

<p>希文 23 春 模式识别与机器学习 23 年 6 月 6 日</p> <p>分类器评价指标</p> <p>测试错误率作为真实错误率的可靠估计是有条件的。 最重要的条件是：测试样本与训练样本独立，测试样本与未来样本独立同分布，测试集样本数足够大。此时测试错误率就是对真实错误率的无偏估计，但估计的置信范围取决于测试样本数目，有可能会很大。</p> <p>解决训练样本集与测试样本集矛盾：n-fold 交叉验证；留一法交叉验证 (LOOCV)(交叉验证的极端情况，样本少时被广泛应用)，自举法与 B.632 估计若有病D^*，检测阳性T^+，则对应 FP(假阳性)，检测阴性T^-，对应 FN 假阴性。若无病D^*，检测阳性T^+，则对应 TP(真阳性)，检测阴性T^-，对应 TN 真阴性。</p> <p>混淆矩阵对应关系混淆对</p> <p>Sensitivity 灵敏度: $P(T^+ D^+) = TP / (TP + FN)$ 为 (1-犯第二类错误概率) Specificity 特异度: $P(T^- D^-) = TN / (TN + FP)$ 为 (1-犯第一类错误概率) Prevalence 患病率: 人群中患病的概率 $D^+ / (D^+ + D^-)$ Discovery rate 发现率: $P(D^+ T^+) = TP / (TP + FP)$ False discovery rate 误发现率: $P(D^- T^+) = FP / (TP + FP)$ Accuracy 精度/准确率: $(TP + TN) / (D^+ + D^-)$</p> <p>第一类错误: 假阳性; 第二类错误: 假阴性</p> <p>计算采用贝叶斯公式, sensitivity Sn, Specificity Sp (等于 $1 - fp$)</p> $\frac{P(D^+ T^-)}{P(D^+ T^+)} = \frac{\frac{P(T^-)}{P(T^+)}}{\frac{P(D^+ T^-) + Sp \cdot (1 - P(D^+))}{Sn \cdot P(D^+) + (1 - Sp)(1 - P(D^+))}}$ <p>两次患病: $P(A (B \cap C)) = \frac{P((B \cap C) A) \times P(A)}{P((B \cap C) A) \times P(A) + P((B \cap C) A^c) \times P(A^c)}$</p> <p>Other Jargons: Precision: $\frac{TP}{TP + FP}$; Recall: $\frac{TP}{TP + FN}$; Accuracy: $\frac{TP + TN}{TP + TN + FP + FN}$</p> <p>ROC: Receiver Operating Characteristic. 横坐标为 False Positive Rate, 纵坐标为 True Positive Rate. AUC: Area Under Curve, higher the better, 0.5: random guess.</p>	<p>Joint probability: $p(\mathbf{x}, \mathbf{h}) = e^{-E(\mathbf{x}, \mathbf{h})} / Z$ 其中 $Z = \sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}$</p> <p>MLE: $P(\mathbf{x}; \theta) = \sum_{\mathbf{h}} p_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})} / Z$</p> <p>则 $\frac{\partial \log P(\mathbf{x}; \theta)}{\partial \theta} = \mathbb{E}_{P(\mathbf{x}, \mathbf{h})} \left[\frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} \right] - \mathbb{E}_{P(\mathbf{x}, \mathbf{h})} \left[\frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} \right]$ 其中 $\mathbb{E}_{P(\mathbf{x}, \mathbf{h})}$ 未知, 需要估计.</p> <p>$p(\mathbf{h} \mathbf{x}) = \prod_j \text{sigm} \left((2\hat{z}_j - 1) (b_j + w_j^T \cdot \mathbf{x}) \right)$</p> <p>$p(\mathbf{x} \mathbf{h}) = \prod_k \text{sigm} \left((2x_k - 1) (c_k + hW_k) \right)$</p> <p>1980 年代三种主要类型的神经网络 前馈神经网络 Feedforward NN (代表性方法: 多层感知器; 反馈型神经网络 Feedback NN (代表性方法: Hopfield NN); 竞争学习神经网络 Competitive Learning NN (代表性方法: 自组织映射 Self-organizing map)</p> <p>近邻法</p> <p>分段线性判别函数: 把各类划分为若干子类, 以子类中心作为代表点. 考查新样本与各类代表点的距离, 极端情况: 近邻法.</p> <p>K-近邻法: 找出 x 的前 k 个近邻, 看其中多数属于哪一类, 就把 x 分到哪一类. K 越大, 分界线越平滑.</p> <p>KNN 常见问题: 存储量和计算量: 要数接近时风险较大, 有噪声时风险加大; 样本无穷多时性能优异, 有样本下性能如何? 改进: 减少计算量和存储量; 引入拒绝机制, 根据实际问题修正投票方式.</p> <p>近邻法在计算上的问题: 需存储所有训练样本; 新样本需与每个样本做比较. 快速算法基本思想: 把样本集分成若干子集 (树状结构), 每个子集 (节点) 可用较少几个样本代表; 通过将新样本与各节点比较排除大量候选样本; 只有最后的节点 (子集) 中这个样本比找出最近邻.</p> <p>分枝界定算法 (Branch-Bound Algorithm): 已知样本划分成多个子集, 形成一个树状结构, 每个节点是一个子集, 每个子集用若干个的几量代表; 新样本按顺序与各节点进行比较来排除不可能包含近邻的样本, 最后只在少数节点上与每个样本进行比较</p> <p>剪枝近邻法: 处在两类交界处分布重合区的样本可能误导决策, 应将其从样本集中去掉. 若思考: 考查样本是否可能的误导样本, 若是则从样本集中去掉 (剪枝). 考查方法是: 通过交叉验证, 认为错分样本为误导样本.</p> <p>压缩近邻法: 主要通过减少存储量. 将训练样本分为 X_0 和 X_c 两个集合, 开始时 X_0 中只有一个样本, X_c 中为其余样本. 考查 X_c 中每个样本, 若用 X_0 可正确分类则保留, 否则移入 X_0. 最后用 X_0 作分类的样本集.</p> <p>原型近邻法: 用若干个原型 (prototype) 样本来代表训练数据. 通常原型的数量远少于训练样本数 (如等于样本数则为 1-近邻法).</p> <p>定义原型点: 利用 k-means 聚类. 对于第 1 类样本, 进行聚类分析, 获得 R 个原型 (该类样本密度高的点). 定义这 R 个原型的类别为 i; 总共获得 CR 个原型样本 (C 为类别数).</p> <p>对于待分类的样本 x, 我们将其类别为距离它最近的原型的标签.</p> <p>原型近邻法的优点有: 模型测试时计算开销小; 节省内存. 没有: 模型训练时计算开销小, 可视化效果好; 能更好捕捉数据分布的细节信息.</p> <p>可做拒绝决策的近邻法: 从简单多数变为绝对多数. 拒绝决策同样可引入改进的近邻法中, 比如剪枝近邻法.</p> <p>高维情况下需要更多的训练样本, 当标准维度较高时, “最近邻”可能相差很远. 使用欧氏距离通常会对应每一维特征做标准化.</p> <p>按距离加权的 KNN(LOESS): 距离近的邻居加权大, 距离远的邻居加权小. 一种常用的权值函数 $w^{(i)} = \exp \left[-\frac{(x^{(i)} - x)^2}{\tau^2} \right]$. τ 为带宽. 训练样本少, 带宽应该取得大一些; 模型回归系数 θ^T 与 τ 取值有关.</p>
---	--

