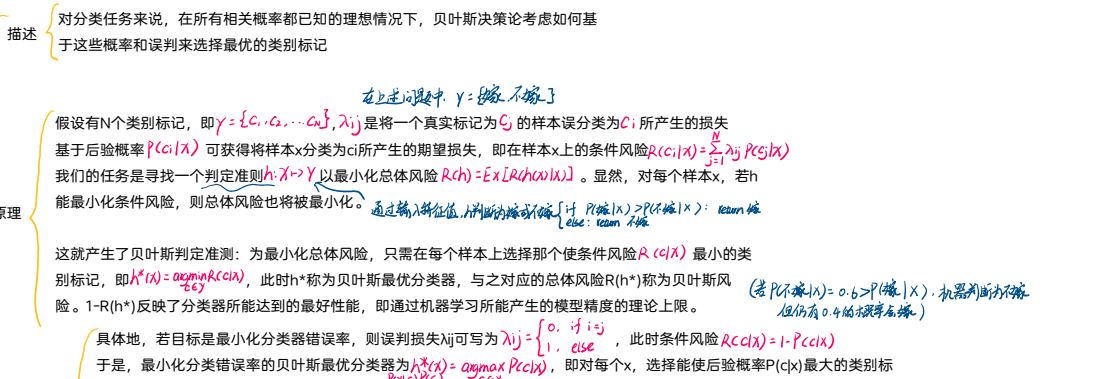


# 贝叶斯分类器

## 贝叶斯定理



## 贝叶斯决策论



极大似然估计

解决问题

极大似然估计

Maximum Likelihood Estimation

公式转换

险。1-R(h\*)反映了分类器所能达到的最好性能，即通过机器学习所能产生的模型精度的理论上限。

具体地，若目标是最小化分类器错误率，则误判损失 $\lambda_{ij}$ 可写为 $\lambda_{ij} = \begin{cases} 0, & \text{if } i=j \\ 1, & \text{else} \end{cases}$ ，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最优分类器为 $\hat{y}(x) = \underset{c \in Y}{\operatorname{argmax}} P(c|x)$ ，即对每个x，选择能使后验概率P(c|x)最大的类别标记。基于贝叶斯定理，P(c|x)可写为 $\frac{P(c)P(x|c)}{P(x)}$ 。其中，P(c)是类先验概率，P(x|c)是样本x相对于类标记c的类条件概率，或称“似然”；P(x)是用于归一化的“证据”因子。

对给定样本x，证据因子P(x)与类标记无关，因此估计P(c|x)的问题就转化为如何基于训练数据D来估计先验P(c)和似然P(x|c)。

根据大数定律，当训练集包含充足的独立分布样本使，P(c)可通过各类样本出现的频率来进行估计。

对类条件概率P(x|c)来说，由于它涉及关于x所有属性的联合概率，直接根据样本出现的频率来估计将会遇到严重的困难，直接使用频率来估计P(x|c)显然不行，因为“未被观测到”和“出现概率为零”通常是不同的。

似然

假设现有一枚特殊的硬币，每次抛出后得到是花或字的概率可能并不为1/2。在这种情况下，通过抛硬币的方式，如果抛100次硬币的结果都是花，则可认为这枚硬币的两面都是花的可能性最大。这种通过事实反过来来猜测硬币的情况就是似然。通过事实推断出最有可能的硬币情况，就是最大似然估计。

似然和概率

已知硬币具有花和字两面，推测抛硬币的各种情况的可能性，称为概率。

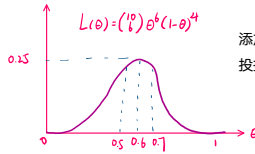
如果对硬币的两面是什么并不清楚，需要通过抛硬币的情况来确定，称为似然。

假设在一次实验中，10次抛硬币，有6次是花，则硬币抛出后是花的概率(以下简称为硬币的参数)有多大？

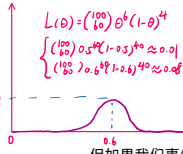
所谓极大似然估计，就是假设硬币的参数，然后计算实验结果的概率是多少。通过对比设置不同参数下得到的概率，概率越大的越可能是正确的。

用0.5作为参数，概率为 $\binom{10}{6} 0.5^6 (1-0.5)^4 \approx 0.21$ ，再用0.6作参数： $\binom{10}{6} 0.6^6 (1-0.6)^4 \approx 0.25$ ，可见 $0.25 > 0.21$ ， $\frac{0.25}{0.21} \approx 1.2$ ，可以认为0.6作为参数的可能性是0.5作为参数的可能性的1.2倍。

