

模型表达: $\hat{y}_i^{(b)} = \sum_{b=1}^{B-1} f_b(x_i) + f_B(x_i)$

$$\text{obj}^{(b)} = \sum_{i=1}^N \ell(y_i, \hat{y}_i^{(b)}) + \sum_{j=1}^b \Omega(f_j)$$

$\Omega(f_j)$ is regularization term

$$\Omega(f_j) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T C_j^2 \quad (\gamma, \lambda \text{ are hyperparams})$$

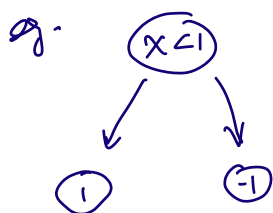
↑
of DT 节点

$\sum_{j=1}^{b-1} \Omega(f_j)$ is known by the time we train tree b .

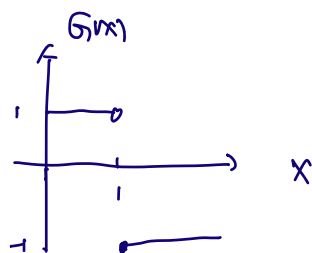
$$\sum_{j=1}^b \Omega(f_j) \Rightarrow \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T C_j^2$$

优化时, 正则项只与当前 tree 的 T 和 C_j 有关

梯度下降不适合树模型



$$G(x) = \begin{cases} 1 & x < 1 \\ -1 & x \geq 1 \end{cases}$$



λ continuous, can't take derivative

$$\text{def. } L(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2$$

$$\text{obj}^{(b)} = \sum_{i=1}^N L(y_i, \hat{y}_i^{(b)}) + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T C_j^2$$

$$\Rightarrow \text{obj}^{(b)} = \gamma T + \sum_{j=1}^T \left[\sum_{x_i \in R_j} L(y_i, \hat{y}_i^{(b-1)} + C_j) \right] + \frac{1}{2} \lambda \sum_{j=1}^T C_j^2$$

T : # of leaf nodes.

C_j : value assigned to j th leaf node

we don't know exact form of L , so how can we min L step by step? \Rightarrow approx L

$$\begin{cases} \text{GBDT} : -\eta \eta \text{ Taylor} \\ \text{XGBoost} : +\eta \eta \text{ Taylor} \end{cases}$$

$$f(x) = f(x_0) + (x - x_0) f'(x_0) + \frac{1}{2} (x - x_0)^2 f''(x_0)$$

For individual sample i :

$$\underbrace{L(y_i, \hat{y}_i^{(t)})}_{\text{常量}} = L(y_i, \underbrace{\hat{y}_i^{(t-1)}}_{\text{equiv to 'x'}} + \underbrace{C_j}_{\text{常量}})$$

第 t 颗树对应当前样本 i 的叶节点取值

$$\approx \underbrace{L(y_i, \hat{y}_i^{(t-1)})}_{\text{常量}} + C_j \cdot L'(y_i, \hat{y}_i) \Big|_{\hat{y}_i = \hat{y}_i^{(t-1)}} + \frac{1}{2} C_j^2 L''(y_i, \hat{y}_i) \Big|_{\hat{y}_i^{(t-1)}}$$

$$\Rightarrow \min_{\underbrace{C_j}_{g_i := -\eta \eta / g^{(t-1)}}} \underbrace{L'(y_i, \hat{y}_i) \Big|_{\hat{y}_i = \hat{y}_i^{(t-1)}}}_{g_i} + \frac{1}{2} C_j^2 \underbrace{L''(y_i, \hat{y}_i) \Big|_{\hat{y}_i^{(t-1)}}}_{h_i := \eta \eta / g^{(t-1)}}$$

$$= \min C_j g_i + \frac{1}{2} C_j^2 h_i$$

Aggregate: $\min_{j=1}^T \sum_{x_i \in R_j} \frac{1}{2} (y_i - \hat{y}_i^{(t-1)} + C_j)^2 + \delta T + \frac{1}{2} \lambda \sum C_j^2$

$\approx \min_{j=1}^T \left[C_j \sum_{x_i \in R_j} y_i + \frac{1}{2} C_j^2 \sum_{x_i \in R_j} h_i + \frac{1}{2} \lambda C_j^2 \right] + \delta T$

$\Rightarrow \text{obj}^{(b)} = \sum_{j=1}^T \left[C_j \underbrace{\sum_{x_i \in R_j} y_i}_{G_j} + \frac{1}{2} C_j^2 \left(\underbrace{\sum_{x_i \in R_j} h_i}_{H_j} + \lambda \right) \right] + \delta T$

G_j, h_i can be computed directly H_j j : j -th leaf node

$(C_1^* \dots C_T^*) = \arg \min_{C_j} \text{obj}^{(b)}$ * Each leaf node can be computed individually to min L , and then \sum results

