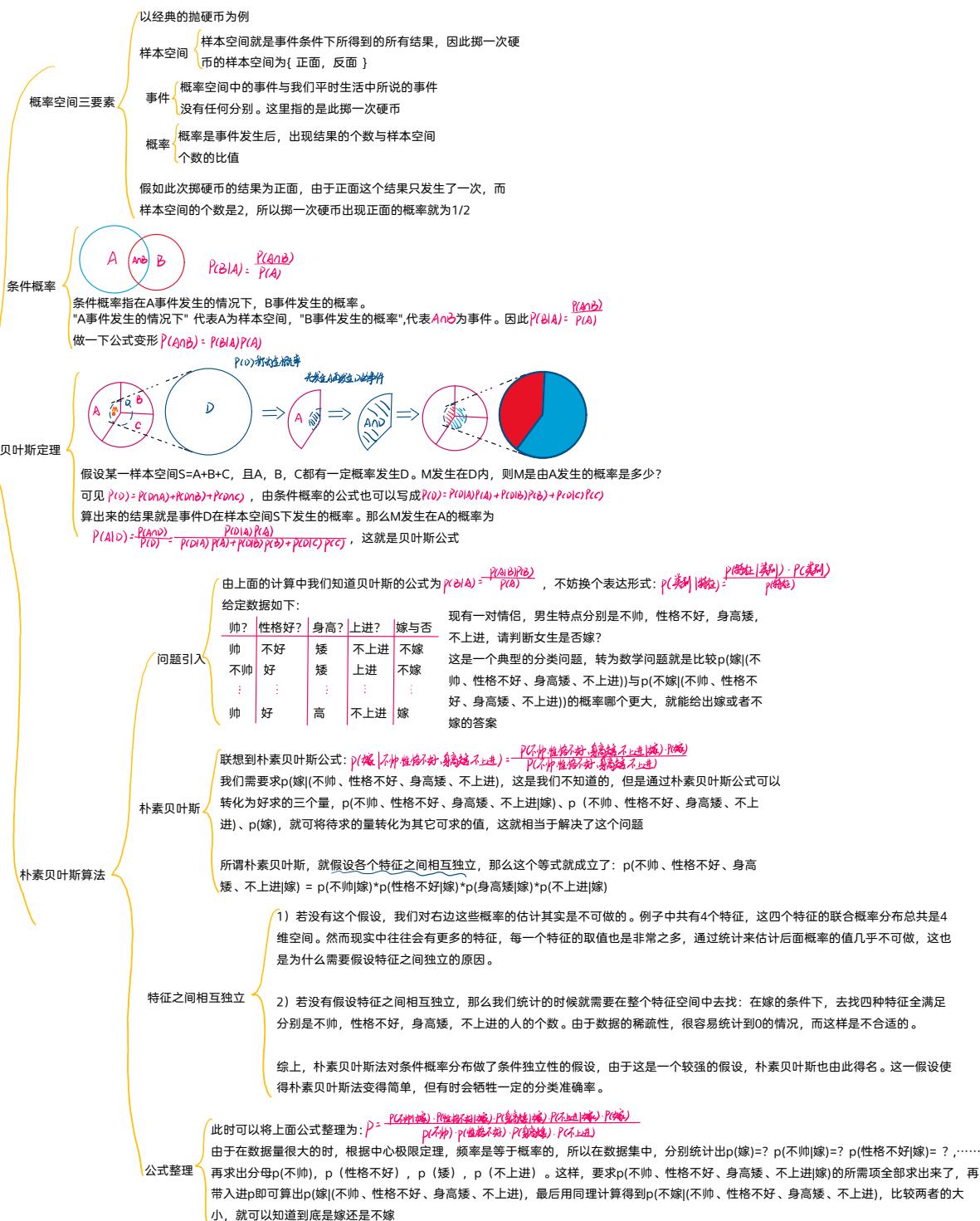
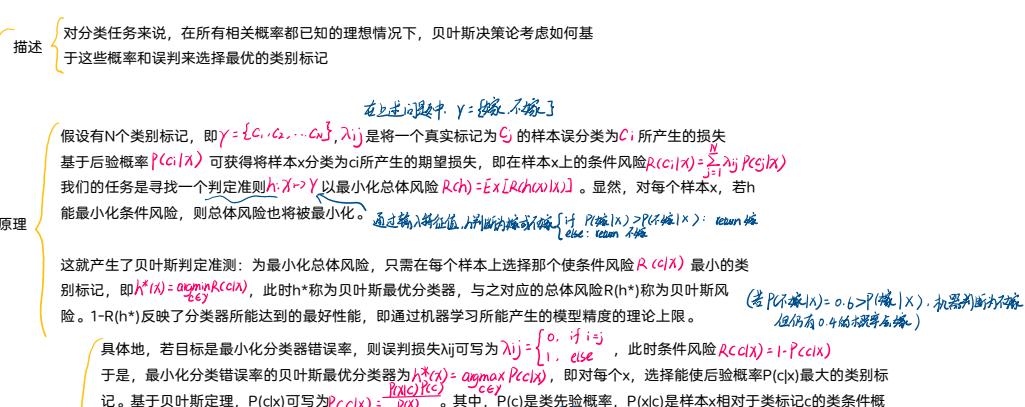


