

5.8 集成学习

1. 个体与集成

弱分类器 (weak learner): 分类准确率只稍好于随机猜测的分类器
将多个学习器组合后产生一个强学习器。具有多样性、鲁棒性

$$H(x) = \text{sign}\left(\sum_{i=1}^T h_i(x)\right)$$

boosting: 个体学习器间存在强依赖关系, 必须串行生成序列化方法

bagging, random forest: 个体学习器间不存在强依赖关系, 可并行生成并行化方法

为什么需要 ensemble learning?

(1) 弱分类器间有一定差异性, 导致误差不同。若多个弱分类器合并后, 可得到更低误差
减少整体误差率。

(2) 对于过大/小数据集, 可随机进行划分和有效回采样产生不同的数据集

然后根据子集训练不同分类器，最后合并成一个分类器

(3) 若也经过多次，可以训练多个模型，再进行模型融合

(4) 对于多个异构特征集，可以进行融合，为每个数据集构建一个分类模型，然后将多个模型融合。

2. Boosting

原理：先从初始训练集训练一个基学习器，再根据表现进行调整，使后续的训练样本多关注，然后基于调整后的样本分布来训练下一个基学习器，直至基学习器数目达到指定值下，将T进行加权融合。

根据损失函数的梯度——梯度提升 (gradient boosting)

Adaboost: 通过错误预测，将权重升高，使迭代过程直至 error rate 小或达到一定的迭代次数。

$$f(x) = \sum_{m=1}^M \alpha_m G_m(x) \quad G_m(x) = \text{sign}(f_m(x)) = \text{sign}\left[\sum_{n=1}^M \alpha_n G_n(x)\right]$$

$$\text{损失函数 loss} = \frac{1}{n} \sum_{i=1}^n I(G(x_i) \neq y_i) \leq \frac{1}{n} \sum_{i=1}^n e^{-y_i f(x_i)}$$

Adaboost过程:

(1) 假设训练数据 $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

(2) 初始化训练数据权重分布 $D_1 = (w_{11}, w_{12}, \dots, w_{1n})$

(3) 用具有权重分布 D_m 的训练数据，基分类器 $G_m(x): X \rightarrow \{-1, +1\}$

估计 $G_m(x)$ 在训练集上的分类误差 $\epsilon_m = P(G_m(x_i) \neq y_i) = \sum_{i=1}^n w_{mi} \cdot I(G_m(x_i) \neq y_i)$

(4) 估计 $G_m(x)$ 模型上的权重 $\alpha_m = \frac{1}{2} \log 2 \left(\frac{1 - \epsilon_m}{\epsilon_m} \right)$

(5) 权重训练数据上的权重分布: $D_{m+1} = (w_{m+1,1}, \dots, w_{m+1,n})$ $w_{m+1,i} = \frac{w_{mi}}{Z_m} e^{-\alpha_m y_i G_m(x_i)}$

(6) Z_m 规范化因子 $Z_m = \sum_{i=1}^n w_{mi} e^{-\alpha_m y_i G_m(x_i)}$

(7) 构建基分类器的线性组合 $f(x) = \sum_{m=1}^M \alpha_m G_m(x)$

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

Adaboost 只能二分类

3. bagging, random forest
bagging: 适用于分类、回归。自助采样 (out-of-bag estimate), 有放回抽样
bagging 降低方差 bias-variance

RF: 加入随机属性选择.

RF步骤: (1) n 样中用 bootstrap 有放回重采样选出 m 个样本.
(2) 使用子数据集训练决策树, 从所有特征中取 k 个, 再选择最佳分割点并以此
节点来递归构造决策树

(3) 重复 m 次, 建立 m 棵决策树 (4) 通过投票表决定属于哪一类

RF优点: 并行化速度快, 随机划分特征对训练数据敏感, 给出各特征重要性
列表, 由于随机抽样, 模型方差小泛化能力强, 实现简单, 对部分特征
缺失不敏感

缺点: 某些数据容易过拟合, 取值比较多的划分特征对 RF 决策影响大

4 结合策略

(1) 平均法 简单, 加权平均法

(2) 投票法 绝对多数 (过半), 相对多数 (得票最多标记, 若同时多个标记
取最高票, 从中随机取一个), 加权.

(3) 学习法: stacking, BMA (Bayes Model Averaging)

5 GBDT

要求的模型必须是 回归 CART 模型. 要求预测值与样本损失反方向, 逐层为回归模型

GBDT: DT (Regression Decision Tree), GB (Gradient Boosting), Shrinkage (衰减)

Loss function: $L(y, F(x)) = \frac{1}{2}(y - F(x))^2$ $L(y, F(x)) = |y - F(x)|$

最优解: $F^*(x) = \arg\min_F L(y, F(x))$

GBDT 回归与分类算法区别

平方损失函数 绝对误差损失函数
 $L(y, F_m(x)) = \frac{1}{2}(y - F_m(x))^2$ $L(y, F_m(x)) = |y - F_m(x)|$ 损失函数
 $\alpha_{im} = y_i - F_{m-1}(x)$ $\alpha_{im} = \text{sign}(y_i - F_{m-1}(x))$ 负梯度
 均值 中值 初始值

分类
 对数损失函数(二分类) 对数损失函数(多分类) 损失函数 $\exp(f_m(x))$
 $L(y, F_m(x)) = -[y \ln(p_m) + (1-y) \ln(1-p_m)]$ $L(y, F_m(x)) = -\sum_{k=1}^K y_k \ln p_k(x)$ $p_k(x) = \frac{\exp(f_k(x))}{\sum_{l=1}^K \exp(f_l(x))}$
 $p_m = \frac{1}{1 + e^{-F_m(x)}}$
 $\alpha_{im} = y_i - p_m$ $\alpha_{im} = y_{ik} - p_{mk}(x)$ 负梯度
 \ln 正样本数/负样本数 0 初始值

GBDT 优点:
 可处理连续/离散值, 少过拟合效果不错 robust
 缺点:
 训练点之间存在关联关系, 对异常值训练, 速度慢

★ bagging, boosting 区别

bagging		boosting
样本选择	有放回随机采样	训练集不变, 对弱模型权重/γ变化 根据误差调整
样本权重	随机为权重	根据错误率不断调整权重, error 权重↑
预测函数	所有权重相等	误差小, 分类权重更大
并行计算	并行生成各基模型	顺序生成, 需上一个结果
bias-var	↓ var 方差	↓ bias 偏差
	需降低过拟合, 泛读	防过拟合, 欠拟合