

Boosting

① 概念:

*问题: 随机森林中 m 棵决策树独立采样建立, 互相之间没有关系.

*期望: ① 能否把前 $m-1$ 棵树的信通用到第 m 棵树的建立上

② 可否在建立 m 棵决策树时, 得到投票时各棵树的权重

*思路: ① 在生成 m 个弱预测模型时, 每一棵都依据损失函数的梯度方向 (Gradient Boosting) \Rightarrow XGBoost 的思想,

② 直观理解: 一个样本第一次分对, 第二次分错, 第三次这个分错的样本权重被调高。如此反复, 总能收敛到一个效果不错的分类器 \Rightarrow AdaBoost 的思路

③ LR 的梯度: $\theta^* = \theta - \alpha \nabla J(\theta)$

Boost 的梯度: 给决策树做一个目标函数, 沿着偏导方向下降, 得到另一棵决策树.

② Boosting 算法的目标.

*目标: 对训练样本 $(\vec{x}_1, y_1), (\vec{x}_2, y_2), (\vec{x}_3, y_3), \dots, (\vec{x}_n, y_n)$ 目标是找到近似函数 $F(\vec{x})$, 使得损失函数 $L(y, F(\vec{x}))$ 的损失值最小.

$L(y, F(\vec{x}))$ 的典型定义 $\begin{cases} L(y, F(\vec{x})) = \frac{1}{2} (y - F(\vec{x}))^2 \leftarrow \text{最小二乘} \\ L(y, F(\vec{x})) = |y - F(\vec{x})| \end{cases}$

用公式来表达该目标, 就是

$$F^*(\vec{x}) = \arg \min_F E_{(x,y)} [L(y, F(\vec{x}))]$$

*表示最优 F 找到使得 L 最小的 F 使得. 损失函数在 (x,y) 值域上的数学期望

假设误差服从高斯分布 (见线性回归)

其中 $L(y, F(\vec{x})) = \frac{1}{2} (y - F(\vec{x}))^2$ 或 $|y - F(\vec{x})|$

$$F(\vec{x}) = \sum_{i=1}^M \alpha_i f_i(\vec{x}) + \text{const}$$

弱分类器 弱分类器预估值函数 权重.

③ XGBoost Boosting 算法框架: 用梯度提升的方法找到最优解 $F_0(x), F_1(x), \dots, F_n(x)$ 使得损失函数在训练集上的数学期望最小

* 起始步骤: 给定常函数 $F_0(x) = C$, 求 $F_0(\vec{x}) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$

\downarrow 找到最优的 γ \downarrow 损失 \downarrow 样本 y_i 的预估值, 其实就是 C^* 常数 γ 就是 C

这样的函数总能找到, 以 $L(y, F(x)) = |y - F(x)|$ 为例, C^* 就是样本 y 值的中位数.

↓ 解法: 不妨假设样本 y_1, y_2, \dots, y_n 递增排序, $x_k \leq \gamma, x_{k+1} > \gamma$

$$J(\gamma) = \sum_{i=1}^n |x_i - \gamma| = \sum_{i=1}^k (\gamma - x_i) + \sum_{i=k+1}^n (x_i - \gamma)$$

优化目标

$$\frac{\partial J(\gamma)}{\partial \gamma} = \sum_{i=1}^k (1) + \sum_{i=k+1}^n (-1) \xrightarrow{\text{定义}} 0$$

若要让偏导为 0, 前 k 个样本数目与后 $n-k$ 个样本数目必须为 0, 即 γ 为中位数

* 迭代步骤: 给定前 $t-1$ 个弱分类器 (或回归器) 生成的 $f_{t-1}(\vec{x})$, 推导出 $f_t(\vec{x})$

$$J(f_t) = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t) + C$$

优化目标 损失函数 分类器 $0 \sim t-1$ 的预测结果 分类器 t 的预测结果 弱分类器 t 的罚项 (防止过拟合) 公式推导时析出的常数项

把损失函数看作增量形式, 用泰勒展开式展开, 忽略三阶无穷小

$$\textcircled{1} f(x + \Delta x) \approx f(x) + f'(x)\Delta x + \frac{1}{2}f''(x)\Delta x^2$$

