

## 朴素贝叶斯 (Naive Bayesian Classifier)

限制条件: ①  $X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$  每个维度离散  
②  $X$  的每个维度独立

应用: 垃圾邮件分类

输入: 一个文件  $d$

输出:  $d \in C_1$  或  $d \in C_2$

训练样本:  $\{(d_i, c_i) \mid i=1 \sim N\}$

$\swarrow$   
 $d = \{w_1, w_2, \dots, w_n\}$   
 $\uparrow$   
(单词)

学习  $P(d|C_1)$  与  $P(d|C_2)$

$$P(d|C)$$

$$= P(\{w_1, \dots, w_n\} | C)$$

$$= \prod_{i=1}^n P(w_i | C) \quad (\text{独立假设})$$

$$P(w | C_j) = \frac{\text{count}(w, C_j) + 1}{\sum_{\substack{w \in V \\ (\text{种类})}} \text{count}(w, C_j) + |V|}$$

$(j=1 \sim 2)$

$$p(c_1) = \frac{c_1 \text{的个数}}{\text{总个数}}$$

$$p(c_2) = \frac{c_2 \text{的个数}}{\text{总个数}}$$

若  $p(c_1) \cdot p(d|c_1) > p(c_2) \cdot p(d|c_2)$  则  $d \in c_1$   
 否则  $\dots$   $d \in c_2$