线性学习机器

线性判别函数: $y = \text{sgn}(\sum_{i=1}^n w_i x_i + w_0)$

在 X 空间, 类均值向量: $m_1 = \frac{1}{N_1} \sum_{x_i \in x_1} x_i$, $i = 1, 2$ 类内离散度矩阵: $S_1 = \sum_{x_i \in x_1} (x_i - m_1)(x_i - m_1)^T$, $i = 1, 2$. 总类内离散度矩阵 $S_w = S_1 + S_2$.

类间离散度矩阵: $S_b = (m_1 - m_2)(m_1 - m_2)^T$

在 Y 空间, 类均值: $\tilde{m}_1 = \frac{1}{N_1} \sum_{y_i \in y_1} y_i$, $i = 1, 2$. 类内离散度: $\tilde{S}_1 = \sum_{y_i \in y_1} (y_i - \tilde{m}_1)(y_i - \tilde{m}_1)^T$, $i = 1, 2$. 总类内离散度: $\tilde{S}_w = \tilde{S}_1 + \tilde{S}_2$

类间离散度矩阵: $\tilde{S}_b = (\tilde{m}_1 - \tilde{m}_2)(\tilde{m}_1 - \tilde{m}_2)^T$

Fisher 准则: $\max_{\mathbf{w}} J_F(\mathbf{w}) = \frac{(\tilde{m}_1 - \tilde{m}_2)^T \mathbf{w}}{\sqrt{\tilde{S}_w \mathbf{w} \mathbf{w}^T}}$, where $y_1 = \mathbf{w}^T x_1$. 代价函数: $J_F(\mathbf{w}) = \frac{1}{2} \frac{\mathbf{w}^T S_w \mathbf{w}}{\mathbf{w}^T S_b \mathbf{w}}$

