

为了避免其他属性携带的信息

被训练集中未出现的属性值“抹去”，在估计概率值时通常要进行“平滑”，常用“拉普拉斯修正”

$N$  表示训练集中可能的类别数。

$$P(c) = \frac{|D_c|}{|D|}$$

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

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

朴素贝叶斯分类器

对“属性条件独立性假设”进行一定程度的放松



### III.

EM算法

$Z$  表示隐变量集

$\theta$  表示模型参数

$X$  表示已观测变量

$$LL(\theta | X, Z) = \ln P(X, Z | \theta)$$

$$LL(\theta | X) = \ln P(X | \theta) = \ln \sum_Z P(X, Z | \theta)$$

EM常用来估计参数隐变量，迭代式的方法。

以初始值  $\theta^0$  为起点：

E① 基于  $\theta^t$  推断隐变量  $Z$  的期望，记为  $Z^t$

M② 基于已观测变量  $X$  和  $Z^t$  对参数  $\theta$  做极大

似然估计，记为  $\theta^{t+1}$

直到  $\theta^t$  和  $Z^t$  稳定

基于  $\theta^t$  计算隐变量  $Z$  的概率分布  $P(Z | X, \theta^t)$

E.① 以当前  $\theta^t$  推断  $P(Z | X, \theta^t)$ ，并计算对数似然

$$LL(\theta | X, Z)$$

$$Q(\theta | \theta^t) = E_{Z | X, \theta^t} [LL(\theta | X, Z)]$$