

## 集成学习

Bagging (bootstrap aggregating的缩写)  $\Rightarrow$  降低方差.

### 1) Random Forest

- 行随机: 使用bootstrap方法有返回的采样  $m$  个包含  $n$  个样本的数据集.
- 列随机: 从所有属性中选择  $m$  个属性进行比较, 选举胜出.
- 每棵树都可以生长, 没有剪枝过程.

特点:

- 1° 有较好的鲁棒性.
- 2° 可以有效处理高维特征, 不需要降维.
- 3° 对缺失值不敏感.

Random Forest 错误率相关因素:  $\rightarrow$  如何选择  $m$  个属性, RF中的唯一参数.

- 1° 森林中任意两树之间的相关性, 相关性低, 错误率低.
- 2° 每棵树的分类能力, 错误率低.

选择  $m$   $\Rightarrow$  袋外错误率 (oob error)

随机森林不需要交叉验证或额外的测试集, 因为可以内部估计.

bootstrap采样方法使得对每棵树, 约有 0.368 数据没有用于训练, 它们称为第  $k$  棵树的 oob 样本.

- 1) 对于每个样本, 计算其作为 oob 样本的树对它的分类情况.
- 2) 使用多数投票作为预测结果.
- 3) 用该样本个数 / 总样本个数作为模型的 oob error.

$\hookrightarrow$  oob error 是 RF 的近似无偏估计, 近似于重复大量计算的  $k$  折交叉验证.

# Boosting $\Rightarrow$ 降低偏差

1) AdaBoost 思想: 比较容易地找到弱分类器, 然后通过反复学习得到一个强分类器。

训练模型: 强分类器由若干弱分类器的线性组合得到:  $F_m = \sum_{n=1}^M \alpha_n h(x)$

前向分类算法: T-轮迭代产生的分类器由上一轮迭代得到:  $F_{m+1} = F_m + \alpha_{m+1} h(x)$

## 1) AdaBoost.

1° 初始时, 假设所有样本权重相同,

2° 每次训练完成后, 对分类错误的样本赋予较大的权重, 让模型更关注该样本。

3° 根据分类器的错误率赋予权重, 采用多数表决的方法。

假设数据集:  $T = \{(x_1, y_1), \dots, (x_n, y_n)\}$ ,  $y \in \{-1, +1\}$ ,

(1) 初始化权重分布:

$$w_1 = (w_{11}, w_{12}, \dots, w_{1n}), \quad w_{1i} = \frac{1}{n}$$

(2) 对于 m 次迭代:

(a) 得到分类器  $G_m(x)$ .

(b) 计算  $G_m(x)$  的错误率:  $e_m = \sum_{i=1}^n w_{mi} I(G_m(x_i) \neq y_i)$

分类器分类错误率  
权重变小。

(c) 计算  $G_m(x)$  在弱分类器中的权重:  $\alpha_m = \frac{1}{2} \log \frac{1-e_m}{e_m} = \frac{1}{2} \log \left( \frac{1}{e_m} - 1 \right)$

(d) 更新样本的权重分布:

$$w_{m+1,i} = \frac{w_{m,i}}{Z_m} \exp(-\alpha_m y_i G_m(x_i)) \quad \left\{ \begin{array}{l} w_{m,i} e^{-\alpha_m} \Rightarrow \text{分类对, } -\alpha_m < 0, \\ w_{m,i} e^{+\alpha_m} \Rightarrow \text{分类错, } -\alpha_m > 0, \end{array} \right.$$

规范化因子。

$w_{m,i} e^{-\alpha_m}$  变小,

(3) 得到强分类器

$$F(x) = \text{sign} \left( \sum_{m=1}^M \alpha_m G_m(x) \right)$$

则分类对的样本在下次迭代中权重变小。

Adaboost 的训练误差上界:

$$\frac{1}{N} \sum_{i=1}^N I(G(x_i) \neq y_i) \leq \frac{1}{N} \sum_{i=1}^N \exp(-y_i f(x_i)) = \frac{1}{N} Z_M$$

证明:

1) 若  $G(x_i) \neq y_i$ , 则  $y_i f(x_i) < 0$ , 有  $\exp(-y_i f(x_i)) \geq 1$

$$\therefore \frac{1}{N} \sum_{i=1}^N I(G(x_i) \neq y_i) \leq \frac{1}{N} \sum_{i=1}^N \exp(-y_i f(x_i)) \text{ 成立.}$$

