

生成模型必然考虑：

$$P(c|x) = \frac{P(x, c)}{P(x)}$$

$$P(c|x) = \frac{P(c) P(x|c)}{P(x)}$$

$P(c)$ ：样本空间中各类样本所占的比例

极大似然估计：源自频率主义学派，根据数据采样来估计概率分布参数。

$$P(D_c | \theta_c) = \prod_{x \in D_c} P(x | \theta_c) \quad (\text{使用 } \theta_c \text{ 找到 } D_c \text{ 中所有样本的概率})$$

在所有 θ_c 的可能取值中，找到一个使 $P(D_c | \theta_c)$ 最大的值。

$$\mathcal{L}(\theta_c) = \log P(D_c | \theta_c) = \sum_{x \in D_c} \log P(x | \theta_c)$$

$$P(x|c) \sim N(\mu_c, \sigma_c^2)$$

II

朴素贝叶斯分类器

利用“属性条件独立性假设”：对已知类别，假设所有属性相互独立，每个属性独立地对分类结果产生影响。

$$P(c|x) = \frac{P(c) P(x|c)}{P(x)}$$

x_1 = 青绿, 大, 西瓜
 x_2 = 青绿
 x_3 = 大
 x_4 = 西瓜

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

$$P(x|c) = P(x_1|c) P(x_2|c) P(x_3|c)$$

$$= P(x_1 x_2 x_3 | c)$$

朴素贝叶斯判定规则：

$$h_{nb}(x) = \arg \max_{c \in Y} P(c) \prod_{i=1}^d P(x_i | c)$$

基于训练集 D ，估计类先验概率 $P(c)$

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

D_{c, x_i} 表示 D_c 中在第 i 个属性上取值为 x_i 的样本组成的集合。

$$P(x_i | c) = \frac{|D_{c, x_i}|}{|D_c|}$$