设置硬币的参数为 $\theta$ ，可以得到似然函数：



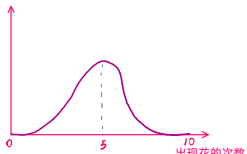
添加更多的实验结果：  
投掷100次，出现了60次花



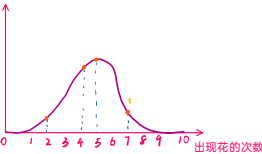
此时，0.6作为参数的可能性是0.5作为参数的可能性的8倍，新的实验结果更加支持0.6这个参数。

图像有明显的缩窄，可理解为可选参数的分布更集中了。通过投入更多的是实验结果，参数会越来越明确。

实际上，如果这是一枚普通的硬币，我们知道，抛出后是花的概率即硬币的参数 $\theta=0.5$ ，这里绘制出抛10次硬币后出现花的次数的概率分布图



但如果我们事先并不知道它是否是普通的硬币，只能通过实验的方式来推断。我们进行了6次实验后（一次实验抛10次硬币，统计花的次数），实验结果为 $\{4, 5, 5, 2, 7, 4\}$ ，表示第一次实验花的次数是4，第二次实验花的次数是5，以此类推，可在图中画出对应的点



我们用 $x_1, x_2, \dots, x_n$ 表示每次实验结果，因为每次实验都是独立的，所以似然函数可写作 $L(\theta) = f(x_1|\theta)f(x_2|\theta)\dots f(x_n|\theta)$ ，其中 $f(x_i|\theta)$ 表示在同一个参数下的结果。随着试验次数的增加， $\theta$ 的值会逐渐逼近真实值。但由于实验本身具有二项随机性，可能会导致它与真实值存在细微的误差。

记关于类别c的类条件概率为P(x|c)，假设P(x|c)具有确定的形式并且被参数向量 $\theta_c$ 唯一确定，则我们的任务就是利用训练集D估计参数 $\theta_c$ 。这个概率模型训练的过程就是参数估计过程，通过极大似然估计法解决问题。

令 $D_c$ 表示训练集D中第c类样本组成的集合，假设这些样本是独立同分布的，则参数 $\theta_c$ 对于数据集 $D_c$ 的似然是 $P(D_c|\theta_c) = \prod_{x \in D_c} P(x|\theta_c)$

对 $\theta_c$ 进行极大似然估计，就是去寻找能最大化似然P( $D_c|\theta_c$ )的参数 $\hat{\theta}_c$ 。直观上看，极大似然估计是试图从 $\theta_c$ 所有可能的取值中，找到一个能使数据出现的“可能性”最大的值。为了防止上式连乘操作造成下溢，通常使用对数似然： $LL(\theta_c) = \log P(D_c|\theta_c) = \sum_{x \in D_c} \log P(x|\theta_c)$

此时参数 $\theta_c$ 的极大似然估计 $\hat{\theta}_c$ 为 $\hat{\theta}_c = \underset{\theta_c}{\operatorname{argmax}} LL(\theta_c)$ 。

例如，在连续属性情形下，假设概率密度函数 $p(x|\mu, \sigma^2) \sim N(\mu, \sigma^2)$ ，则参数 $\mu_c$ 和 $\sigma_c^2$ 的极大似然估计为 $\hat{\mu}_c = \frac{1}{|D_c|} \sum_{x \in D_c} x$ ， $\hat{\sigma}_c^2 = \frac{1}{|D_c|} \sum_{x \in D_c} (x - \hat{\mu}_c)(x - \hat{\mu}_c)^T$

也就是说，通过极大似然估计法得到的正态分布均值就是样本均值，方差就是 $(x - \hat{\mu}_c)(x - \hat{\mu}_c)^T$ 的均值，这显然是一个符合直觉的结果。在离散属性情形下，也可通过类似的方式估计类条件概率。

应注意，由于估计结果的准确性严重依赖于所假设的概率分布形式是否符合潜在的真实数据分布，所以在现实应用中，还需要利用关于应用任务本身的经验知识，否则若仅凭“猜测”来假设概率分布形式，很可能产生误导性的结果。

朴素贝叶

公式

属性条件假设

基于属性条件独立性假设P(c|x)可重写为

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

其中d为属性数目， $x_i$ 为x在第i个属性上的取值

朴素贝叶斯分类器表达式

由于对所有类别来说P(x)相同，因此可将贝叶斯判定准则改写为

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

求解关键

显然，朴素贝叶斯分类器的训练过程就是基于训练集D来估计类先验概率P(c)，并为每个属性估计条件概率P(x<sub>i</sub>|c)。

令 $D_c$ 表示训练集D中第c类样本组成的集合，若有充足的独立同分布样本，则可容易地估计先验概率 $P(c) = \frac{|D_c|}{|D|}$

则条件概率P(x<sub>i</sub>|c)可估计为 $P(x_i|c) = \frac{|D_{c,i}|}{|D_c|}$  [为离散属性]