\downarrow

$\forall j \in [T] \quad C_j^* = \arg \min (C_j G_j + \frac{1}{2} C_j^2 (H_j + \lambda)) + \delta T$

$(\because \text{当 tree grow 时, 每个 } (x_i, y_i) \text{ 落入哪个 } R \text{ is determined by } x_i, \text{ deterministic})$

$C_j^* = \arg \min (C_j G_j + \frac{1}{2} C_j^2 (H_j + \lambda)) + \delta T$

\uparrow H is \sum of 2nd derivative of L fn, which is convex. $\Rightarrow H > 0$

So this $-y_i = \hat{y}_i$ to min is convex too

$$\begin{cases} C_j^* = -\frac{b}{2a} = -\frac{G_j}{H_j + \lambda} \\ \text{obj}_j^*(b) = \delta T - \frac{1}{2} \frac{G_j^2}{H_j + \lambda} \end{cases}$$

树的划分:

以上推导都建立在树已构造好的基础上, 那树如何构造呢?

① Exact Greedy Algo

eg.

$A B C D E$

$G \quad H$

$$w^* = - \frac{G}{H+\gamma}$$

$$obj_{before}^* = \gamma \cdot 1 - \frac{1}{2} \frac{G^2}{H+\gamma}$$

known & deterministic before split

$X \geq 10$

$A B$

G_L
 H_L

$C D E$

G_R
 H_R

$$obj_{after}^* = \gamma \cdot 2 - \frac{1}{2} \left(\frac{G_L^2}{H_L+\gamma} + \frac{G_R^2}{H_R+\gamma} \right)$$

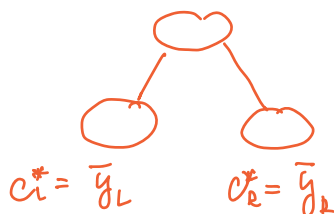
for diff split, will have diff obj_{after}^*

$$gain = obj_{before}^* - obj_{after}^* \quad (\text{loss reduction after split})$$

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

use split that max gain

Relation to CART



$$L_{before} = \left(\sum_i (y_i - \bar{y})^2 \right) / N = \text{Stdev}_{total}^2 \quad \text{known & deterministic before split.}$$

$$L_{after} = \left(\sum_{x_i \in R_L} (y_i - \bar{y}_L)^2 + \sum_{x_i \in R_R} (y_i - \bar{y}_R)^2 \right) / (N_L + N_R)$$

$$= \text{Stdev}_L^2 \cdot P_L + \text{Stdev}_R^2 \cdot P_R$$

- when to stop:
- ① $\max \text{gain} \leq 10e-05$
 - ② 叶子节点包含样本个数 ≤ 1
 - ③ depth / # of leaf nodes

Algo: Exact Greedy Search for Split Finding.

Algorithm 1: Exact Greedy Algorithm for Split Finding

Input: I , instance set of current node

Input: d , feature dimension

$\text{gain} \leftarrow 0$

$G \leftarrow \sum_{i \in I} g_i, H \leftarrow \sum_{i \in I} h_i$

for $k = 1$ to d do

$G_L \leftarrow 0, H_L \leftarrow 0$

for j in sorted(I , by x_{jk}) do

$G_L \leftarrow G_L + g_j, H_L \leftarrow H_L + h_j$

$G_R \leftarrow G - G_L, H_R \leftarrow H - H_L$

$\text{score} \leftarrow \max(\text{score}, \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{G^2}{H + \lambda})$

end

end

Output: Split with max score

pro: exact

con: computation cost

time complexity: $d \cdot n \log n$

init: 对 feature / feature取值进行筛选

$\left\{ \begin{array}{l} \text{compromise 精度} \\ \text{boost computation} \end{array} \right. \Rightarrow \text{近似 algo}$

② 近似 algo

a. feature selection

1. 随机 rand. (split 之前就 rand 出 subset of features)

2. 按层 rand (每次 split 都 rand 新的 subset of feature)

b. feature value selection

1. 分桶

理想状态: assume all samples uniform distribute (出现 prob. 一样)

eq.	X_1	Sample
	1	B D
bin 1	2	A E
	3	C
bin 2	4	H L
	5	I
	6	G
bin 3	7	M N
	8	T

样本: 12

桶数: 3

avg: 4 samples / bin

result: split on $X_1=2$ and $X_1=5$

2. 加权方法

$$\begin{aligned}
 & \sum_{i=1}^N (y_i f_{\theta}(x_i) + \frac{1}{2} h_i f_{\theta}^2(x_i)) + \Omega(f_{\theta}) \\
 &= \sum_{i=1}^N \frac{1}{2} h_i \left(\frac{y_i^2}{h_i^2} + 2 \frac{y_i}{h_i} f_{\theta}(x_i) + f_{\theta}^2(x_i) \right) + \Omega(f_{\theta}) \\
 &= \sum_{i=1}^N \frac{1}{2} h_i \left[f_{\theta}(x_i) - \left(-\frac{y_i}{h_i} \right) \right]^2 + \Omega(f_{\theta})
 \end{aligned}$$