由 $\textcircled{1}, \textcircled{2}$ 得到

$$\textcircled{2} g_i \stackrel{\text{def}}{=} \frac{\partial L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}} \quad \textcircled{3} h_i \stackrel{\text{def}}{=} \frac{\partial^2 L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)^2}}$$

$\hat{y}_i^{(t-1)}$ $\hat{y}_i^{(t-1)}$ 定义 g_i, h_i

$$\textcircled{3} J(f_t) \approx \sum_{i=1}^n \left[L(y_i, \hat{y}_i^{(t-1)}) + g_i \cdot f_t(\vec{x}_i) + \frac{1}{2} h_i f_t^2(\vec{x}_i) \right] + \Omega(f_t) + C$$

相当于 $f(x)$ 相当于 $f'(x)\Delta x$ 相当于 $\frac{1}{2}f''(x)\Delta x^2$

求解式子 $\textcircled{3}$ 即可推导出 $f_t(\vec{x})$

先解决 $f_t(\vec{x}_i)$, $\Omega(f_t)$ 求值的问题

* 求解 $f_t(\vec{x}_i)$: 假定决策树 t 的叶子节点数为 T , 每个叶子节点权值为 $\vec{w} = (w_1, w_2, \dots, w_T)$, 样本 \vec{x}_i 落在叶节点 $q(x)$ 中, 定义

$$f_t(\vec{x}_i) = w_{q(x)} \text{ 即节点 } q(x) \text{ 的权值}$$

* 求解 $\Omega(f_t)$: 有很多种定义, 一种方法是考虑 叶节点数 和 叶权值平方和 正则化系数

事先给定的超参数

$$\Omega(f_t) = \gamma \cdot |T_t| + \lambda \cdot \frac{1}{2} \sum_{j=1}^T w_j^2$$

叶子个数 叶节点权值的平方和

不希望过拟合
不希望树太深 (4~8层为宜)

不希望叶节点太多
太多平方和降用叶节点

控制模型复杂度

* 计算式子 (3):

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

目标函数 前 $t-1$ 棵树的预测值, 与 $f_t(\vec{x}_i)$ 无关, 扔到常数项 C 中

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

用上两步的推导结论求解

$$= \sum_{i=1}^n \left[g_i w_{q(x_i)} + \frac{1}{2} h_i w_{q(x_i)}^2 \right] + \left[\gamma \cdot |T_t| + \lambda \cdot \frac{1}{2} \sum_{j=1}^T w_j^2 \right] + C$$

↓ 加法结合率: 由按样本分组改为按叶节点分组

$$= \sum_{j=1}^T \left[\left[\sum_{i \in I_j} g_i \right] w_j + \frac{1}{2} \left[\sum_{i \in I_j} h_i \right] w_j^2 \right] + \gamma \cdot |T_t| + \lambda \cdot \frac{1}{2} \sum_{j=1}^T w_j^2 + C$$

$$= \sum_{j=1}^T \left[\left[\sum_{i \in I_j} g_i \right] w_j + \frac{1}{2} \left[\sum_{i \in I_j} h_i + \lambda \right] w_j^2 \right] + \gamma \cdot |T_t| + C$$

↓ 命名为

↓ 命名为

$$= \sum_{j=1}^T \left([G_j] w_j + \frac{1}{2} [H_j] w_j^2 \right) + \gamma \cdot |T_t| + C$$

第 t 棵决策树叶节点 j 的权重

到这一步, 优化目标已经可以由第 t 棵决策树叶节点的权重 w_j ($j=1, 2, \dots, T$) 表示

* 求解最优值. 对优化目标 (式子③) 关于 $\vec{W} = (w_1, w_2, \dots, w_m)$ 求偏导

$$\frac{\partial J(f_t)}{\partial w_j} = G_j + (H_j + \lambda) w_j \stackrel{\text{令}}{=} 0 \Rightarrow w_j = -\frac{G_j}{H_j + \lambda}$$