对连续属性可考虑概率密度函数，假定 $p(x_i|\mu_c, \sigma_c^2) \sim N(\mu_c, \sigma_c^2)$ ，其中 $\mu_{c,i}$ 和 $\sigma_{c,i}^2$ 分别是第c类样本在第i个属性上取值的均值和方差，则有 $p(x_i|\mu_c, \sigma_c^2) = \frac{1}{\sqrt{2\pi}\sigma_{c,i}} \exp(-\frac{(x_i - \mu_{c,i})^2}{2\sigma_{c,i}^2})$

西瓜数据集（1-8号为好瓜，9-17号为坏瓜）

编号	色泽	根蒂	敲声	纹理	脐部	触感	密度	含糖率	好瓜
1	青绿	蜷缩	浊响	清晰	凹脐	硬滑	0.691	0.460	是
2	青绿	蜷缩	...	...	...	...	...	...	...
...	...	...	...	...	...	...	...	...	...

# 朴素贝叶斯分类器

例：西瓜数据集

编号	色泽	根蒂	敲声	纹理	脐部	触感	密度	含糖率	好瓜
1	青绿	蜷缩	浊响	清晰	凹陷	硬滑	0.697	0.460	是
2	青绿	蜷缩	浊响	清晰	凹陷	硬滑	0.697	0.460	是
...	...	...	...	...	...	...	...	...	...
17	浅白	...	...	...	...	...	...	...	否

目标：根据属性判断是否是好瓜

编号	色泽	根蒂	敲声	纹理	脐部	触感	密度	含糖率	好瓜
Test	青绿	蜷缩	浊响	清晰	凹陷	硬滑	0.697	0.460	?

注意：某属性值为连续值

$$P(c) \begin{cases} P(\text{好瓜}=是) = \frac{8}{17} \approx 0.471 \\ P(\text{好瓜}=否) = \frac{9}{17} \approx 0.529 \end{cases}$$

离散属性

为每个属性估计条件概率  $P(x_i|c)$   
 $P(\text{青绿}|是) = P(\text{色泽}=\text{青绿}|\text{好瓜}=是) = \frac{3}{8} \approx 0.375$   
 $P(\text{青绿}|否) = P(\text{色泽}=\text{青绿}|\text{好瓜}=否) = \frac{4}{9} \approx 0.444$   
 $P(\text{蜷缩}|是) = P(\text{根蒂}=\text{蜷缩}|\text{好瓜}=是) = \frac{3}{8} \approx 0.375$   
 $P(\text{蜷缩}|否) = P(\text{根蒂}=\text{蜷缩}|\text{好瓜}=否) = \frac{4}{9} \approx 0.444$   
.....  
 $P(\text{硬滑}|否) = P(\text{触感}=\text{硬滑}|\text{好瓜}=否) = \frac{4}{9} \approx 0.444$

连续属性

$p(\text{密度}=0.697|是) = p(\text{密度}=0.697|\text{好瓜}=是)$  均值  
 $= \frac{1}{\sqrt{2\pi} \times 0.129} \exp\left(-\frac{(0.697-0.697)^2}{2 \times 0.129^2}\right) \approx 1.959$  正态分布  
 $p(\text{密度}=0.697|否) = p(\text{密度}=0.697|\text{好瓜}=否)$  均值  
 $= \frac{1}{\sqrt{2\pi} \times 0.175} \exp\left(-\frac{(0.697-0.460)^2}{2 \times 0.175^2}\right) \approx 1.200$   
.....  
 $p(\text{含糖}=0.460|否) = p(\text{含糖}=0.460|\text{好瓜}=否)$   
 $= \frac{1}{\sqrt{2\pi} \times 0.148} \exp\left(-\frac{(0.460-0.460)^2}{2 \times 0.148^2}\right) \approx 0.066$  概率密度

$$P(c) * P(x_i|c) \begin{cases} \text{于是就有} \\ P(\text{好瓜}=是) * P(\text{青绿}|是) * P(\text{蜷缩}|是) * \dots * p(\text{含糖}: 0.460|是) \approx 0.038 \\ P(\text{好瓜}=否) * P(\text{青绿}|否) * P(\text{蜷缩}|否) * \dots * p(\text{含糖}: 0.460|否) \approx 6.80 * 10^{-5} \end{cases}$$

得出结论  $\begin{cases} \text{由于 } 0.038 > 6.80 * 10^{-5}, \text{ 因此朴素贝叶斯分类器将测试样} \\ \text{本Test判别为 "好瓜"} \end{cases}$