由于只与方向有关, 构造 $\mathbf{w}^T S_w \mathbf{w} = 1$, 则将问题转化为 $\min \mathbf{w}^T S_b \mathbf{w}$, s.t. $\mathbf{w}^T S_w \mathbf{w} = 1$. 定义 $L(\mathbf{w}) = -\mathbf{w}^T S_b \mathbf{w} + \lambda(\mathbf{w}^T S_w \mathbf{w} - 1)$, 求导为 0 可得 $S_b \mathbf{w} = \lambda S_w \mathbf{w}$, 即 \mathbf{w} 为矩阵 $S_b^{-1} S_w$ 的特征向量, 则有 $\mathbf{w} = \frac{1}{\lambda} S_b^{-1} S_w \mathbf{w} = \frac{1}{\lambda} S_w^{-1} (S_b - \mu_0)(S_b - \mu_0)^{-1}$, 只考虑方向, 则有 $\mathbf{w}^* \propto S_w^{-1} (m_1 - m_2)$. 还需确定决策的分界点, 如 $w_0 = -\frac{1}{2} (m_1 + m_2)$ 或 $w_0 = -\bar{m}$

Perceptron: $y = \text{sgn}(\sum_{i=1}^n w_i x_i + w_0)$. 若样本线性不可分, 则可以容忍错误, 使错误最小 (比如强制收敛, MSE), 或寻求非线性方法 (比如神经网络, SVM); 若样本线性可分时多解, 则可寻求最优分类器, 如 SVM. 若样本为多类, 可采用多类方法, 或用两类分类器完成多类分类.

线性回归:

$$\min_w E(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N (\hat{y}_j - y_j)^2 = \frac{1}{N} \sum_{j=1}^N (f(\mathbf{x}_j) - y_j)^2 = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

$$= \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

解析解: $VE(\mathbf{w}) = \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = \frac{2}{N} \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{y} = 0$, 有 $\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$

若可逆, 则 $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$. 可逆的条件是特征线性独立. 若特征间不是线性独立, 仍可以计算伪逆, 但解不唯一. 解决方案: 通过特征选择或变换去除冗余, 或通过引入其他准则对解加以限制 (如 SVD 或正则化)

评价指标: $R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$. 在 OLS 模型中, $R^2 = \frac{\sum_{i=1}^N (y_i - \bar{y})^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$. $R^2 = 1$ indicates perfect regression, while $R^2 = 0$ indicates baseline model. 意义: percentage of dependent variable variations that the linear model explains.

Pearson Correlation Coefficient: $r_{XY} = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)} \sqrt{\text{var}(Y)}} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2} \sqrt{\sum (Y_i - \bar{Y})^2}}$

MSE 准则: 求 α 使得 $\alpha^T y_i > 0, i = 1, \dots, N$. 当样本本集线性不可分时, 使尽可能多的样本满足不等式. 为方便求解, 为每个样本引入 b_i , 并令 $\alpha^T y_i = b_i > 0, i = 1, \dots, N$. 这样就可用最小二乘法解决.

优化问题: $\alpha^* = \min_{\alpha} J_S(\alpha)$, where $J_S(\alpha) = \|\mathbf{Y}\alpha - \mathbf{b}\|^2 = \sum_{i=1}^N (\alpha^T y_i - b_i)^2$.

问题: 如何给定 $\mathbf{b} = [b_1, b_2, \dots, b_N]^T$?

可以证明, 如果 \mathbf{b} 选为 $b_i = \frac{1}{N(N/2, \text{if } y_i \in \omega_1)} \frac{1}{N(N/2, \text{if } y_i \in \omega_2)}$, 则 MSE 解等价于 $\mathbf{w}_0 = -\bar{\mathbf{m}}$ 的 Fisher Linear Discrimination 解.

如果 \mathbf{b} 选为 $b_i = 1, i = 1, \dots, N$, 则当 $N \rightarrow \infty$ 时, MSE 解以最小均方误差最优逼近贝叶斯判别函数 $g_0(\mathbf{x}) = P(\omega_1 | \mathbf{x}) - P(\omega_2 | \mathbf{x})$.