常量, 不影响目标优化

so 目标函数为真实值是 $-y_i/h_i$, 权重为 h_i 的平方损失, 因此, 使用 = 价
梯度加权

eg.

feature x_1	Sample	h_i
<u>1</u>	B	2
	D	1
<u>2</u>	A	2
	E	2
<u>3</u>	C	<u>1</u> 3
<u>4</u>	H	4
	L	1
<u>5</u>	I	3
<u>6</u>	G	5
	M	2
<u>7</u>	N	<u>1</u> 7
8	T	2

eg. $\epsilon = \frac{9}{27}$

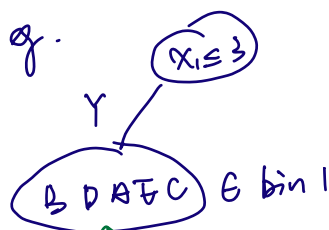
$\sum h_i = 27$

hyperparam: $\epsilon = \frac{n}{27}$ put as many samples in bin s.t. $\frac{\sum \text{weights}}{\sum h_i} \leq \epsilon$

加权分位法 分为

- 全局: 对于每个 tree 的每个 feature, 只进行一次分位
- 局部: 每次 split 都对选中 feature 进行分位.

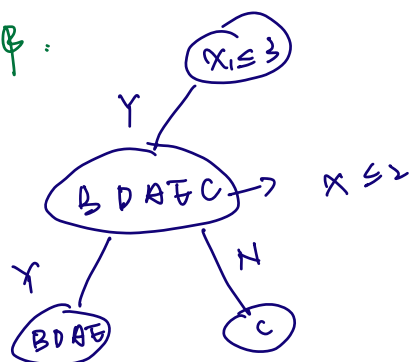
• 全局 above eg.



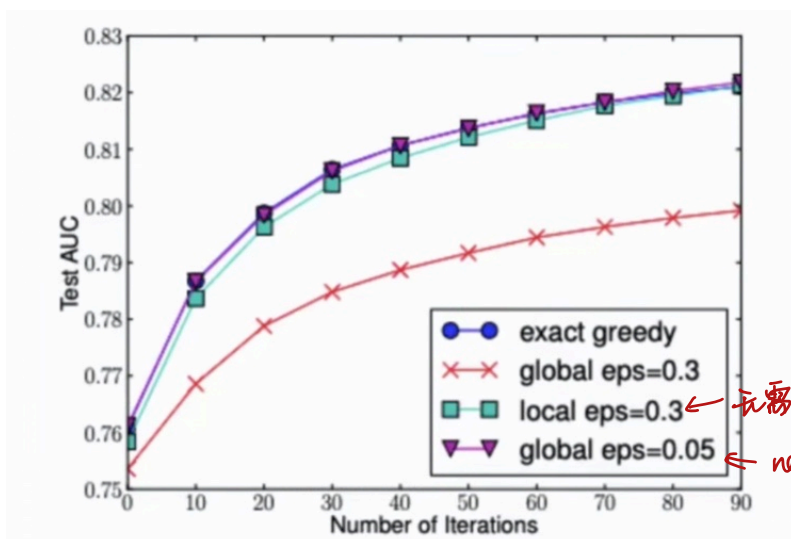
此时无法再对 x_1 进行分位, $\therefore 3, 0, 7$ 已固定

不够 flexible, 往往需要 ϵ 更小 (更多 bins) to perform well

局部:



more flexible. 每次不需过多 bins



需要太多 bins

need ϵ small to perform well

· 缺失值处理.

eg,

feature x_i	Sample
<u>1</u>	<u>B</u>
<u>2</u>	<u>D</u>
<u>3</u>	<u>A</u>
<u>4</u>	<u>E</u>
<u>5</u>	<u>C</u>
<u>6</u>	<u>H</u>
<u>7</u>	<u>L</u>
<u>8</u>	<u>I</u>
<u>9</u>	<u>G</u>
<u>10</u>	<u>M</u>
<u>11</u>	<u>N</u>
<u>12</u>	<u>J</u>
<u>13</u>	<u>K</u>
<u>14</u>	<u>F</u>
<u>15</u>	<u>O</u>
<u>16</u>	<u>P</u>
<u>17</u>	<u>Q</u>
<u>18</u>	<u>R</u>
<u>19</u>	<u>S</u>
Nan	

可简单, 但 $O(N^2)$

论文采用方法: 将 NM 全体放到

左 or 右, 比较 GainL and GainR.

if GainL 大, 则所有 NM 放 left.

缺失值不参与排序和分位

· 学习率 Shrinkage

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

可防 overfitting, 通常 $\eta = 0.1$