

朴素贝叶斯

W.J.Z

摘要：朴素贝叶斯法是基于贝叶斯原理和特征条件独立假设的分类方法。对于给定的输出集，首先基于特征条件独立假设学习输入/输出之间的联合概率分布，然后基于此模型，利用贝叶斯定理求出后验概率最大的输出。

0.1 朴素贝叶斯原理

朴素贝叶斯法采用期望风险最小化（误分类的样本尽可能少），选择 0-1 损失函数：

$$L(Y, f(X)) = \begin{cases} 1, & Y \neq f(X) \\ 0, & Y = f(X) \end{cases} \quad (1)$$

期望风险函数为：

$$E[L(Y, f(X))] \quad (2)$$

注意公式 2 是针对的联合分布，而我们需要的是条件期望，根据 0-1 分布期望公式，可得条件期望公式为：

$$\sum_{i=1}^n (L(c_i, f(X)) P(c_i|X)) \quad (3)$$

$c_i \in Y$ 我们要使期望风险最小化，也就是使公式 3 最小化，将公式 1 带入公式 3 得：

$$\begin{aligned} f(x) &= \arg \min \sum_{i=1}^n (L(c_i, f(X)) P(c_i|X = x)) \\ &= \arg \min \sum_{i=1}^n (1 \times P(y \neq c_i|X = x) + 0 \times P(y = c_i|X = x)) \\ &= \arg \min \sum_{i=1}^n (P(y \neq c_i|X = x)) \\ &= \arg \min (1 - (P(y = c_i|X = x))) \\ &= \arg \max P(y = c_i|X = x) \end{aligned}$$

后验概率最大化准则：

$$f(x) = \arg \sum_{c_i} (P(Y = c_i|X = x)) \quad (4)$$

0.2 朴素贝叶斯公式

朴素贝叶斯进行了条件独立性假设，条件独立性假设为：

$$P(X = x|Y = c_k) = P(x_1, x_2, \dots, x_n|Y = c_k) \quad (5)$$

$$= \prod_{i=1}^n P(x_i|Y = c_k) \quad (6)$$

后验概率公式：

$$P(Y = c_i|X = x) = \frac{P(X = x|Y = c_i) P(Y = c_i)}{\sum_i P(X = x|Y = c_i) P(Y = c_i)} \quad (7)$$

将公式 6 带入公式 7 得：

$$P(Y = c_k|X = x) = \frac{P(Y = c_k) \prod_i P(x_i|Y = c_k)}{\sum P(Y = c_k) \prod_i P(x_i|Y = c_k)} \quad (8)$$

由于分母对于任何 c_k 都是一样的，综上可得朴素贝叶斯分类器可表示为：

$$y = \arg \max_{c_k} P(Y = c_k) \prod_i P(x_i|Y = c_k) \quad (9)$$

0.3 极大似估计

先验概率 $P(Y = c_k)$ 的极大似然估计为：

$$P(Y = c_k) = \frac{\sum_{i=1}^N I(y_i = c_k)}{N} \quad (10)$$

条件概率的极大似然估计：

$$P(x_i|Y) = \frac{\sum_{i=1}^N I(x_i, y_i = x_k)}{\sum_{i=1}^N I(y_i = c_k)} \quad (11)$$

参考：统计学习方法