

集成学习

集成学习和方差的模型组合效果很好

输入 x 和输出 y 满足 $y = h(x)$

对于 M 个不同模型 $f_i(x)$

$$R(f_m) = E_x[(f_m(x) - h(x))^2] \\ = E_x[\epsilon_m(x)^2]$$

$\epsilon_m(x) = f_m(x) - h(x)$ 是模型 m 在样本 m 上的误差

$$\bar{R}(f) = \frac{1}{M} \sum_{m=1}^M E_x[\epsilon_m(x)^2]$$

通过某种策略集成多个模型，群策群力来提高准确率。

{ 直接平均
加权平均

投票: $F(x) = \frac{1}{M} \sum_{m=1}^M f_m(x)$

$$R(F) = E_x[(\frac{1}{M} \sum_{m=1}^M f_m(x) - h(x))^2] \\ = E_x[(\frac{1}{M} \sum_{m=1}^M \epsilon_m(x))^2]$$

$$= \frac{1}{M^2} E_x[\sum_{m=1}^M \sum_{n=1}^M \epsilon_m(x) \epsilon_n(x)]$$

$$= \frac{1}{M^2} \sum_{m=1}^M \sum_{n=1}^M E_x[\epsilon_m(x) \epsilon_n(x)]$$

其中 $E_x[\epsilon_m(x) \epsilon_n(x)]$ 是两模型误差的相关性

不相关 $E_x[\epsilon_m(x) \epsilon_n(x)] = 0 \quad \forall m \neq n$

相同 $E_x[\epsilon_m(x) \epsilon_n(x)] = E_x[\epsilon_m(x)^2] \quad \forall m = n$

且 $\epsilon_m \geq 0 \quad \forall m$ 得

$$\bar{R}(f) \geq R(F) \geq \frac{1}{M} \bar{R}(f)$$

要求每个模型有差异性，随着模型数量增多，误差率下降。→ 尽可能不

Bagging类

随机抽取训练样本，随机选择特征
提高每个模型的独立性

① Bagging (Bootstrap Aggregating):

原始训练集上进行有放回的采样(随机)
得到 M 个彼此独立的训练集训练 M 个模型
然后投票

② 随机森林

引入随机特征。
基模型是决策树

Boosting类

先后训练不同的基模型，每个模型针对前序模型的错误进行专门训练。

根据前序模型的结果，来调整训练样本的权重。从而增加差异性

① AdaBoost.

加权模型

$$F(x) = \sum_{m=1}^M \underbrace{\alpha_m}_{\text{强分类}} \underbrace{f_m(x)}_{\text{弱分类}}$$

学习第 m 个弱分类器后，增加其分类样本权重，使得第 $m+1$ 个弱分类器更专注于分类错误样本。

线性化模型

$$L(F) = \exp(-yF(x)) \\ = \exp(-y \sum_{m=1}^M \alpha_m f_m(x))$$

$$y, f_m(x) \in \{-1, +1\}$$

假设经过 $m-1$ 次迭代

$$F_{m-1}(x) = \sum_{t=1}^{m-1} \alpha_t f_t(x)$$

第 m 次迭代目标找 $\alpha_m, f_m(x)$ 使得

$$L(\alpha_m, f_m(x)) = \sum_{n=1}^N \exp(-y^n (F_{m-1}(x^n) + \alpha_m f_m(x^n)))$$

$$\text{令 } w_m^n = \exp(-y^n F_{m-1}(x^n))$$

则

$$L(\alpha_m, f_m(x)) = \sum_{n=1}^N w_m^n \exp(-y^n \alpha_m f_m(x^n))$$

$$\because y, f_m(x) \in \{-1, +1\},$$

$$\therefore y f_m(x) = 1 - 2I(y \neq f_m(x))$$

将损失函数在 $-\alpha_m y^n f_m(x^n)$ 处展开

$$L(\alpha_m, f_m(x)) = \sum_{n=1}^N w_m^n (1 - \alpha_m y^n f_m(x^n) + \frac{1}{2} \alpha_m^2)$$

$$\alpha \alpha_m \sum_{n=1}^N w_m^n I(y^n \neq f_m(x^n))$$

当 $\alpha_m > 0$ 时, 最优分类是使样本权重为 w_m^n 时加权错误率最小分类

$$L(\alpha_m, f_m(x)) = \sum_{y^n \neq f_m(x)} w_m^n \exp(-\alpha_m) + \sum_{y^n = f_m(x)} w_m^n \exp(\alpha_m)$$

$$\propto (1 - \epsilon_m) \exp(-\alpha_m) + \epsilon_m \exp(\alpha_m)$$

$f_m(x)$ 加权错误率

$$\epsilon_m = \frac{\sum_{y^n \neq f_m(x^n)} w_m^n}{\sum_n w_m^n} \quad (1)$$

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