

# Boosting 减少偏差

弱可学习算法  $\rightarrow$  强学习算法

Adaboost

整体数据如：提升该分类器权重  
单个样本如：降低该样本权重

算法流程

1. 初始化权重分布 (均匀分布)  
 $D_1 = (w_{11}, \dots, w_{1i}, \dots, w_{1N}) \quad w_{1i} = \frac{1}{N}$
  2. 对  $m = 1, \dots, M$ 
    - (1) 使用权重分布  $D_m$  训练得到基本分类器  $G_m$
    - (2) 计算  $G_m$  在训练数据上的误差率  $e_m = \sum_{i=1}^N I(G_m(x_i) \neq y_i) w_{mi}$
    - (3) 计算  $G_m(x)$  在训练数据集上的分类误差  $\alpha_m = \frac{1}{2} \log \frac{1-e_m}{e_m}$
    - (4) 更新样本权重分布  $D_{m+1} = (w_{m+1,1}, \dots, w_{m+1,N})$   

$$w_{m+1,i} = \frac{w_{mi}}{Z_m} \exp(-\alpha_m y_i G_m(x_i)) \quad \left| \begin{array}{l} \frac{w_{mi}}{Z_m} e^{-\alpha_m} \quad G_m = y_i \\ \frac{w_{mi}}{Z_m} e^{\alpha_m} \quad G_m \neq y_i \end{array} \right.$$
- $Z_m$ : 归一化因子  $= \sum_{i=1}^N w_{mi} \exp(-\alpha_m y_i G_m(x_i))$   
 本轮样本权重之和 (权重之和 = 1)
3. 构造基本分类器线性组合  

$$f(x) = \sum_{m=1}^M \alpha_m G_m(x) \rightarrow G(x) = \text{sign}(f(x))$$

Bagging  $\rightarrow$  减少方差

$$\sigma^2 = \rho \sigma^2 + (1-\rho) \frac{\sigma^2}{n}$$

相关系数

提升树 Boosting tree

• 基分类器：分类 / 回归树

• 提升树模型  

$$f_m(x) = \sum_{m=1}^M T(x; \theta_m)$$

• 前向分步算法

$$f_m(x) = f_{m-1}(x) + T(x; \theta_m)$$

消除残差

$$\hat{\theta}_m = \arg \min_{\theta} \sum_{i=1}^N L(y_i, \underbrace{f_{m-1}(x_i)}_{\text{固定}} + T(x_i; \theta_m))$$

$$f_0(x) = 0$$

$$f_1(x) = f_0(x) + T(x; \theta_1)$$

$$\frac{1}{N} \sum_{i=1}^N L(y_i, f_1(x_i))$$

对于回归问题：

$$L(y_i, f_1(x_i)) = \frac{1}{2} (y_i - f_1(x_i))^2$$

## 梯度提升算法

输入: 训练数据集  $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$   
损失函数  $L(y, f(x))$

输出: 回归树  $f(x)$

1) 初始化

$$f_0(x) = \arg \min_c \sum_{i=1}^N L(y_i, c)$$

2) 对于  $m=1, \dots, M$

(a) 对  $i=1, \dots, N$  计算

$$r_{mi} = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right] \quad f(x) = f_{m-1}(x)$$

残差近似

$$- \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \approx [y_i - f(x_i)] \quad \text{残差}$$

(b) 对  $r_{mi}$  拟合一个回归树, 得到第  $m$  棵树的叶节点区域  $R_{mj} \quad j=1, \dots, J$

(c) 对  $j=1, \dots, J$  计算  $C_{mj} = \arg \min_c \sum_{x \in R_{mj}} L(y_i, f_{m-1}(x) + c)$

(d) 更新  $f_m(x) = f_{m-1}(x) + \sum_{j=1}^J C_{mj} I(x \in R_{mj})$

3) 得到回归树

$$\hat{f}(x) = f_M(x) = \sum_{m=1}^M \sum_{j=1}^J C_{mj} I(x \in R_{mj})$$

## Adaboost

### 前向分步算法

$$\begin{aligned} f_m(x) &= f_{m-1}(x) + \alpha_m G_m(x) \\ (\alpha_m, G_m(x)) &= \arg \min_{\alpha, G} \sum_{i=1}^N \exp[-y_i (f_{m-1}(x_i) + \alpha G_m(x_i))] \\ &= \arg \min_{\alpha, G} \sum_{i=1}^N \bar{w}_{mi} \exp[-y_i \alpha G_m(x_i)] \\ &= \arg \min_{\alpha, G} \sum_{i=1}^N \bar{w}_{mi} \exp[-\alpha_m] \\ &\quad + \arg \min_{\alpha, G} \sum_{i=1}^N \bar{w}_{mi} \exp[\alpha_m] \\ &= \arg \min_{\alpha, G} \sum_{i=1}^N \bar{w}_{mi} \exp[-\alpha_m] \\ &\quad + \arg \min_{\alpha, G} \sum_{i=1}^N \bar{w}_{mi} \exp[\alpha_m] \\ &= \exp[-\alpha_m] \sum_{i=1}^N \bar{w}_{mi} + (e^{\alpha_m} - e^{-\alpha_m}) \sum_{i=1}^N \bar{w}_{mi} I(y_i \neq G_m(x_i)) \quad \text{①} \end{aligned}$$

$G$  的最优解

$$G_m^* = \arg \min_G \sum_{i=1}^N \bar{w}_{mi} I(y_i \neq G(x_i))$$

$\alpha$  的最优解

对 0 求导,  $\frac{\partial 0}{\partial \alpha} = 0$

$$-e^{-\alpha} \sum_{i=1}^N \bar{w}_{mi} + (e^{\alpha} + e^{-\alpha}) \sum_{i=1}^N \bar{w}_{mi} I(y_i \neq G_m(x_i)) = 0$$

$$\frac{e^{\alpha} + e^{-\alpha}}{e^{-\alpha}} = \frac{1 / \sum \bar{w}_{mi} I(y_i \neq G_m(x_i))}{1 / \sum \bar{w}_{mi}} = \frac{1}{e^{\alpha}}$$

$$e^{\alpha} = \frac{1}{e^{\alpha}} - 1$$

$$\alpha = \frac{1}{2} \log \frac{1}{e^{\alpha}}$$

权重更新

$$\begin{aligned} \bar{w}_{mi} &= \exp(-y_i f_{m-1}(x_i)) \\ &= \exp(-y_i \sum_{j=1}^{m-1} \alpha_j G_j(x_i)) \\ &= \prod_j \exp(-y_i \alpha_j G_j(x_i)) \end{aligned}$$

$\therefore \bar{w}_{m+1,i}$

$$= \bar{w}_{mi} \cdot \exp(-y_i \alpha_m G_m(x_i))$$