2)  $\because w_{m+1,i} = \frac{w_m}{Z_m} \exp(-\alpha_m y_i G_m(x_i))$

$$\therefore Z_m w_{m+1,i} = w_m \exp(-\alpha_m y_i G_m(x_i))$$

$$\therefore \frac{1}{N} \sum_{i=1}^N \exp(-y_i f(x_i))$$

$$= \frac{1}{N} \sum_{i=1}^N \exp(-\sum_{m=1}^M \alpha_m y_i G_m(x_i))$$

$$= \sum_{i=1}^N w_{1,i} \prod_{m=1}^M \exp(-\alpha_m y_i G_m(x_i)) \Rightarrow w_{1,i} = \frac{1}{N}$$

$$= \sum_{i=1}^N w_{1,i} \exp(-\alpha_1 y_i G_1(x_i)) \prod_{m=2}^M \exp(-\alpha_m y_i G_m(x_i))$$

$$= \sum_{i=1}^N Z_1 w_{2,i} \prod_{m=2}^M \exp(-\alpha_m y_i G_m(x_i)) \Rightarrow Z_m w_{m+1,i} = w_{m,i} \exp(-\alpha_m y_i G_m(x_i))$$

$$= Z_1 \sum_{i=1}^N w_{2,i} \prod_{m=2}^M \exp(-\alpha_m y_i G_m(x_i))$$

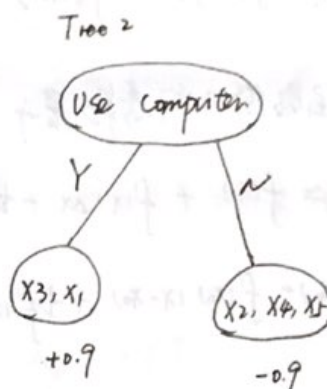
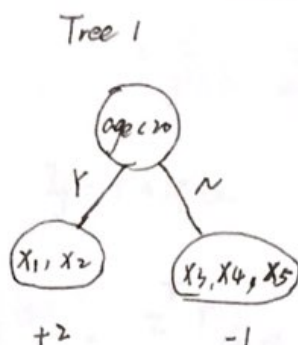
$$= Z_1 Z_2 \cdots Z_{m-1} \sum_{i=1}^N w_{m,i} \exp(-\alpha_m y_i G_m(x_i))$$

$$= \prod_{m=1}^M Z_m$$



## 2) XGBOOST

### 1° Decision Tree Ensembles



$$\Rightarrow f(x_1) = 2 + 0.9 = 2.9 \quad f(x_5) = -1 - 0.9 = -1.9$$

模型:  $\hat{y}_i = \sum_{k=1}^K f_k(x_i), f_k \in F$

目标函数:  $obj(\theta) = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$

### 2° Tree Boosting

目标函数:  $obj^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t)}) + \sum_{i=1}^t \Omega(f_i)$

→ 需要学习(t)棵树, 但是无法一次学习所有树结构是很难的

→ 迭代模型

$$\left\{ \begin{array}{l} \hat{y}_i^{(0)} = 0 \\ \hat{y}_i^{(1)} = f_1(x_i) = \hat{y}_i^{(0)} + f_1(x_i) \\ \hat{y}_i^{(2)} = f_1(x_i) + f_2(x_i) = \hat{y}_i^{(1)} + f_2(x_i) \\ \vdots \\ \hat{y}_i^{(t)} = \sum_{k=1}^t f_k(x_i) = \hat{y}_i^{(t-1)} + f_t(x_i) \end{array} \right.$$

$$\therefore \text{obj}^{(t)} = \sum_{i=1}^n \ell(y_i, \hat{y}_i^{(t)}) + \sum_{i=1}^t \omega(f_i)$$

$$= \sum_{i=1}^n \ell(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \omega(f_t) + \text{constant}$$

泰勒展开  $\Rightarrow$  将损失函数做二阶泰勒展开:

$$f(x+\Delta x) \approx f(x) + f'(x)\Delta x + \frac{1}{2}f''(x)\Delta x^2$$

$$f(x) \approx f(x_0) + f'(x_0)(x-x_0) + \frac{1}{2}f''(x_0)(x-x_0)^2$$



$$\therefore \text{obj}^{(t)} = \sum_{i=1}^n \left[ \ell(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \omega(f_t) + \text{constant}$$

$$g_i = \partial_{\hat{y}_i^{(t-1)}} \ell(y_i, \hat{y}_i^{(t-1)})$$

$$h_i = \partial_{\hat{y}_i^{(t-1)}}^2 \ell(y_i, \hat{y}_i^{(t-1)})$$

$$\Rightarrow \text{obj}^{(t)} = \sum_{i=1}^n \left[ g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \omega(f_t)$$

### 3° Model Complexity

$\rightarrow$  将模型看作一个叶节点的集合

假设定义  $f_t(x) = w q_t(x)$ ,  $w \in \mathbb{R}^T$ ,  $q: \mathbb{R}^d \rightarrow \{1, 2, \dots, T\}$ .