从而得到第 i 棵树各个叶子节点的权重 (第 $i-1$ 棵树 T_{i-1} 的叶子节点权重也是这么来的, 当 $i=0$ 时 T_0 是常函数不存在权重问题)

* 把求解得到的 w_j 代入到目标函数 $J(f_t)$ 中

$$\begin{aligned} J(f_t) &= \sum_{j=1}^m (G_j w_j + \frac{1}{2} H_j w_j^2) + \gamma |T_t| + C \\ &= \sum_{j=1}^m \left(-\frac{G_j^2}{H_j + \lambda} + \frac{1}{2} \cdot \frac{H_j G_j^2}{(H_j + \lambda)^2} \right) + \gamma |T_t| + C \\ &= -\frac{1}{2} \sum_{j=1}^m \frac{G_j^2}{H_j + \lambda} + \gamma |T_t| \end{aligned}$$

即为最小损失值公式

* 能够计算损失值, 那么给定分割点 ϕ (把样本分成 L, R 两部分), 就能计算信息增益.

$$\begin{aligned} \underset{\text{分割点}}{\text{Gain}(\phi)} &= J(f_t) - [J_R(f_t) + J_L(f_t)] \\ &= \frac{1}{2} \left[-\frac{(G_L + G_R)^2}{H_L + H_R + \lambda} - \left[-\frac{G_L^2}{H_L + \lambda} - \frac{G_R^2}{H_R + \lambda} \right] \right] - \gamma \\ &= \frac{1}{2} \left[\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma \end{aligned}$$

枚举所有可分割点, 增益大于阈值, 选增益最大的点进行分割,

增益小于阈值, 不再分割; 或者遇到纯节点也不再分割

* 注意, 在得到式子③时用的 $f_t(\vec{x}_i) = w_{g(x)}$ 是假设我们知道 $g(x)$ 和 $w_{g(x)}$ 的值, 而现在我们已经能得到 $g(\vec{x})$ 和 $w_{g(x)}$ 的真实值, 因此可用于

下一轮迭代 n

$$J(f_{t+1}) = \sum_{i=1}^n \mathcal{L}(y_i, \hat{y}_i^{(t)} + f_{t+1}(\vec{x}_i)) + \Omega(f_{t+1}) + C \quad \text{以推导 } f_{t+1}(\vec{x}_i)$$

* 模型调参

- <a>正则项 $\Omega(f_i)$: 对复杂模型增加惩罚项, 如模型复杂度正比于叶节点数目或叶节点预测值的平方和等
- 叶节点数目控制了树的层数, 一般选择 $4 \leq J \leq 8$
- <c>叶节点包含最少样本数目限制: 防止过小叶节点, 降低预测精度
- <d>梯度提升迭代次数 M : 增加 M 可降低训练集损失, 但有过拟合风险

* 其它训练技巧

- <a>衰减 Shrinkage $F_m(\vec{x}) = F_{m-1}(\vec{x}) + \nu \cdot \rho_m f_m(\vec{x}_i)$ $0 < \nu \leq 1$

ν 称为学习率, $\nu=1$ 即为前面推导的原始模型, 推荐 $\nu < 0.1$ 的学习率, 但太小也会导出迭代次数增加。

- 随机梯度提升: Stochastic Gradient Boosting

每次迭代对伪残差样本采用无放回的概率采样, 用部分样本来训练基函数的参数。令训练样本数占所有伪残差样本的比例为 f : $f=1$ 即为原始模型; 推荐 $0.5 \leq f \leq 0.8$

⊗ 较小的 f 能够增强随机性, 防止过拟合, 并且收敛的快。

降采样的好处是能够使用剩余样本做模型验证。

* 例子: 见PPT.

