

XGBoost (Extreme Gradient Boosting)

Additive model $f(x)$ is base learner, estimation is \hat{y}

$$\hat{y} = \sum_{m=1}^M f_m(x), \quad f_m \in \mathcal{F} \quad \text{①}$$

The target function

$$L^{(m)} = \sum_{i=1}^N L(y_i, \hat{y}_i^{(m-1)} + f_m(x_i)) + Q_r(f_m) \quad \text{②}$$

where Q_r is the regularization term and GBDT doesn't have.

$$Q_r(f) = \gamma J + \frac{1}{2} \lambda \sum_{j=1}^J b_j^2$$

where J is the number of leaf nodes and b_j is the value
 Q_r aims at restricting the complexity of tree, which means decrease the number of leaves and values on the leaves.
 b_j larger, learn faster \Rightarrow overfitting.

Second-order taylor expansion of $L^{(m)}$ on $\hat{y}_i^{(m-1)}$

$$L^{(m)} \approx \sum_{i=1}^N [L(y_i, \hat{y}_i^{(m-1)}) - g_i f_m(x_i) + \frac{1}{2} h_i f_m^2(x_i)] + (J + \frac{1}{2} \lambda \sum_{j=1}^J b_j^2) \quad \text{③}$$

$$\text{where } g_i = \frac{\partial L(y_i, \hat{y}_i^{(m-1)})}{\partial \hat{y}_i^{(m-1)}}, \quad h_i = \frac{\partial^2 L(y_i, \hat{y}_i^{(m-1)})}{(\partial \hat{y}_i^{(m-1)})^2}$$

DT splits feature space into different isolated area and each sample belongs to one of the regions. Single DT:

$$f(x) = \sum_{j=1}^J b_j I(x \in R_j). \quad \text{If we sum all samples up, } \sum_{i=1}^N f(x_i) = \sum_{j=1}^J \sum_{x \in R_j} b_j$$

Plug it to ③ and remove the constant term $L(y_i, \hat{y}_i^{(m-1)})$

$$L^{(m)} = \sum_{j=1}^N \left[g_j f_m(x_j) + \frac{1}{2} h_j f_m^2(x_j) \right] + \gamma J + \frac{1}{2} \lambda \sum_{j=1}^J b_j^2$$

$$= \sum_{j=1}^J \left[\sum_{x_i \in R_j} g_j b_j + \frac{1}{2} \left(\sum_{x_i \in R_j} h_j + \lambda \right) b_j^2 \right] + \gamma J$$

$$= \sum_{j=1}^J (G_j b_j + \frac{1}{2} (H_j + \lambda) b_j^2) + \gamma J \quad \textcircled{4}$$

where $G_j = \sum_{x_i \in R_j} g_i$, $H_j = \sum_{x_i \in R_j} h_i$

New gain calculation method

If the structure of a tree is confirmed, compute the deriv of ④ and optimal value:

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

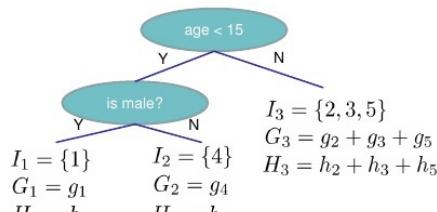
Plug it to ④:

$$L^{(m)} = -\frac{1}{2} \sum_{j=1}^J \frac{G_j^2}{H_j + \lambda} + \gamma J \quad \textcircled{5}$$

It can be thought as score and match different loss func.

(can be derivative
in first-order
and second-order)

Instance index	gradient statistics
1	g_1, h_1
2	g_2, h_2
3	g_3, h_3
4	g_4, h_4
5	g_5, h_5



$$Obj = -\sum_j \frac{G_j^2}{H_j + \lambda} + 3\gamma$$

The smaller the score is, the better the structure is

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

Gain larger \Rightarrow split - the tree.

Pseudocode :

① initialize $f_0(x)$

② for $m=1$ to M :

(a') compute 1st deriv and 2nd deriv of loss func over each sample:

$$f_i = \frac{\partial L(y_i, \hat{y}_i^{(m-1)})}{\partial \hat{y}_i^{(m-1)}}, h_i = \frac{\partial^2 L(y_i, \hat{y}_i^{(m-1)})}{(\partial \hat{y}_i^{(m-1)})^2}$$

(b) recursively use **Algorithm 2** to generate a DT for $f_m(x)$

(c) add DT to models: $\hat{y}_i^{(m)} = \hat{y}_i^{(m-1)} + f_m(x)$

Algorithm 2: Exact Greedy Algorithm for Split Finding

Input: I , instance set of current node

Input: d , feature dimension

$\text{Gain} = 0$

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

For $k = 1$ to m **do:**

$$G_L = 0, H_L = 0$$

For j in $\text{sorted}(I, \text{by } x_{j,k})$ **do:**

$$G_L = G_L + g_j, H_L = H_L + h_j$$

$$G_R = G - G_L, H_R = H - H_L$$

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

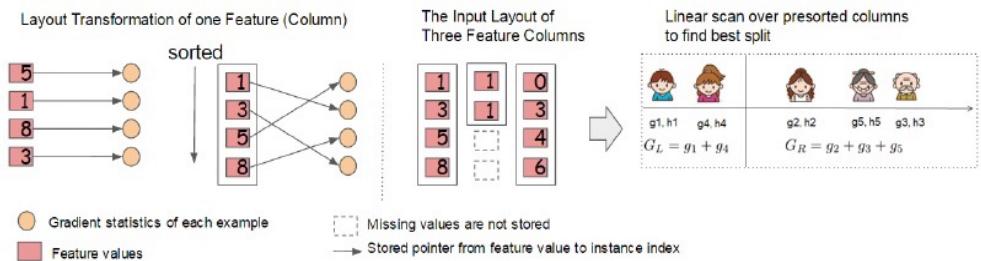
End

End

Output: Split with max score

Column Block for Parallel Learning.

It's hard to split the continuous value for DT. because we have many values. So we can sort the samples by feature values and sequentially compute the Gain to choose the split point. saved as block and each cat stored in ascending order.



To improve training speed.

Algorithm 2: Approximate Algorithm for Split Finding

```
for k = 1 to m do
    | Propose  $S_k = \{s_{k1}, s_{k2}, \dots, s_{kl}\}$  by percentiles on feature k.
    | Proposal can be done per tree (global), or per split(local).
end
for k = 1 to m do
    |  $G_{kv} \leftarrow= \sum_{j \in \{j | s_{k,v} \geq x_{jk} > s_{k,v-1}\}} g_j$ 
    |  $H_{kv} \leftarrow= \sum_{j \in \{j | s_{k,v} \geq x_{jk} > s_{k,v-1}\}} h_j$ 
end
Follow same step as in previous section to find max score only among proposed splits.
```

