

$GBDT \Rightarrow$ Gradient Boost Decision Tree. (梯度提升树)

Boost : many weak predicted model ensemble into one Strong model.

Gradient Descent + Boosting \Rightarrow 弱学习机

Boost Tree

Scale:

Boosting $\xrightarrow{\text{Specify}}$ Gradient Boosting $\xrightarrow{\text{Specify}}$ GBDT

Difference between GBDT and Random Forest: 集成模型

① Drop single decision tree.
In Random Forest, each decision tree is trained independently and it doesn't affect the performance of the model if we randomly / selectively drop one tree.

In GBDT, a sequence of decision tree each improve the decision of all future trees. And if we drop one tree, the performance of the model will be affected (the first tree affect more than the last tree).

Δ AdaBoost 算法 \Rightarrow 最具代表性的提升树算法

② Combine two models:

RandomForest can be constructed in parallel
while GBDT model can only be constructed in
sequence manner.

③ When the number of decision tree increase.

For RandomForest, it doesn't affect a lot.
Because the trees don't have relationship among
them.

For GBDT, it does matter. The model
will be more powerful if the tree ↑.

提升的思想： F_{m+1} 通过提升值纠正 F_m 的误差变得更完美。

GBDT思想：提升树由多棵基础决策树构成，每一棵决策树均依赖于前一棵决策树分类结果而设置不同的权重（预测样本错误 \Rightarrow 增大样本权重；正确 \Rightarrow 降低样本权重。）每次构造新的提升树均是对原本的决策树的优化。

Adaboost Algorithm:

$$F_m(x) = \sum_{m=1}^M \alpha_m \cdot f_m(x)$$

$$F_m(x) = F_{m-1}(x) + h_m(x) = F_{m-1}(x) + \underbrace{\underbrace{\alpha_m}_{\text{学习率}} \underbrace{f_m(x)}_{\text{权重}}}_{\text{估计量}}$$

第 m 棵基础决策树
黑板以为 |

△ Why Gradient Boost \Leftrightarrow Gradient Descent (以平方损失函数为例)

$$\frac{\text{估计值}}{\text{提升值}} = \frac{y - F_m(x)}{h_{m+1}(x)} = \frac{\underbrace{\frac{1}{2} \cdot (y - F_m(x))^2}_{\text{MSE}}}{\alpha_m}$$

即损失函数的导数，因为梯度下降为 $- \alpha_m \frac{\partial J(\theta)}{\partial \theta}$

而梯度提升为 $F_m(x) + \alpha_{m+1} \cdot h_{m+1}(x)$ ，因为提升其实
是 MSE 的负梯度。

因此 GBDT 模型的估计值/提升值为损失函数的负梯度值，即间接地求得了损失函数梯度值。（当我们损失函数为一般函数，梯度值不好求时，原本回归等机器学习模型不再适用，因此引入 GBDT。）

样本点权重 w_m

确定的量 \swarrow 基础决策树 $f(x)$ 选择

决策树权重 α_m .

How:

Cost function:

$$\begin{aligned} L(y, F(x)) &= e^{-yF(x)} \quad \text{提升函数} \\ &= e^{(-y \sum_{m=1}^M \alpha_m f_m(x))} \\ &= e^{-y(F_{m-1}(x) + \alpha_m f_m(x))} \end{aligned}$$

确定 $\alpha_m - f_m(x)$:

$$\begin{aligned} (\alpha_m, f_m(x)) &= \underset{\alpha, f(x)}{\operatorname{argmin}} \sum_{i=1}^n e^{-y_i(F_{m-1}(x_i) + \alpha_m \cdot f_m(x_i))} \quad \text{所有样本数} \\ &= \underset{\alpha, f(x)}{\operatorname{argmin}} \sum_{i=1}^n e^{-y_i \cdot \alpha_m \cdot f_m(x_i)} \cdot W_{mi} \quad \text{已知, 与 } \alpha_m - f_m(x) \text{ 无关} \end{aligned}$$

$$= \begin{cases} \underset{\alpha, f(x)}{\operatorname{argmin}} \sum_{i=1}^n e^{-\alpha_m} \cdot W_{mi}, & \text{第 } m \text{ 棵棵树预测准确} \\ \underset{\alpha, f(x)}{\operatorname{argmin}} \sum_{i=1}^n e^{\alpha_m} \cdot W_{mi}, & \text{第 } m \text{ 棵棵树预测错误} \end{cases} \Rightarrow y_i \cdot f_m(x_i) = 1 \\ \Rightarrow y_i \cdot f_m(x_i) = -1 \end{math>$$

此时仅与 α_m 有关, 因此得先求最佳 $f_m(x)$ 使 cost function

最小.

$$f_m^*(x) = \underset{\text{从所有决策树空间 } \mathcal{T} \text{ 中}}{\arg\min} \sum_{i=1}^n w_{mi} \cdot I(y_i \neq f_m(x))$$