使 $e^2 = \int (\alpha^T y - g_0(\mathbf{x}))^2 p(\mathbf{x}) d(\mathbf{x})$ 最小

Logistic Function: $P(y | \mathbf{x}) = \frac{e^{\alpha^T y + \beta x}}{1 + e^{\alpha^T y + \beta x}}$. Odds 几率 $\frac{P(y | \mathbf{x})}{1 - P(y | \mathbf{x})} = e^{\alpha^T y + \beta x}$

$h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$, $\theta(s) = \frac{e^s}{1 + e^s}$

回归目标: 求最大似然函数. 参数为 \mathbf{w} 的模型在数据 $\{(x_1, y_1), \dots, (x_N, y_N)\}$, $x_i \in R^{d+1}$, $y_i \in \{-1, 1\}$ 上的似然函数为 $L(\mathbf{w}) = \prod_{i=1}^N P(y_i | x_i) = \prod_{i=1}^N \theta(y_i \mathbf{w}^T x_i)$

最大似然函数, 等价于最小化 $-\frac{1}{N} \sum_{i=1}^N \ln(1 + e^{-y_i \mathbf{w}^T x_i})$, 也即:

$$\min E(\mathbf{w}) = -\frac{1}{N} \ln L(\mathbf{w}) = -\frac{1}{N} \ln \left(\prod_{i=1}^N \theta(y_i \mathbf{w}^T x_i) \right)$$

梯度下降求解: $VE = -\frac{1}{N} \sum_{i=1}^N \frac{y_i x_i}{1 + e^{y_i \mathbf{w}^T x_i}}$

人工神经网络

Softmax 函数: $P(y = j | \mathbf{x}) = \frac{e^{w_j^T \mathbf{x}}}{\sum_{k=1}^K e^{w_k^T \mathbf{x}}}$. 交叉熵: $J(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{j=1}^K y_j \log \hat{y}_j$

有 $J(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{j=1}^K y_j \log \frac{\exp(w_j^T \mathbf{x})}{\sum_{k=1}^K \exp(w_k^T \mathbf{x})} = \log \sum_{k=1}^K \exp(w_k^T \mathbf{x}) - \sum_{j=1}^K y_j w_j^T \mathbf{x}$

QDA: $g_i(\mathbf{x}) = k_i^2 - (\mathbf{x} - \hat{\mathbf{m}}_i)^T \sum_{i=1}^K (\mathbf{x} - \hat{\mathbf{m}}_i)$, $i = 1, \dots, C$.

decision: $\text{class}(\mathbf{x}) = \text{argmax}_i g_i(\mathbf{x})$

BP 算法训练中可能遇到的问题: 收敛慢, 学习曲线震荡, 过拟合/推广能力差. 可能的原因: 网络结构或激活函数问题, 训练样本不充分或不合适, 初值不合适, 学习率不合适, 缺乏恰当的预处理 (比如归一化), 需要更好的训练策略, 运气问题. **策略:** 激活函数: Sigmoid, tanh, ReLU; 特征尺度调整或归一化 (有效值范围, 相对重要性等); 目标调整: 引入“伪样本”增大训练样本数. 如加噪声或者数据增强; 隐层节点数, 层数, 网络剪枝.

引入动量(Momentum), Weight Decay, 引入提示输出(Hints), 终止条件.

Hopfield Network: An array of binary threshold units fully connected with symmetric weights. Each binary “configuration” of the network has an energy. The dynamic procedure converges to a “memory”. (Associative memory) How does Hopfield Net capture relations? Hebb-rule training: finding associations among pixels.

Boltzmann Machines consist of Visible Nodes and Hidden Nodes, without the output layer. 设 Visible nodes 中的 Offsets 为 c, Hidden Nodes 中的 Offsets 为 b, x 之间的权重为 U, h 之间的权重为 w. x 与 h 之间的权重为 W, 则 Hopfield Network is defined as:

$$E(\mathbf{x}, \mathbf{h}) = -c^T \mathbf{x} - b^T \mathbf{h} - h^T W \mathbf{x} - x^T U \mathbf{h} - h^T V \mathbf{h}$$

A fully-connected Boltzmann machine can represent complex probabilistic relations, but parameter estimation is hard.

Restricted Boltzmann Machines (RBM): Restricted to connections between visible and hidden nodes. Energy function: $E(\mathbf{x}, \mathbf{h}) = -c^T \mathbf{x} - b^T \mathbf{h} - h^T W \mathbf{x}$

Joint probability: $p(\mathbf{x}, \mathbf{h}) = e^{-E(\mathbf{x}, \mathbf{h})} / Z$ 其中 $Z = \sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}$

MLE: $P(\mathbf{x}; \theta) = \sum_{\mathbf{h}} p_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})} / Z$

则 $\frac{\partial \log P(\mathbf{x}; \theta)}{\partial \theta} = \mathbb{E}_{P(\mathbf{x}, \mathbf{h})} \left[\frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} \right] - \mathbb{E}_{P(\mathbf{x}, \mathbf{h})} \left[\frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} \right]$ 其中 $\mathbb{E}_{P(\mathbf{x}, \mathbf{h})}$ 未知, 需要估计.

