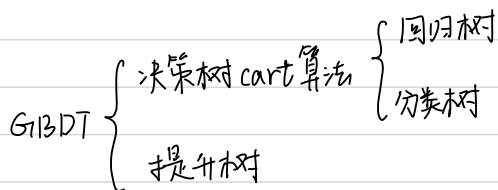


GBDT



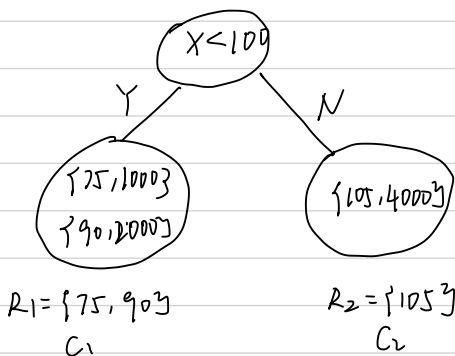
1. 回归树

1.1. $f(x) = \sum_{m=1}^M C_m I(x \in R_m)$ m: 叶子节点中的序号, 第几个叶子节点,
Cm: 每个叶子节点对应的树

x: 房屋面积

y: 租金

$$D = \{(75, 1000), (90, 2000), (105, 4000)\}$$



注: 每个叶子节点, 都会输出一个预测值

预测值一般是该叶子所含训练样本在该点上输出的均值

1.2. 如何构建树:

- (1) 树的深度: { 确定叶子节点个数 or 树的深度
子节点包含样本数
给定精度 (loss)

(2) 划分节点, 如何选 找到损失最小的树的划分情况

(3) 叶子节点代表的 C_m 如何定 C_m 为平均值 (每个叶子节点)

划分条件 — C_m — 损失

1.3. 损失函数

优化求解

$$\frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 \rightarrow \text{按样本方式遍历 (同样也可以按节点)}$$

按样本遍历

$$J = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$$

$$= \frac{1}{n} \sum_{m=1}^M (C_m I(x_i \in R_m) - y_i)^2$$

$$= \frac{1}{n} \sum_{m=1}^M \sum_{x_i \in R_m} (C_m - y_i)^2 \rightarrow \text{其中 对于第 } m \text{ 个节点, 遍历所有的样本}$$

按节点遍历

$$C_m^* = \min_{C_m} \frac{1}{n} \sum_{x_i \in R_m} (C_m - y_i)^2 \quad (\text{求导令导数}=0, \text{从而求解 } C_m^*)$$

$$\frac{\partial J}{\partial C_m} = \frac{\partial \frac{1}{n} \sum_{x_i \in R_m} (C_m - y_i)^2}{\partial C_m} = \frac{\partial \sum_{x_i \in R_m} (C_m - y_i)^2}{\partial C_m} \quad (C_m \text{ 解完后, } \sum_{x_i \in R_m} \text{ 为常数})$$

$$= 2 \sum_{x_i \in R_m} (C_m - y_i)$$

设叶子节点 R_m 包含样本 N_m 个

$$= N_m C_m - \sum_{x_i \in R_m} y_i = 0$$

$$C_m = \frac{\sum_{x_i \in R_m} y_i}{N_m} \Rightarrow \text{每个叶子节点的含有样本的均值}$$

即当每个叶子节点的 C_m 取值为该节点所有样本 y_i 的均值时, 损失最小

2. 分类树

· 二叉树, 使用基尼指数作为损失函数

基尼指数: (K 类, 样本点属于第 k 类的概率为 p_k)

$$Gini(p) = \sum_{k=1}^K p_k (1 - p_k) = 1 - \sum_{k=1}^K p_k^2$$

对于二分类问题, 样本点属于第 1 类的概率为 p , 则概率分布的基尼指数为

$$Gini(p) = 2p(1-p)$$

对于给定样本集合 D

$$Gini(D) = 1 - \sum_{k=1}^K \left(\frac{|A_k|}{|D|} \right)^2$$

$Gini$ 很小 \rightarrow 样本基本都属于同一类别

$$Gini(D, A) = \frac{|D_1|}{|D|} Gini(D_1) + Gini(D_2) \frac{|D_2|}{|D|}$$

回归树, 分类树本质都是构建二叉树; 计算损失不同

回归 — MSE

分类 — 基尼系数

3. GBDT: 提升树被认为是统计学习中性能最好的方法之一

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

$\xrightarrow{M \text{ 为树的个数}}$
 $\xrightarrow{\theta_m: \text{决策树的参数}}$
 $T(x; \theta_m): \text{决策树}$

模型: 加法模型

损失函数:

- 回归问题 — MSE
- 分类问题
 - 二分类 指数损失
 - 多分类 softmax
- 一般决策: 自定义损失函数

优化方法: 前向分步算法

提升树算法采用前向分步算法。首先确定初始提升树 $f_0(x) = 0$, 第 m 步的模型是

$$f_m(x) = f_{m-1}(x) + T(x; \theta_m) \quad (8.25)$$

其中, $f_{m-1}(x)$ 为当前模型, 通过经验风险极小化确定下一棵决策树的参数 θ_m :

$$\hat{\theta}_m = \arg \min_{\theta_m} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + T(x_i; \theta_m)) \quad (8.26)$$