$\stackrel{y_i \neq f_m(x)}{=} 1$
 $\stackrel{y_i = f_m(x)}{=} 0$
 $y_i = f_m(x)$

求出 $f_m^*(x)$ 后，上式：

$$L(y, F(x)) = \sum_{\substack{y_i = f_m(x_i)}} w_{mi} \cdot e^{-d_m} + \sum_{\substack{y_i \neq f_m(x_i)}} w_{mi} \cdot e^{d_m}$$

第 m 棵棵树预测准确

 $y_i \cdot f_m(x_i) = 1$

第 m 棵棵树预测错误

 $\Rightarrow y_i \cdot f_m(x_i) = -1$

$$= \sum_{i=1}^n w_{mi} \cdot e^{-d_m} + (e^{d_m} - e^{-d_m}) \sum_{\substack{y_i \neq f_m(x_i)}} w_{mi}$$

$y_i \neq f_m(x_i)$

$$\sum_{i=1}^n w_{mi} \cdot I(y_i \neq f_m(x_i))$$

$$= \sum_{i=1}^n w_{mi} \cdot e^{-\alpha_m} + (e^{\alpha_m} - e^{-\alpha_m}) \cdot \sum_{i=1}^n w_{mi} \cdot I(y_i \neq f_m(x_i))$$

对 $L(y, F(x))$ 求偏导：



$$\frac{\partial L(y, F(x))}{\partial \alpha_m} = 0 \quad \Rightarrow \quad \alpha_m^* = \frac{1}{2} \log \left(\frac{1-e_m}{e_m} \right)$$

e_m : 错误率

$$e_m = \sum_{i=1}^n w_{mi} \cdot I(y_i \neq f_m(x_i))$$

因此，我们可以得到：

$$\left\{ \begin{array}{l} f_m^*(x) = \operatorname{argmin} \sum_{i=1}^n w_{mi} \cdot I(y_i \neq f_m(x_i)) \\ e_m = \sum_{i=1}^n w_{mi} \cdot I(y_i \neq f_m(x_i)) \\ \alpha_m^* = \frac{1}{2} \cdot \log \left(\frac{1-e_m}{e_m} \right) \end{array} \right.$$

\Rightarrow 根据 $f_m(x)$ 的预测结果，得到下一轮样本点

权重:

$$W_{m,i}^* = \left\{ \begin{array}{l} \frac{w_{mi} \cdot e^{(-dm^*)}}{\sum_{i=1}^n w_{mi} \cdot e^{(-dm^*)}}, f_m(x_i)^* = y_i \\ \frac{w_{mi} \cdot e^{dm^*}}{\sum_{i=1}^n w_{mi} \cdot e^{dm^*}}, f_m(x_i)^* \neq y_i \end{array} \right.$$

GBT Algorithm:

$$F_m(x) = F_{M-1}(x) + h_m(x)$$

$$= F_{M-1}(x) + \alpha^m \cdot f_m(x)$$

$$= F_{M-1}(x) + \sum_{j=1}^J c_{mj} \cdot I(x_j \in R_{mj})$$

m 层提升树

$$= \sum_{m=1}^M$$

$$\sum_{j=1}^J$$

$f_m(x)$ 在样本 j 上预测值
结点

J 样本个数

上采样

SMOTE Algorithm \Rightarrow 合成少数过采样技术 (Synthetic Minority Over-sampling Technique)

过采样
欠采样
组合采样

Steps:

I. Choose minority sample a .

Use KNN to get k nearest neighbour.

II. According to sampling rate N ,

choose $B = N \cdot \frac{\#a}{\text{num of } a}$ from k sample

III. For each b in B , create new sample

$$c = a + |a-b| * \text{rand}(0, 1), \text{ for } b \in B.$$

Do until we balance the data.

In python:

随机数种子

SMOTE (ratio='auto', random_state=None, k_neighbours=5,

all \Rightarrow 过采样

近邻,即 k

not minority \Rightarrow 少数
majority \Rightarrow 对多数的样本采样

m -neighbours = 10)

即 B, 近邻中再选取个数.

GBDT

Improve \Rightarrow

XGBoost

Advantage:

XGBoost是由传统的GBDT模型发展而来的，GBDT模型在求解最优化问题时应用了一阶导技术，而XGBoost则使用损失函数的一阶和二阶导，而且可以自定义损失函数，只要损失函数可一阶和二阶求导。

XGBoost算法相比于GBDT算法还有其他优点，例如支持并行计算，大大提高算法的运行效率；XGBoost在损失函数中加入了正则项，^③用来控制模型的复杂度，进而可以防止模型的过拟合；^④ XGBoost除了支持CART基础模型，还支持线性基础模型；^⑤ XGBoost采用了随机森林的思想，对字段进行抽样，既可以防止过拟合，也可以降低模型的计算量。^⑥