$p(\mathbf{h} | \mathbf{x}) = \prod_j \text{sigm} \left((2\hat{z}_j - 1) (b_j + w_j^T \cdot \mathbf{x}) \right)$

$p(\mathbf{x} | \mathbf{h}) = \prod_k \text{sigm} \left((2x_k - 1) (c_k + hW_k) \right)$

1980 年代三种主要类型的神经网络 前馈神经网络 Feedforward NN (代表性方法: 多层感知器; 反馈型神经网络 Feedback NN (代表性方法: Hopfield NN); 竞争学习神经网络 Competitive Learning NN (代表性方法: 自组织映射 Self-organizing map)

近邻法

分段线性判别函数: 把各类划分为若干子类, 以子类中心作为代表点. 考查新样本与各类代表点的距离, 极端情况: 近邻法.

K-近邻法: 找出 x 的前 k 个近邻, 看其中多数属于哪一类, 就把 x 分到哪一类. K 越大, 分界线越平滑.

KNN 常见问题: 存储量和计算量: 要数接近时风险较大, 有噪声时风险加大; 样本无穷多时性能优异, 有样本下性能如何? 改进: 减少计算量和存储量; 引入拒绝机制, 根据实际问题修正投票方式.

近邻法在计算上的问题: 需存储所有训练样本; 新样本需与每个样本做比较. 快速算法基本思想: 把样本集分成若干子集 (树状结构), 每个子集 (节点) 可用较少几个样本代表; 通过将新样本与各节点比较排除大量候选样本; 只有最后的节点 (子集) 中这个样本比找出最近邻.

分枝界定算法 (Branch-Bound Algorithm): 已知样本划分成多个子集, 形成一个树状结构, 每个节点是一个子集, 每个子集用若干个的几量代表; 新样本按顺序与各节点进行比较来排除不可能包含近邻的样本, 最后只在少数节点上与每个样本进行比较

剪枝近邻法: 处在两类交界处分布重合区的样本可能误导决策, 应将其从样本集中去掉. 若思考: 考查样本是否可能的误导样本, 若是则从样本集中去掉 (剪枝). 考查方法是: 通过交叉验证, 认为错分样本为误导样本.

压缩近邻法: 主要通过减少存储量. 将训练样本分为 X_0 和 X_c 两个集合, 开始时 X_0 中只有一个样本, X_c 中为其余样本. 考查 X_c 中每个样本, 若用 X_0 可正确分类则保留, 否则移入 X_0 . 最后用 X_0 作分类的样本集.

原型近邻法: 用若干个原型 (prototype) 样本来代表训练数据. 通常原型的数量远少于训练样本数 (如等于样本数则为 1-近邻法).

定义原型点: 利用 k-means 聚类. 对于第 1 类样本, 进行聚类分析, 获得 R 个原型 (该类样本密度高的点). 定义这 R 个原型的类别为 i; 总共获得 CR 个原型样本 (C 为类别数).

对于待分类的样本 x, 我们将其类别为距离它最近的原型的标签.

原型近邻法的优点有: 模型测试时计算开销小; 节省内存. 没有: 模型训练时计算开销小, 可视化效果好; 能更好捕捉数据分布的细节信息.

可做拒绝决策的近邻法: 从简单多数变为绝对多数. 拒绝决策同样可引入改进的近邻法中, 比如剪枝近邻法.

高维情况下需要更多的训练样本, 当标准维度较高时, “最近邻”可能相差很远. 使用欧氏距离通常会对应每一维特征做标准化.

按距离加权的 KNN(LOESS): 距离近的邻居加权大, 距离远的邻居加权小. 一种常用的权值函数 $w^{(i)} = \exp \left[-\frac{(x^{(i)} - x)^2}{\tau^2} \right]$. τ 为带宽. 训练样本少, 带宽应该取得大一些; 模型回归系数 θ^T 与 τ 取值有关.

决策树

节点的分裂需要定义一个度量指标.