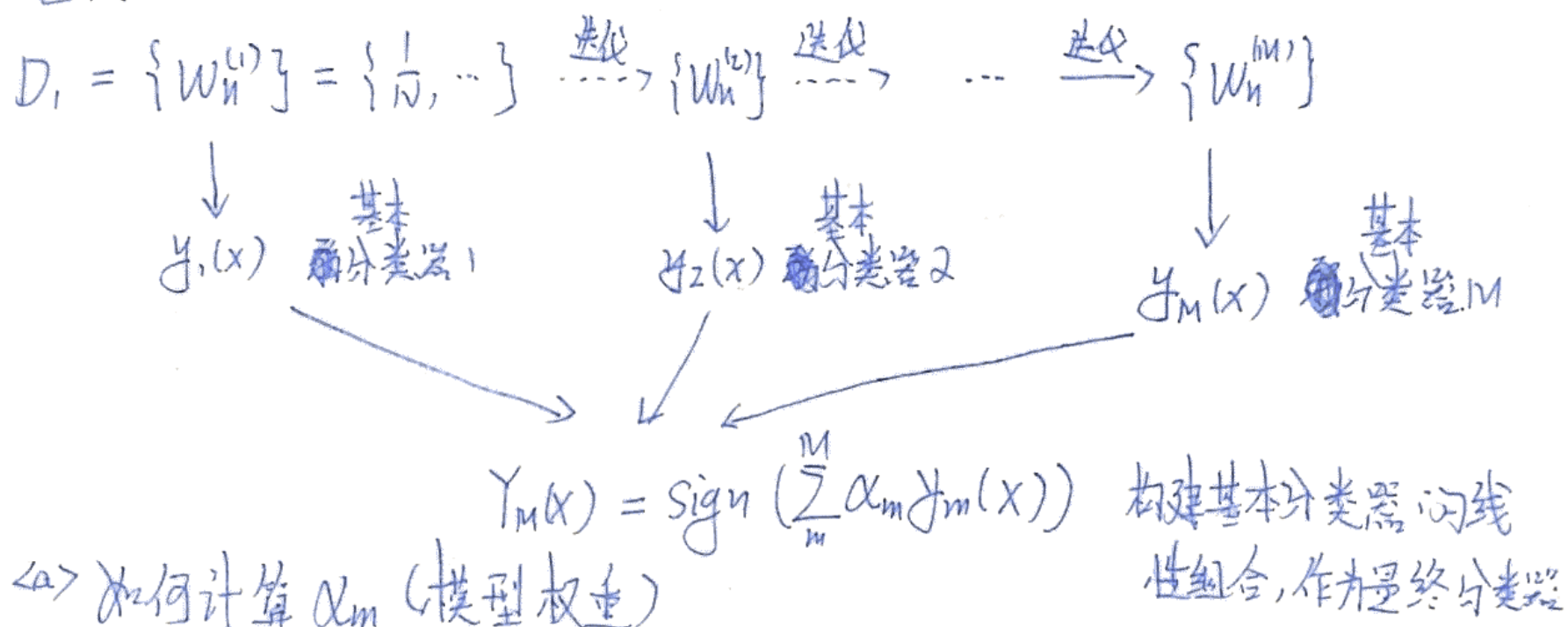
* 小结: 以上为 XGBoost 使用的核心推导过程, 相比传统的 GBDT, XGBoost 使用了二阶导信息, 可以更快地在训练集上收敛

XGBoost (属于“随机森林族”) 因此也具有抗过拟合的优势
| 使用了并行/多核计算, 训练速度更快,

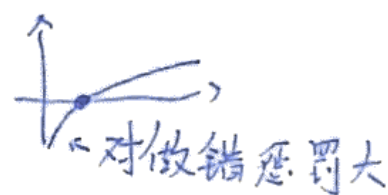
④ AdaBoost

* 思想: 对于 N 个样本, 给每个样本一个权重, 初始权重都是 $1/N$
 第 i 轮迭代中: 如果分错, 该权重会被调高; 如果分对, 该权重会被调低
 第 $i+1$ 轮迭代中, 之前分错的样本会得到更多重视
 每轮迭代都会生成一个弱分类器及权重, 最终组合成一个强分类器

