

详细版本：

https://www.julyedu.com/question/big/kp_id/23/ques_id/2595

给定一个训练数据集 $T = \{(x_1, y_1), (x_2, y_2) \cdots (x_N, y_N)\}$ ，其中实例 $x \in \mathcal{X}$ ，而实例空间 $\mathcal{X} \subset \mathbb{R}^n$

y_i 属于标记集合 $\{-1, +1\}$ 。Adaboost 的目的就是从训练数据中学习一系列弱分类器或基本分类器，然后将这些弱分类器组合成一个强分类器。

Adaboost 的算法流程如下：

步骤1：首先，初始化训练数据的权值分布。每一个训练样本最开始时都被赋予相同的权值： $1/N$

$$D_1 = (w_{11}, w_{12} \cdots w_{1i} \cdots, w_{1N}), \quad w_{1i} = \frac{1}{N}, \quad i = 1, 2, \cdots, N$$

步骤2：进行多轮迭代，用 $m=1, 2, \dots, M$ 表示迭代的多少轮

1. 使用具有权值分布 D_m 的训练数据集学习，得到基本分类器（选取让误差率最低的阈值来设计基本分类器）：

$$G_m(x): \mathcal{X} \rightarrow \{-1, +1\}$$

2. 计算 $G_m(x)$ 在训练数据集上的分类误差率

$$e_m = P(G_m(x_i) \neq y_i) = \sum_{i=1}^N w_{mi} I(G_m(x_i) \neq y_i)$$

由上述式子可知， $G_m(x)$ 在训练数据集上的误差率 e_m 就是被 $G_m(x)$ 误分类样本的权值之和。

3. 计算 $G_m(x)$ 的系数， α_m 表示 $G_m(x)$ 在最终分类器中的重要程度（目的：得到基本分类器在最终分类器中所占的权重）：

$$\alpha_m = \frac{1}{2} \log \frac{1 - e_m}{e_m}$$

由上述式子可知， $e_m \leq 1/2$ 时， $\alpha_m \geq 0$ ，且 α_m 随着 e_m 的减小而增大，意味着分类误差率越小的基本分类器在最终分类器中的作用越大。

4. 更新训练数据集的权值分布（目的：得到样本新的权值分布），用于下一轮迭代

$$D_{m+1} = (w_{m+1,1}, w_{m+1,2}, \dots, w_{m+1,i}, \dots, w_{m+1,N})$$

$$w_{m+1,i} = \frac{w_{mi}}{Z_m} \exp(-\alpha_m y_i G_m(x_i)), \quad i = 1, 2, \dots, N$$

使得被基本分类器 $G_m(x)$ 误分类样本的权值增大，而被正确分类样本的权值减小。就这样，通过这样的方式，Adaboost方法能“重点关注”或“聚焦于”那些较难分的样本上。

其中， Z_m 是规范化因子，使得 D_{m+1} 成为一个概率分布：

$$Z_m = \sum_{i=1}^N w_{mi} \exp(-\alpha_m y_i G_m(x_i))$$

步骤3：组合各个弱分类器

$$f(x) = \sum_{m=1}^M \alpha_m G_m(x)$$

从而得到最终分类器，如下：

$$G(x) = \text{sign}(f(x)) = \text{sign}\left(\sum_{m=1}^M \alpha_m G_m(x)\right)$$