Entropy (used in ID3 and C4.5): $H(D) = -\sum_{k=1}^K \frac{|c_k|}{|D|} \log \frac{|c_k|}{|D|}$

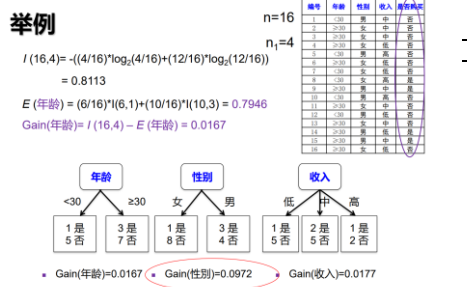
Gini index (used in CART): $Gini(D) = 1 - \sum_{k=1}^K \left(\frac{|c_k|}{|D|} \right)^2$

当 Entropy 最大为 1 的时候, 是分类效果最差的状态, 当它最小为 0 的时候, 是完全分开的状态. 因为熵等于是零理想状态, 一般实际情况下, 嫡介于 0 和 1 之间. 在分裂时, 我们衡量父节点的嫡与两个儿子的嫡的差, 即嫡量分裂后是否能够减少嫡 $H(D) - (H(D_1) + H(D_2))/2$

或者定义**信息增益**: $\text{Gain}(N, a) = I(N, N^{(1)}) - E(N, a)$

上层节点样本数 N, 第一类样本数 $N^{(1)}$, 信息嫡 $I(N, N^{(1)})$. 若按属性 a 分类, 则下层节点, 在属性 a 上有 k 个取值, 取值 i 下有 m_i 个样本, 其中 m_{k1} 属于第一类. $\text{Gain}(N, a)$

$$I(N, N^{(1)}) = -\left(\left(N^{(1)} / N \right) \log_2 (N^{(1)} / N) + \left(1 - N^{(1)} / N \right) \log_2 (1 - N^{(1)} / N) \right)$$

$$E(N, a) = (m_1 / N) (I(m_1, m_{11}) + (m_2 / N) (I(m_2, m_{21}) + \dots + (m_k / N) (I(m_k, m_{k1}))$$


ID3 存在一个问题, 那就是越细小的分割分类错误率越小. 所以为了避免分割太细, C4.5 进行了改进, 优化项要除以分割太细的代价, 比值叫做信息增益率, 分割太细导致增加, 信息增益率会降低. 除此之外, 其他的原理和 ID3 相同.

$\text{Gain}_{\text{ratio}}(N, a) = \frac{\text{Gain}(N, a)}{\frac{-\log_2 \frac{N^{(1)}}{N}}{N}}$ 其中 $P_{ij} = \frac{m_{ij}}{N}$; $v = 1, \dots, k$

CART 与 C4.5 算法的最大相异之处是在每一个节点上都是采用二分法, 也就是一次只能有两个子节点, C4.5 则在每一个节点上可以产生不同数量的分枝. CART 和 ID3 一样, 存在偏向细小分割, 则过度学习 (过度拟合的问题), 为了解决, 对特别长的树进行剪枝处理, 直接剪掉.

如果正负两类样本从参数完全相同的多维正态分布中抽样产生, 把其中一部分用作训练集, 一部分用作测试集, 利用决策树来构建分类器, 则在训练集上可以得到高正确率, 在测试数据集上不能得到高正确率, 因为过拟合.

剪枝 (Pruning): 避免过学习. 先剪枝 (prepruning): 数据划分方法: 利用训练集决定节点划分, 根据测试集的分裂性质决定是否停止; 阈值法: 预先设定信息增益阈值. 信息增益的统计显著性分析. 后剪枝 (Postpruning): 减少分类错误修正法: 用复杂性折减完整的树, 在独立剪枝集上减少分类错误的修剪法: 最小代价与复杂性折中: 合并分枝后错误率增加和复杂度减少折中考虑; 最小描述长度准则: 最简单的树就是最好的树.

集成学习: 个体学习器之间不存在强依赖关系, 可以并行生成学习器. Bagging, Random Forest: 个体学习器之间存在强依赖关系, 串行生成学习器. Boosting.

Bagging (Bootstrap Aggregating): 对训练数据 D 有放回抽样生成 m 个数据集 D_i , 在每一个数据集 D_i 上训练分类器 h_i , 对各个分类器的结果集成投票. 防止过拟合, 减小模型方差.

Random Forests: 利用 bootstrapping 的办法 (有放回抽样). 训练多棵树 (每次从 N 个训练样本中有放回地随机抽取 N 个样本作为当前训练集, 每棵树随机抽取 k (k<p) 个特征作为当前节点下决策的备选特征, 从中选择特征进行划分 (split), 最后对每一颗树的预测结果进行汇总投票.