* 算法框架:



$$\alpha_m = \frac{1}{2} \log \frac{1 - e_m}{e_m} \quad \begin{array}{l} \text{做对概率} \\ \text{做错概率} \end{array}$$



误差率 $e_m = P(G_m(\vec{X}_i) \neq y_i) = \sum_{i=1}^N \frac{W_{m,i}}{N} I(G_m(\vec{X}_i) \neq y_i)$

$\frac{W_{m,i}}{N}$ 样本权重 $I(\text{test}) = \begin{cases} 1 & \text{test 为 true} \\ 0 & \text{test 为 false} \end{cases}$

$G_m(\vec{X}_i)$ 弱分类器预测值: -1 或 +1

 如何更新 $D_m = \{W_n^{(m)}\}$ (样本权重)

$$W_{m+1,i} = \frac{W_{m,i}}{Z_m} \cdot e^{-\alpha_m \cdot y_i \cdot G_m(\vec{X}_i)}$$

\downarrow \downarrow \downarrow \downarrow \downarrow

基本分类器序号 样本 基本分类器权重 取与 取值 取与 取值 取与 取值

规范化因子, 为了使 $W_{m+1,i}$ 成为一个概率值.

$$\text{其中 } Z_m = \sum_{i=1}^N W_{m,i} e^{-\alpha_m y_i G_m(\vec{X}_i)}$$

* 例子: 见PPT: 模型实例推导

* 例子: 见PPT: AdaBoost 代码

* AdaBoost 误差上限: AdaBoost 是 Adaptive Boost 缩写, 含义是只要基础分类器误差 < 0.5 , 我可以让 AdaBoost 误差下降, 其理论依据就是 AdaBoost 误差上限

$$\frac{1}{N} \sum_{i=1}^N I(G(x_i) \neq y_i) \leq \frac{1}{N} \sum_i e^{-y_i \cdot f(x_i)} = \prod_m Z_m = \prod_m \left(\sum_{j=1}^N w_{mj} e^{-\alpha_m y_j G_{mj}} \right)$$

*证明左半部分 ② : $G(x_i) \neq y_i$ 时, $y_i f(x_i) < 0$.

*证明右半部分 (b) :

$$\begin{aligned}
 \frac{1}{N} \sum_i e^{-y_i \cdot f(\vec{x}_i)} &= \sum_i \frac{1}{N} e^{-\sum_{m=1}^M \alpha_m \cdot y_i \cdot G_m(\vec{x}_i)} & \leftarrow f(\vec{x}_i) = -\sum_{m=1}^M \alpha_m G_m(\vec{x}_i) \\
 &= \sum_i w_{1i} \cdot e^{-\sum_{m=1}^M \alpha_m \cdot y_i \cdot G_m(\vec{x}_i)} = \sum_i w_{1i} \prod_{m=1}^M e^{-\alpha_m \cdot y_i \cdot G_m(\vec{x}_i)} & \leftarrow w_{11} = w_{12} = \dots = w_{1n} = \frac{1}{N} \text{ 初始值} \\
 &= \sum_i w_{1i} \cdot e^{-\alpha_1 \cdot y_i \cdot G_1(\vec{x}_i)} \prod_{m=2}^M e^{-\alpha_m y_i \cdot G_m(\vec{x}_i)} \\
 &= \sum_i z_1 \cdot w_{2i} \prod_{m=2}^M e^{-\alpha_m y_i \cdot G_m(\vec{x}_i)} & \leftarrow z_1 = \frac{w_{1i}}{z_1} \cdot e^{-\alpha_1 y_i G_1(\vec{x}_i)} \\
 &= z_1 \sum_i w_{2i} \prod_{m=2}^M e^{-\alpha_m y_i \cdot G_m(\vec{x}_i)} & \leftarrow w_{2i} = \frac{w_{2i}}{z_2} e^{-\alpha_2 y_i G_2(\vec{x}_i)} \\
 &= z_1 z_2 \sum_i w_{3i} \prod_{m=3}^M e^{-\alpha_m y_i \cdot G_m(\vec{x}_i)} & \vdots \\
 &= z_1 z_2 \dots z_{m-1} \sum_i w_{mi} e^{-\alpha_m y_i \cdot G_m(\vec{x}_i)} & \vdots \\
 &= z_1 z_2 \dots z_{m-1} z_m = \prod_{m=1}^M z_m & \leftarrow w_{mi} = \frac{w_{mi}}{z_m} e^{-\alpha_m y_i G_m(\vec{x}_i)}
 \end{aligned}$$