$\downarrow$   
T维的向量, 表示每个叶节点的得分.

则模型复杂度的定义为:  $\omega(f) = \|w\|_1 + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$

~~表示你使用~~

#### 4° Structure Score

$$\begin{aligned} obj^{(t)} &\approx \sum_{i=1}^n [g_i w_{q(x_i)} + \frac{1}{2} h_i w_{q(x_i)}^2] + YT + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \\ &= \sum_{j=1}^T [(\sum_{i \in I_j} g_i) w_j + \frac{1}{2} (\sum_{i \in I_j} h_i + \lambda) w_j^2] + YT \end{aligned}$$

$I_j = \{i | q(x_i) = j\}$  表示属于第  $j$  个叶节点的样本下标。

$$\Rightarrow obj^{(t)} = \sum_{j=1}^T [G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2] + YT$$

求导得到： $w_j^* = -\frac{G_j}{H_j + \lambda} \Rightarrow$  实际使用时会乘以学习率防止过拟合。

$$obj^* = -\frac{1}{2} \sum_{j=1}^T \frac{G_j^2}{H_j + \lambda} + YT$$

#### 5° 属性选择

$$Gain = \frac{1}{2} \left[ \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - Y$$

↓ 节点划分后期望收益比增加一个节点的代价 ( $\lambda$ ) 高

### 使用例子说明:

假设有10个样本, 2个特征,  $y \in \{0, 1\}$ 。使用 CART 作为基学习器, 树的节点最大深度为3, 模型学习率为0.1, 正则化参数  $\lambda = 1, \gamma = 0$ 。损失函数使用  $\log$  loss。

1. 损失函数

$$l(y_i, \hat{y}_i) = y_i \ln(1 + e^{-\hat{y}_i}) + (1 - y_i) \ln(1 + e^{\hat{y}_i})$$

$$\begin{aligned} \therefore \frac{\partial l(y_i, \hat{y}_i)}{\partial \hat{y}_i} &= y_i \frac{-e^{-\hat{y}_i}}{1 + e^{-\hat{y}_i}} + (1 - y_i) \frac{e^{\hat{y}_i}}{1 + e^{\hat{y}_i}} \\ &= y_i \frac{-1}{1 + e^{\hat{y}_i}} + (1 - y_i) \frac{1}{1 + e^{-\hat{y}_i}} \end{aligned}$$

$$= y_i \left( \frac{1}{1 + e^{\hat{y}_i}} - 1 \right) + (1 - y_i) \frac{1}{1 + e^{-\hat{y}_i}}$$

$$= y_i (y_{i, \text{pred}} - 1) + (1 - y_i) y_{i, \text{pred}}$$

$$= y_{i, \text{pred}} - y_i, \Rightarrow g_i$$

$$\therefore \frac{\partial^2 l(y_i, \hat{y}_i)}{\partial \hat{y}_i^2} = y_{i, \text{pred}} (1 - y_{i, \text{pred}}) \Rightarrow h_i$$

2. 计算每一个样本的  $g_i$  和  $h_i$ 。

1° 每个样本之间相互独立, 可以并行  $\Rightarrow$  XGBOOST 快。

2° 每次迭代不需重复计算  $g_i$  和  $h_i$

\* 第一次迭代时,  $y_{i, \text{pred}}$  初始化为 0.5,  $\Rightarrow$  至少是  $\{0, 1\}$  中任意一个

3. 遍历每一个特征，对特征值排序，然后依次从左右两侧为划分点计算 Gain 值。

时间复杂度： $m \cdot n \log n$ 。找到 Gain 最大的划分。

不断重复上述步骤，直到达到最大深度。

4. 对于叶子节点，其对应的权重值为  $w^* = -\frac{G_j}{H_j + \lambda}$

实际上，叶子节点的取值还需要乘以一个学习率  $\Rightarrow w^* \cdot \text{学习率}$ 。

这是为了防止过拟合。

5. 接着构造第二棵树。它的样本对应的初始  $y_{i, \text{pred}}$  由下式给出：

$$\begin{aligned} \text{~~Y_i~~} &= \hat{y}_i^{(1)} = f_0(x_i) + f_1(x_i) = 0 + w_{g_1}(x_i) \\ y_{i, \text{pred}} &= \frac{1}{1 + e^{-\hat{y}_i^{(1)}}} \\ 0.5 &= \frac{1}{1 + e^{-f_0(x_i)}} \Rightarrow f_0(x_i) = 0 \end{aligned}$$