随机选取特征的好处: 用越少的特征去构造决策树, 决策树就不容易过拟合; 在 bootstrap samples 中, 如果这个特征非常重要, 则它会在不同的 sample 中都出现 (或者说出现的频率比较高), 则这样就会带来 Correlation, 而随机森林的原则是避免相关性, 因此随机选取 feature.

为什么效果好? 每决策树的决策面都是平行坐标轴的立方体, 多颗决策树组合的结果可形成复杂的决策面. 随机森林获得的每棵树之间结构不同, 参数不同. KNN 中需要对各维度特征进行 Normalization, 使用决策树时不需要.

Boosting: 通过迭代对学习器的输入输出进行加权组合, additive training, 训练新学习器拟合残差. 开始时所有训练样本同样权重, 每一轮学习一个“弱”学习器, 根

据当前“弱”学习器的组合预测结果进行调整, 对预测偏差给予更多关注. 直到“弱”学习器数目达到预先规定的数目 k, 最终将 k 个学习器的结果加权组合.

AdaBoost: Adaptive Boosting. 计算当前轮次对应的错误率与正确率, 调整正确样本与错误样本的权重, 使得下一轮分类器的训练样本一半来自分对的, 一半来自分错的. 训练过程是串行的. 计算错误率 e_m , 权重 $\alpha_i = \log((1 - \epsilon_i) / \epsilon_i)$. 更新 $D(x_i) = D(x_i) \cdot \exp(\alpha_i \cdot [y_i \neq h_i(x_i)])$

Gradient Boosting: Gradient Descent + Boosting. 输入: $\{(x_i, y_i)\}_{i=1}^n$; 损失函数 $L(y, F(x))$; 迭代次数 k_{max} .

算法流程: 初始化 $F_0(x) = \arg \min_{\eta} \sum_{i=1}^N L(\eta_i, y_i)$. 循环 $k = 1 \dots k_{max}$: 计算 $\eta_k = -\frac{\partial L(y, F(x))}{\partial F(x)}|_{F(x)=F_{k-1}(x)}$; 用数据集 $\{(x_i, \eta_k)\}_{i=1}^n$ 训练新学习器 $f_k(x)$; 计算更新权重: $\eta_k = \arg \min_{\eta} \sum_{i=1}^N L(y_i, F_{k-1}(x_i) + \eta f_k(x_i))$; 更新模型: $F_k(x) = F_{k-1}(x) + \eta_k f_k(x)$

XGBoost: Extreme Gradient Boosting. 改进 1: 目标函数引入正则项 (防止过拟合); 改进 2: 引入导数项简化目标函数. 随机森林, Gradient Boosting 和 XGBoost 都是由学习器集成得到最终结果, 适用于分类和回归. 随机森林对样本和特征降采样, 不易过拟合, 多模型使用均值, 对异常值不敏感. 未设计稀疏数据预处理策略; Gradient Boosting 使用全部样本, 容易过拟合, 模型加约束和, 对异常值敏感, 未设计稀疏数据预处理策略; XGBoost 对样本和特征降采样, 正则项防止过拟合, 模型加约束和, 对异常值敏感, 适用于稀疏数据. **Boosting 方法的特点:** 快速: 简单: 只有一个参数: 灵活: 可以与任何分类器使用. 只需要把随机机器的“弱”分类器, 缺点: “弱”分类器不能太复杂 (容易过拟合); 对噪声敏感, 例如会给标记错误的样本过大的权重. **神经网络防止过拟合:** Dropout. 每次训练时随机把部分节点参数置零, 相当于利用很多不同结构网络的结果的集合.

支持向量机

若向量集合被超平面 $f(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$ 没有错误地分开, 且离超平面最近的向量与超平面之间的距离 (称作间隔 margin) 是最大的, 则我们说这个向量集合被这个**最优超平面** (或最大间隔超平面) 分开

要留下一定的“裕度”, 故要把尺度固定下来: $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$

点到平面距离 $\frac{|\mathbf{w} \cdot \mathbf{x}_i + b|}{\|\mathbf{w}\|}$. 故 margin 为 $2 / \|\mathbf{w}\|$

优化问题 $\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$, s.t. $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, 1 \leq i \leq n$.