3.1. 二分类问题的提升树

相当于把 Adaboost 中的 $G(x)$ 换成二分类

$$f(x) = T(x; \theta_1) + T(x; \theta_2) + \dots + T(x; \theta_m)$$



相当于 Adaboost 的特殊情况

- ① 基分类器 $G(x)$ 限制为二分类树
- ② 权重 α_m 置为 1

使用指数损失函数更新权重 $e^{-y f(x)}$

3.2. 回归问题的提升树

损失函数: 平方误差根 $L(y, f(x)) = (y - f(x))^2$

$$\begin{aligned}
 &= [y - f_{m-1}(x) - T(x; \theta_m)]^2 \\
 &= [r - T(x; \theta_m)]^2 \quad (\text{残差} = \text{真实值} - \text{预测值})
 \end{aligned}$$

前向分步算法: 利用残差数据构建训练样本, 拟合残差新样本

$$\hat{\theta}_m = \arg \min_{\theta_m} \sum_{i=1}^N L(y^{(i)}, f_{m-1}(x^{(i)}) + T(x^{(i)}; \theta_m))^2$$

$$= \arg \min_{\theta_m} \sum_{i=1}^N (r_m^{(i)} - T(x^{(i)}; \theta_m))^2$$

用残差训练当前学习器

3.3. GBDT (梯度提升树)

分类 — 指数损失 f 回归 — MSE — 残差 通解?

目标: 每增加一棵, 损失减小

$$L(y^{(i)}, f_m(x^{(i)})) < L(y^{(i)}, f_{m-1}(x^{(i)}))$$

即 $L(y^{(i)}, f_m(x^{(i)})) - L(y^{(i)}, f_{m-1}(x^{(i)})) > 0$ 满足此式说明损失越来越小

一阶泰勒: $f(x) = f(x_0) + f'(x_0)(x - x_0)$

$L(y, f_m(x))$ 中只有 $f_m(x)$ 是未知量

$$L(y, f_m(x)) \approx L(y, f_{m-1}(x)) + \left. \frac{\partial L(y, f_m(x))}{\partial f_m(x)} \right|_{f_m(x)=f_{m-1}(x)} (f_m(x) - f_{m-1}(x))$$

$$\approx L(y, f_{m-1}(x)) + \left. \frac{\partial L(y, f_m(x))}{\partial f_m(x)} \right|_{f_m(x)=f_{m-1}(x)} \cdot T(x; \theta_m)$$

$$\text{即: } L(y, f_m(x)) - L(y, f_{m-1}(x)) \approx - \left. \frac{\partial L(y, f_m(x))}{\partial f_m(x)} \right|_{f_m(x)=f_{m-1}(x)} T(x; \theta_m)$$

$$\text{当 } T(x; \theta_m) \approx - \left. \frac{\partial L(y, f_m(x))}{\partial f_m(x)} \right|_{f_m(x)=f_{m-1}(x)} \text{ 时}$$

$$L(y, f_{m-1}(x)) - L(y, f_m(x)) > 0$$

$$r_m(x, y) = - \left[\left. \frac{\partial L(y, f_m(x))}{\partial f_m(x)} \right|_{f_m(x)=f_{m-1}(x)} \right]$$

将 (x_i, y_i) 代入 $r_m(x, y)$, 即可得 r_{mi}

进而得到第 m 轮的训练数据集

$$T_m = \{(x_1, r_{m1}), (x_2, r_{m2}), \dots, (x_N, r_{mN})\}$$

梯度提升: 1. 计算当前损失函数的负梯度表达式 $r_m(x, y) = - \left[\left. \frac{\partial L(y, f_m(x))}{\partial f_m(x)} \right|_{f_m(x)=f_{m-1}(x)} \right]$
2. 构建新的训练样本 (第 m 轮的训练数据集)

算法流程:

算法 8.4 (梯度提升算法)

输入: 训练数据集 $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, $x_i \in \mathcal{X} \subseteq \mathbf{R}^n$, $y_i \in \mathcal{Y} \subseteq \mathbf{R}$;

损失函数 $L(y, f(x))$;

输出: 回归树 $\hat{f}(x)$ 。

(1) 初始化

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

(2) 对 $m = 1, 2, \dots, M$

(a) 对 $i = 1, 2, \dots, N$, 计算

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

(b) 对 r_{mi} 拟合一个回归树, 得到第 m 棵树的叶结点区域 R_{mj} , $j = 1, 2, \dots, J$ 。

(c) 对 $j = 1, 2, \dots, J$, 计算

$$c_{mj} = \arg \min_c \sum_{x_i \in R_{mj}} L(y_i, f_{m-1}(x_i) + c)$$

(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

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

e.g.

$$\begin{aligned} L &= \sum_{i=1}^N (y_i - c)^2 \\ &= \sum_{i=1}^N -2c(y_i - c) \\ &= \sum_{i=1}^N 2c^2 - 2c \sum_{i=1}^N y_i \\ &= 2NC^2 - 2c \sum_{i=1}^N y_i \\ NC^2 - c \sum_{i=1}^N y_i &= 0 \\ \Rightarrow C &= \dots ? \end{aligned}$$