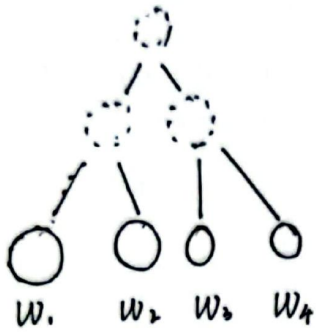




回归树



节点值

L_1, L_2, L_3, L_4 节点样本集合

$L_j = \{i \mid q(x_i) = j\}$ 节点样本集合和样本关系

$W_{q(x_i)} = T(\theta_j x_i)$ x_i 对应的值.

模型表达

$$y_i^{(t)} = \sum_{t=1}^T f_t(x_i) = \underbrace{\sum_{t=1}^{T-1} f_t(x_i)}_{y_i^{(T-1)}} + f_T(x_i)$$

前向分布算法: 贪心, 逐个优化.

5.1

$$\begin{aligned} y_i^{(t)} &= y_i^{(t-1)} + f_t(x_i) \\ &= y_i^{(t-1)} + W_{q(x_i)} \end{aligned}$$

目标函数.

$$\sum_{i=1}^N L(y_i, \hat{y}_i^{(t)}) + \sum_{t=1}^T \lambda(f_t) = Obj^{(t)}$$

优化目标.

$$(w_1, w_2, \dots, w_m) = \arg \min \sum_{i=1}^N L(y_i, \hat{y}_i^{(t)}) + \sum_{t=1}^T \lambda(f_t)$$

$$\lambda(f_t) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T W_j^2$$

T : 叶子节点个数.

original by Alkaid



目标函数

$$\sum_{i=1}^N L(y_i, \hat{y}_i^{(n)}) + \sum_{j=1}^J \mu_j f_j$$

非连续. $= \sum_{j=1}^J \mu_j f_j + \mu_0 f_0$

$$= y^T + \frac{1}{2} \lambda \sum_{j=1}^J w_j^2$$

$$\sum_{i=1}^N L(y_i, \hat{y}_i^{(n)})$$

$$= \sum_{j=1}^J \sum_{i \in I_j} L(y_i, \hat{y}_i^{(n)})$$

∴ 目标函数 $\sum_{j=1}^J \sum_{i \in I_j} L(y_i, \hat{y}_i^{(n)}) + \frac{1}{2} \lambda \sum_{j=1}^J w_j^2 + y^T$

$$\hat{y}_i^{(n)} = \hat{y}_i^{(n-1)} + w_j$$

$$\begin{matrix} \downarrow & \downarrow \\ x_0 & x - x_0 \end{matrix} = \lambda$$

$$\therefore L(y_i, \hat{y}_i^{(n)}) \approx \underbrace{L(y_i, \hat{y}_i^{(n-1)})}_x + \underbrace{L'(y_i, \hat{y}_i^{(n-1)}) w_j}_{\lambda} + \frac{1}{2} \underbrace{L''(y_i, \hat{y}_i^{(n-1)}) w_j^2}_{\lambda^2}$$

$$\therefore obj^{(n)} = y^T + \sum_{j=1}^J \frac{1}{2} \lambda w_j^2 + \sum_{j=1}^J \sum_{i \in I_j} L'(y_i, \hat{y}_i^{(n-1)}) w_j + \frac{1}{2} L''(y_i, \hat{y}_i^{(n-1)}) w_j^2$$

$$= y^T + \sum_{j=1}^J \left[w_j \cdot \underbrace{\sum_{i \in I_j} L'(y_i, \hat{y}_i^{(n-1)})}_{G_j} + \frac{1}{2} w_j^2 \cdot \underbrace{\sum_{i \in I_j} L''(y_i, \hat{y}_i^{(n-1)})}_{H_j} + \frac{1}{2} \lambda w_j^2 \right]$$

$$= y^T + \sum_{j=1}^J w_j + G_j + \frac{1}{2} w_j^2 (\lambda + H_j)$$



$$\therefore (w_1, w_2, \dots, w_j) = \min \text{obj}^{(n)}$$

$$w_1^* = \arg \min (w_1 G_1 + \frac{1}{2} w_1^2 H_1)$$

$$w_2^* = \arg \min (w_2 G_2 + \frac{1}{2} w_2^2 H_2)$$

$$\therefore w_j^* = - \frac{G_j}{H_j + \lambda}$$

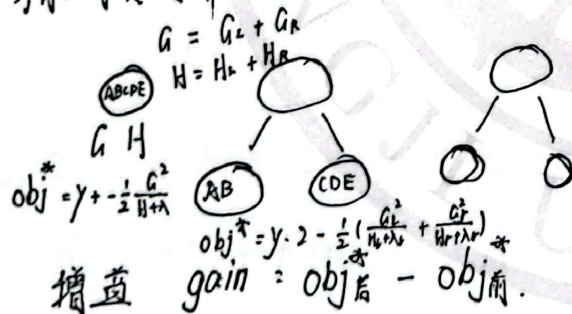
代入目标函数.

$$\text{obj}^{(n)} = \gamma T + (1 - \frac{1}{2}) \sum_{j=1}^T \frac{G_j^2}{H_j + \lambda}$$

确定树的结构.

obj^*

精确贪心算法



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

1. if $\text{gain} < 0$, 或 $\text{len}(lj) = 1$,
stop.
2. or if $\text{gain} < \text{threshold}$,
stop.
3. max-depth or T .