拉格朗日函数 $L(W, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i [y_i(\mathbf{x}_i^T \mathbf{w} + b) - 1]$. 即 $\min_{\mathbf{w}, b, \alpha} \max_{\mathbf{w}, b, \alpha} L(W, b, \alpha)$. 对偶问题: $\max_{\alpha} \min_{\mathbf{w}, b} L(W, b, \alpha)$

求导: $\nabla_{\mathbf{w}} L(W, b, \alpha) = \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$

$\nabla_b L(W, b, \alpha) = -\sum_{i=1}^n \alpha_i = 0 \Rightarrow \sum_{i=1}^n \alpha_i = 0$

代入后结果为 $\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j$

取极大对应取负后极小: $\min_{\alpha} \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \sum_{i=1}^n \alpha_i$

解得 $\hat{\mathbf{w}} = \sum_{i=1}^n \hat{\alpha}_i y_i \mathbf{x}_i$, $\hat{b} = y_j - \sum_{i=1}^n \hat{\alpha}_i y_i \mathbf{x}_i^T \mathbf{x}_j$

判别函数 $f(\mathbf{x}) = \text{sign}(\sum_{i=1}^n \hat{\alpha}_i y_i \mathbf{x}_i^T \mathbf{x}_i + \hat{b})$

软间隔: 引入惩罚参数 $\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i$, s.t. $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$, $\xi_i \geq 0, 1 \leq i \leq n$ 与最优硬平面类似, 广义最优超平面满足下列 KKT 条件: $\alpha_i^* \{[(\mathbf{x}_i - \mathbf{w}_0) + b_0] y_i - 1 + \xi_i\} = 0, i = 1, \dots, l$. 只有部分样本的 α_i 不为 0, 是不等式 $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$

线性不可分情况下的支持向量包括两部分: 正确分类且离分类面最近的样本, 它们的 $0 < \alpha_i < C$; 错误分类的样本, 它们的 $\alpha_i = C$

推广到非线性

挑战 1: 能否对付在特征空间中的运算 (如升到高维是否会出現维度灾难). 实际上, 变换空间中的内积运算仍然对应着原空间中的某种运算. 核函数必须是某个空间中的内积. 在核函数给定的情况下, 可以利用解线性问题的方法求解非线性问题的支持向量机, 此过程是隐式地在特征空间中进行的.

RBF 核: $K(\mathbf{x}, \mathbf{x}_i) = \exp \left[-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{2\sigma^2} \right]$ 线性核函数的 SVM 一定会得到线性分类面.

支持向量回归 (SVR). 原问题: 用函数 $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$ 拟合样本. 如果所有样本均在精度下被拟合, 即 $|y_i - \mathbf{w} \cdot \mathbf{x}_i - b| \leq \epsilon$, 最小化 $\frac{1}{2} \|\mathbf{w}\|^2$ 以获得最优回归.

引入松弛变量 $\xi_i \geq 0, i \geq 0$, 原问题变为 $\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*)$, s.t. $|y_i - \mathbf{w} \cdot \mathbf{x}_i - b| \leq \epsilon + \xi_i^*$

回归函数 $f(\mathbf{x}) = (\mathbf{w} \cdot \mathbf{x}) + b = \sum_{i=1}^l (\alpha_i^* - \alpha_i) (\mathbf{x}_i \cdot \mathbf{x}) + b^*$

支持向量: 只有 $\alpha_i^{(*)} = C$ 的样本落在“ ϵ -不敏感管带”外 (类比 SVM 分类中的错分样本). 满足 $\alpha_i^{(*)} \in (0, C)$ 的样本有 $\xi_i^{(*)} = 0$ 且 $|y_i - f(\mathbf{x}_i)| = \epsilon$ (类比 SVM 分类中的边界样本)

统计学习理论

为什么最大间隔是最优? 大间隔 > 低 VC 维 > 低复杂度 > 高推广能力.

经验风险 $R_{emp}(\alpha) = \frac{1}{N} \sum_{i=1}^N Q(z_i, \alpha)$. 经验风险最小化 Empirical Risk Minimization 原则: 用最小化 $Q(z, \alpha_i)$ 的近似 $Q(z, \alpha_0)$

初心 $\min R(\alpha) = \int Q(z, \alpha) dF(z)$ 实际行 $\min R_{emp}(\alpha) = \frac{1}{N} \sum_{i=1}^N Q(z_i, \alpha)$

一个指示函数集 $Q(z, \alpha)$, $\alpha \in A$ 的 VC 维, 是能被集合中的函数以所有可能的 2^h 种方式分成两类的向量 $Z_1, \dots, Z_n</$

$$f(x) = \frac{1}{1+e^{-x}}, \quad f'(x) = \frac{e^{-x}}{(1+e^{-x})^2}, \quad g(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}, \quad g'(x) = \frac{4}{(e^x + e^{-x})^2}$$