 predictors are selected.

Note that if X_j 's are orthogonal, then the least square estimator of $\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$ (assume $p < n$) is $\hat{\beta} = (X^T X)^{-1} X^T Y = \left(\frac{X_j^T Y}{\|X_j\|^2} \right)_{j=1, \dots, p}$

In 2(b), $\alpha_m = \frac{X_{j_m}^T (Y - \alpha_1 X_{j_1} - \dots - \alpha_{m-1} X_{j_{m-1}})}{\|X_{j_m}\|^2} = \frac{X_{j_m}^T Y}{\|X_{j_m}\|^2} = \hat{\beta}_{j_m}$ if $j_m \notin \{j_1, \dots, j_{m-1}\}$

And for $\alpha_k = \hat{\beta}_{j_k}$, for $k=1, \dots, m-1$, suppose j_1, \dots, j_{m-1} are distinct, then $j_m \notin \{j_1, \dots, j_{m-1}\}$.

Suppose not, WLOG, let $j_m = j_1$, then

$$\begin{aligned} X_{j_m}^T U_{m-1} &= X_{j_m}^T (Y - \alpha_1 X_{j_1} - \dots - \alpha_{m-1} X_{j_{m-1}}) \\ &= X_{j_m}^T Y - \alpha_1 X_{j_m}^T X_{j_1} \\ &= X_{j_1}^T Y - \frac{X_{j_1}^T Y}{\|X_{j_1}\|^2} \|X_{j_1}\|^2 = 0 \end{aligned}$$

$\Rightarrow j_m$ cannot be selected.

Using this idea, Ing and Lai (2011) propose modifying the L_2 -Boosting by keep orthogonalizing the selected variables.

Orthogonal Greedy Algorithm (OGA)

1. Initialize $U_0 = Y$, $I_0 = \emptyset$

2. For $m=1$ to M

(a) Choose $j_m = \arg \max_j \frac{|X_j^T U_{m-1}|}{\|X_j\|^2}$

(b) Update $I_m = I_{m-1} \cup \{j_m\}$ and compute the QR decomposition

$$X_{I_m} = [X_{I_{m-1}} \ X_{j_m}] = [Q_{m-1} \ q_m] \begin{bmatrix} R_{m-1} \\ 0^T \ r_m \end{bmatrix} = Q_m R_m$$

Compute $\alpha_m = q_m^T U_{m-1}$

(c) Update $U_m = U_{m-1} - \alpha_m q_m$

Output $f_m(x) = \hat{\beta}_{j_1} X_{j_1} + \dots + \hat{\beta}_{j_m} X_{j_m}$, where $\hat{\beta}_{j_m}$ is the m th entry of $R_m^{-1} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{pmatrix}$.

Ing and Lai (2011) suggest $M = O\left(\sqrt{\frac{n}{\log p}}\right)$ under some sparsity (11) conditions. Let $J_m = \{j_1, j_2, \dots, j_m\}$, they propose a high dimensional information criterion (HDIC) to further select the variables in J_m .

First consider
$$HDIC(m) = N \log \|U_m\|^2 + mC \log(p)$$
$$\hat{m} = \arg \min_{1 \leq m \leq M} HDIC(m)$$

where C is a constant independent of m .

The final selected set is

$$\hat{J} = \{j_\ell : HDIC(J_{\hat{m}} \setminus \{j_\ell\}) > HDIC(J_{\hat{m}}), 1 \leq \ell \leq \hat{m}\}$$

if $\hat{m} > 1$.

Variable Importance

Variable importance plots can be constructed in exactly the same way as they were for random forests for regression trees.

$$I_\ell^2(T) = \sum_{t=1}^{J-1} \hat{t}_t^2 I(v(t) = \ell) \quad I_\ell^2 = \frac{1}{M} \sum_{m=1}^M I_\ell^2(T_m)$$

For K -classification, K separate models $f_k(x)$, $k=1, 2, \dots, K$ are induced, each consisting of a sum of trees

$$f_k(x) = \sum_{m=1}^M T_{km}(x)$$

In this case, we define $I_{\ell k} = \sqrt{\frac{1}{M} \sum_{m=1}^M I_\ell^2(T_{km})}$, which is the relevance of X_ℓ in separating the class k observations from other classes.

The overall relevance of X_ℓ is obtained by averaging over all the classes

$$I_\ell^2 = \frac{1}{K} \sum_{k=1}^K I_{\ell k}^2$$