* 训练误差上界

$$\prod_{m=1}^M Z_m = \prod_{m=1}^M [2\sqrt{e_m(1-e_m)}] \stackrel{\text{令 } \gamma_m = \frac{1}{2} - e_m}{=} \prod_{m=1}^M \sqrt{1-4\gamma_m^2} \leq e^{-2\sum_{m=1}^M \gamma_m^2} \leq e^{-\frac{1}{2}}$$

$$Z_m = \sum_{i=1}^m w_{mi} e^{-\alpha_m \cdot \gamma_i G_m(\vec{x}_i)}$$

$$= \sum_{j_i = G_m(\vec{x}_i)}^{j=1} w_{m_i} e^{-\alpha_m} + \sum_{j_i \neq G_m(\vec{x}_i)} w_{m_i} e^{\alpha_m}$$

$$= (1 - e_m)e^{-\alpha_m} + e_m e^{\alpha_m} = 2\sqrt{e_m(1 - e_m)}$$

$$-2 \sum_{m=1}^{\infty} \gamma_m^2$$

C_m = 分錯次數

推導見附子

*附录二：证明 $\prod_{m=1}^M \sqrt{1-4\gamma_m^2} \leq e^{-2\sum_{m=1}^M \gamma_m^2}$

$$e^{-x} + x - 1 \geq 0 \Rightarrow e^{-x} \geq 1 - x$$

代入 $x = 4\gamma^2$ 得

$$e^{-4\gamma^2} \geq 1 - 4\gamma^2 \Rightarrow e^{-2\gamma^2} \geq \sqrt{1-4\gamma^2}$$

$$\Rightarrow \begin{cases} e^{-2\gamma_1^2} \geq \sqrt{1-4\gamma_1^2} \\ e^{-2\gamma_2^2} \geq \sqrt{1-4\gamma_2^2} \\ \vdots \\ e^{-2\gamma_m^2} \geq \sqrt{1-4\gamma_m^2} \end{cases} \Rightarrow \prod_{m=1}^M e^{-2\gamma_m^2} \geq \prod_{m=1}^M \sqrt{1-4\gamma_m^2}$$

$$\Rightarrow e^{-2\sum_{m=1}^M \gamma_m^2} \geq \prod_{m=1}^M \sqrt{1-4\gamma_m^2}$$

*前面的 $\frac{1}{N} \sum_i e^{-y_i \cdot f(\vec{x}_i)}$ 看着比较突兀，其实可以把它看作是损失函数取指数函数，认为模型服从若干基础分类器加权，然后用心法（前向分步算法）来推导，发现 AdaBoost 其实正符合这个思路（事后推导）

（推导过程见附录1）

⑤ Bagging 与 Boosting

↓
增强训练集精度 减少偏差 Bias
降低 Variance 增加泛化能力

⑥ GBDT：使用关于分类器的一阶导数进行学习（与 XGBoost 不同）

XGBoost 是 GBDT (GBM) 的一个 C++ 实现

⑦ 数据预处理 — 清洗

*特征：重要性高，缺失率低

- ① 通过计算填充
- ② 通过经验或业务知识估计

*特征：重要性高，缺失率高

- ① 从其它渠道补全
- ② 用其它字段计算或获取
- ③ 去除字段并在结果中体现

低 特征：重要性低，缺失率低

- ① 不做处理或简单处理

+ 特征：不重要，缺失率高

- 去除该字段

低