

第7章 贝叶斯分类器

1.1 贝叶斯决策论

$Y = \{c_1, c_2, \dots, c_N\}$ λ_{ij} : 真实为 c_j 误分为 c_i 的损失
 $P(c_i|x)$ 后验概率

样本分为 c_i 所产生的期望损失:

$$R(c_i|x) = \sum_{j=1}^N \lambda_{ij} P(c_j|x)$$

判定准则 $h: X \rightarrow Y$ 最小化总风险

$$R(h) = E_x[R(h(x)|x)]$$

贝叶斯判定准则:

最小化总风险, 只需在每一样本上选择
 那个能使条件风险最小的类别标记.

$$h^*(x) = \underset{c \in Y}{\operatorname{argmin}} R(c|x)$$

$$\text{若 } \lambda_{ij} \Rightarrow \begin{cases} 0 & \text{if } i=j \\ 1 & \text{otherwise} \end{cases}$$

此时条件风险:

$$R(c|x) = 1 - P(c|x)$$

于是最小化分类错误率的贝叶斯分类器

$$h^*(x) = \underset{c \in Y}{\operatorname{argmax}} P(c|x)$$

对每一样本, 选择能使后验概率 $P(c|x)$ 最大的类别标记

判别式模型: 给定 x , 可通过直接建模 $P(c|x)$ 来预测 c (决策树, BP, SVM)

生成式模型: 先对联合分布 $P(x, c)$ 建模

如: 高斯判别分析 $P(c|x)$

$$\text{生成式: } P(c|x) = \frac{P(x, c)}{P(x)}$$

$$\text{且 } P(c|x) = \frac{P(c)P(x|c)}{P(x)} \rightarrow \text{贝叶斯定理}$$

\downarrow 先验 \downarrow 后验
 证据因子
 给定样本 x
 与类别 c 有关

1.2 极大似然估计

假设 $P(x|c)$ 具有确定的形式并且被参数向量 θ_c 唯一确定

用训练集 D 估计参数 θ_c

样本模型训练 \longleftrightarrow 参数估计

令 D_c 表示训练集 D 中类 c 的样本组成集合
 设样本独立同分布, 则参数 θ_c

$$P(D_c|\theta_c) = \prod_{x \in D_c} P(x|\theta_c)$$

对 θ_c 进行极大似然估计

$$\begin{aligned} LL(\theta_c) &= \log P(D_c|\theta_c) \\ &= \sum_{x \in D_c} \log P(x|\theta_c) \end{aligned}$$

$$\hat{\theta}_c = \underset{\theta_c}{\operatorname{argmax}} LL(\theta_c)$$

1.3 朴素贝叶斯分类器

"属性条件独立性假设" ^{避免} 组合爆炸
 样本: 每一属性独立地对分类结果产生影响

$$P(c|x) = \frac{P(c)P(x|c)}{P(x)} = \frac{P(c)}{P(x)} \prod_{i=1}^d P(x_i|c)$$

$$h_{NB}(x) = \underset{c \in Y}{\operatorname{argmax}} P(c) \prod_{i=1}^d P(x_i|c)$$

朴素贝叶斯常需要“平滑”

$$\hat{p}(c) = \frac{|D_c| + 1}{|D| + N}$$

$$\hat{p}(x_i|c) = \frac{|D_{c,x_i}| + 1}{|D_c| + N_i}$$

7.4 朴素贝叶斯分类器

适当考虑一部分属性间：相互依赖
信息

“树依赖估计”：OPE

$$P(c|x) \propto p(c) \prod_{i=1}^d p(x_i|c, \text{par}_i)$$

par_i 是 x_i 所依赖：属性，父属性

如何确定？

① 假设所有属性都依赖于同一属性

spooE 方法

通过交叉验证方法

② TAN

7.5 贝叶斯网