

XGBoost 算法推导

XGBoost 整体上仍属于 GBDT 算法系统,也是由多个基模型组成的加性模型

所以 XGBoost 可表示为: $\hat{y}_i = \sum_{k=1}^K f_k(x_i)$

前向分步算法第 t 次迭代的基模型为 $f_t(x)$, 有:

$$\hat{y}_i^{(t)} = \sum_{k=1}^t \hat{y}_i^{(k-1)} + f_t(x_i)$$

XGBoost 损失函数基本形式由经验损失项和正则化项组成:

$$L = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_{i=1}^t \Omega(f_i) \quad \dots \text{式(1)}$$

$\sum_{i=1}^n l(y_i, \hat{y}_i)$ 为经验损失项,表示训练数据预测值与真实值之间的损失; $\sum_{i=1}^t \Omega(f_i)$ 为正则化项表示全部 t 棵树的复杂度之和,这也是 XGBoost 控制模型过拟合的方法。

根据前向分步算法,以第 t 步模型为例,假设模型对第 i 个样本 x_i 的预测值为:

$$\hat{y}_i^{(t)} = \hat{y}_i^{(t-1)} + f_t(x_i)$$

式(1)的目标函数可以改写为:

$$\begin{aligned} L^{(t)} &= \sum_{i=1}^n l(y_i, \hat{y}_i^{(t)}) + \sum_{i=1}^t \Omega(f_i) \\ &= \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \sum_{i=1}^t \Omega(f_i) \end{aligned}$$

因为前 $t-1$ 棵树的结树已经确定,所以前 $t-1$ 棵树的复杂度之和也可以表示为常数

$$\sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t) + \text{Constant} \quad \dots \text{式(2)}$$

针对式(2)的前半部分使用二阶泰勒展开式,这里需要用到函数的二阶导数,相应的损失函数经验损失项可以改写为:

$$l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) = l(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)$$

g_i 为损失函数一阶导数, h_i 为损失函数二阶导数。注意这里是对 $\hat{y}_i^{(t-1)}$ 求导。

将上式代入式(2)中得到损失函数的近似表达式:

$$L^{(t)} \approx \sum_{i=1}^n \left[l(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \Omega(f_t) + \text{Constant}$$

去掉相关常数项,得到简化后的损失函数表达式为:

$$L^{(t)} \approx \sum_{i=1}^n \left[g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \Omega(f_t)$$

由上式知,只需要求出损失函数每一步的一阶导数和二阶导数值,并对目标函数进行优化求解,就可以得到前向分步中每一步的模型 $f(x)$,最后根据加性得到XGBoost.

为了计算XGBoost决策树对结点分裂条件,进一步进行推导:

假设一棵决策树是由叶子结点的权重 w 和样本实例到叶子结点的映射关系 q 构成,映射关系可以理解为决策树的分支结构。

$$f_t(x) = wq(x)$$

定义决策树复杂度的正则化项。模型复杂度 Ω 可由单棵决策树的叶子结点数 T 和叶子权重 w 决定,所以模型的复杂度可以表现为:

$$\Omega(f_t) = rT + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

XGBoost的损失函数改写为

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

将属于第 j 个叶子结点的所有样本 x_i 划入一个叶子结点的样本集合中,因而:

$$L^{(t)} = \sum_{i=1}^T \left[\left(\sum_{i \in L_j} g_i \right) w_j + \frac{1}{2} \left(\sum_{i \in L_j} h_i + \lambda \right) w_j^2 \right] + rT$$

定义 $G_j = \sum_{i \in L_j} g_i$, $H_j = \sum_{i \in L_j} h_i$, 其中 G_j 可以理解为叶子结点 j 所包含样本的一阶偏导数累加之和, H_j 可以理解为叶子结点 j 所包含样本的二阶偏导数累加之和,都为常量。

将 G_j 和 H_j 代入上式,损失函数变换为 $L^{(t)} = \sum_{j=1}^T [G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2] + rT$

又对每个叶子结点 j ,将其从目标函数中单独取出: $G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2$

上式是一个只包含一个变量叶子结点权重 w_j 的一元二次函数,可根据最值公式求其最值点。当相互独立的每棵树的叶子结点都达到最优值时,整个损失函数也相应地达到最优。在树结构固定的情况下,对上式求导并令其为0,可得到最优点和最优值为:

$$w_j^* = -\frac{G_j}{H_j + \lambda} \quad L = -\frac{1}{2} \sum_{j=1}^T \frac{G_j^2}{H_j + \lambda} + rT$$

假设决策树模型在某个结点进行了分裂,分裂前的损失函数写为:

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

分裂后的信息增益为:

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

如果增益 $\text{Gain} > 0$, 即分裂为两个叶子结点后, 目标函数下降了, 则考虑此次结果。

XGBoost 的核心就是通过损失函数开展到二阶导数来进一步逼近真实损失。

当定义或者选择 XGBoost 损失函数时, 需要其二阶可导, 以平方损失为例:

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

对其的一阶导数和二阶导数分别为:

$$g_i = \frac{\partial l(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}} = -2(y_i - \hat{y}_i^{(t-1)})$$

$$h_i = \frac{\partial^2 l(y_i, \hat{y}_i^{(t-1)})}{\partial (\hat{y}_i^{(t-1)})^2} = 2$$