同濟大學

TONGJI UNIVERSITY
SHANGHAI

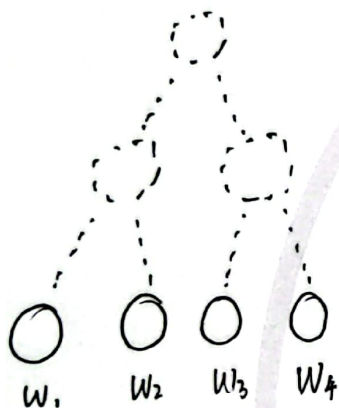
PEOPLE'S REPUBLIC OF CHINA

XGBoost 梯度提升回归树

集成算法，加法模型。

由许多棵回归树共同进行预测。

1. 单棵回归树。



2. 模型表达

$$y_i^{(T)} = \underbrace{\sum_{t=1}^T f_t(x_i)}_{y_i^{(T-1)}} + f_T(x_i) = \sum_{t=1}^T f_t(x_i)$$

T 棵树预测值之和为最终

3. 前向分布算法：逐个优化，优化第 t 棵树时，假定前 t-1 棵已

$$y_i^{(t)} = y_i^{(t-1)} + f_t(x_i)$$

4. 目标函数

$$\sum_{i=1}^N L(y_i, y_i^{(t)}) + \sum_{t=1}^T \lambda(f_t) = obj^{(t)}$$



同 濟 大 學

TONGJI UNIVERSITY
SHANGHAI
PEOPLE'S REPUBLIC OF CHINA

XGBoost.

回归树. CART树. 二叉.

使用方式 sklearn 中 API, 或者 xgboost

目前先用 sklearn 中的 API.

XGBoost 三大版块.

参数 集成算法 弱评估器 其他过程

集成算法: 多个弱评估器

装袋法: bagging. 每棵树之间 (弱评估器) 互不干扰, 最后求均值.

提升法: boosting. 逐次迭代. $1 \rightarrow 2 \rightarrow 3$. 基于前面的树的结果生成新的树, 不断有很多棵这样的树

梯度提升回归树. GBDT

单颗树中每个叶子节点都包含一些样本;

待预测的样本在这颗树中最后会落到某一个叶子节点上, 该样本的预测值即为该叶子节点的样本均值.

输入 x_i , 在第 k 棵树上的结果表示为 $h_k(x_i)$, γ_k 表示树在森林中的权重
则 $\hat{y}_i = \sum_k \gamma_k h_k(x_i)$, $h_k(x_i)$ 为该叶子上所有样本的均值.

但 XGB 作为 GBDT 改进, 有所不同.

XGB 中每颗树上每个叶子节点上的值不再是均值, 而是以更复杂计算后的回归取值. 称为叶子权重. 记为 f_k . \therefore 输入 x_i , 第 k 棵树得
共计 K 棵树. $\therefore \sum_k f_k(x_i) = \hat{y}_i^{(k)}$



同濟大學

TONGJI UNIVERSITY

SHANGHAI

PEOPLE'S REPUBLIC OF CHINA

∴ 在 XGB 中, $\hat{y}_i^{(k)} = \sum_k f_k(x_i)$.

k 表示集成算法中生成树的数目, 是一个超参数.

sk: n_estimators 默认 = 100 表示评估器数量. ①
xgb: num-round 默认 10

sk: silent default: True 是否打印每次训练结果 ②
xgb: silent de: false 可监控建模流程.

W_j 第 j 个节点的值预测

I_j 第 j 个节点包含的样本的集合.

x_i 第 i 个样本

$q(x_i)$ 第 i 个样本落入的节点的序号

$W_{q(x_i)}$ 第 i 个样本的预测值.

J 回归树的总数

$y_i^{(T)}$ 模型对第 i 个样本的预测值

$y_i^{(t)}$ 包含 t 棵树的模型对第 i 个样本的预测值

$f_t(x_i)$ 第 t 棵树对第 i 个样本的预测值.

N 样本总数

$\sum (y_i, y_i^{(t)})$ 包含 t 棵树的模型对 i 样本的预测值与真实值误差.

$$\sum (y_i, y_i^{(t)}) = (y_i - y_i^{(t)})^2.$$

$\gamma(t)$ 第 t 棵树的复杂度