贝叶斯分类器

贝叶斯定理



贝叶斯决策论



朴素贝叶斯

极大似然估计

公式转换

具体地,若目标是最小化分类器错误率,则误判损失 λ_{ij} 可写为 $\lambda_{ij} = \begin{cases} 0, & \text{if } i=j \\ 1, & \text{else} \end{cases}$,此时条件风险 $RCC(x) = 1 - P(c|x)$

于是,最小化分类错误率的贝叶斯最优分类器 $h^*(x) = \operatorname{argmax}_{c \in C} P(c|x)$,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 $\frac{P(c)P(x|c)}{\sum_{c' \in C} P(c')P(x|c')}$ 。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时, $P(c)$ 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0.5。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。这种通过对事反过来自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,就是极大似然估计。

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。通过比对设置不同参数下得到的概率,概率越大的可能性是正确的。

用0.5作为参数,概率为 $(\frac{1}{2})^{10} \cdot (\frac{1}{2})^6 \approx 0.01$,再用0.6作参数: $(\frac{6}{10})^{10} \cdot (\frac{4}{10})^6 \approx 0.25$,可见 $0.25 > 0.01$, $\frac{0.25}{0.01} \approx 25$ 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25倍。

设置硬币的参数为 θ ,可以得到似然函数:

此时如果打开上帝视角,可以观察到抛出后是花的概率即硬币的参数其实是 $\theta=0.5$,所以这里绘制出抛10次硬币后出现花的次数的概率分布图

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 $L(\theta) = f_1(x_1|\theta) f_2(x_2|\theta) \dots f_n(x_n|\theta)$,其中 $f_i(x_i|\theta)$ 表示在同一个参数下的结果。随着试验次数的增加, θ 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。

极大似然估计

Maximum Likelihood Estimation

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 θ_c 唯一确定,则我们的任务就是利用训练集 D 估计参数 θ_c 。这个概率模型训练的过程就是参数估计过程,通过极大似然法解决这个问题。

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 θ_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 θ_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 θ_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。为了防止上式乘积操作造成下溢,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 θ_c 的极大似然估计 $\hat{\theta}_c$ 为 $\hat{\theta}_c = \operatorname{argmax}_{\theta_c} LL(\theta_c)$ 。

例如,在连续属性情形下,假设概率密度函数 $p(x|c) \sim N(\mu_c, \sigma_c^2)$,则参数 μ_c 和 σ_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。

公式

属性条件假设

基于属性条件独立性假设 $P(x|c)$ 可重写为 $P(c|x) = \frac{P(c)P(x|c)}{\sum_{c' \in C} P(c')P(x|c')}$

朴素贝叶斯分类器表达式

其中 d 为属性数目, x_i 为 x 在第 i 个属性上的取值

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 $h_{\theta}(x) = \operatorname{argmax}_{c \in C} P(c) \prod_{i=1}^d P(x_i|c)$

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 $P(x_i|c)$ 。

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 $P(x_i|c)$ 可估计为 $P(x_i|c) = \frac{|D_{ci}|}{|D_c|}$ 。[对于离散属性]

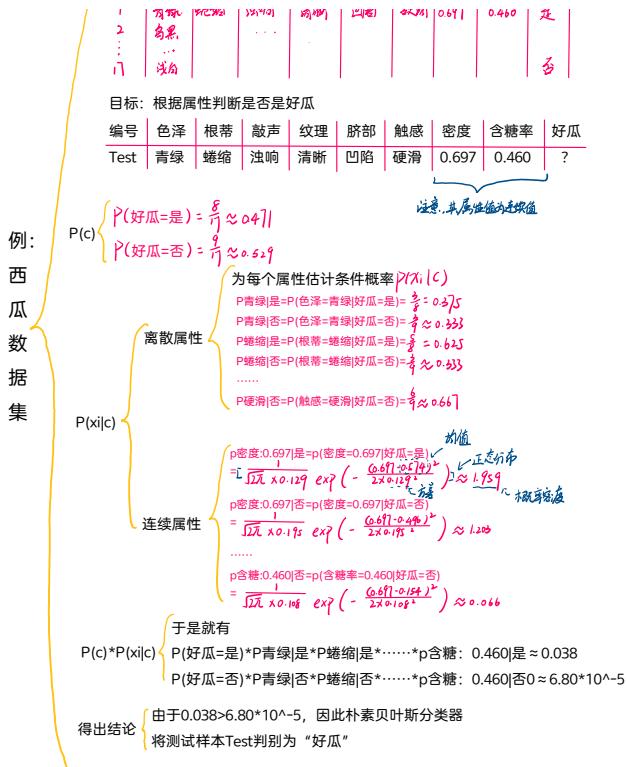
对连续属性可考虑概率密度函数,假定 $p(x_i|c) \sim N(\mu_{ci}, \sigma_{ci}^2)$,其中 μ_{ci} 和 σ_{ci}^2 分别是第 c 类样本在第 i 个属性上取值的均值和方差,则有 $p(x_i|c) = \frac{1}{\sqrt{2\pi}\sigma_{ci}} \exp\left(-\frac{(x_i - \mu_{ci})^2}{2\sigma_{ci}^2}\right)$

西瓜数据集 (1-8号为好瓜, 9-17号为坏瓜)

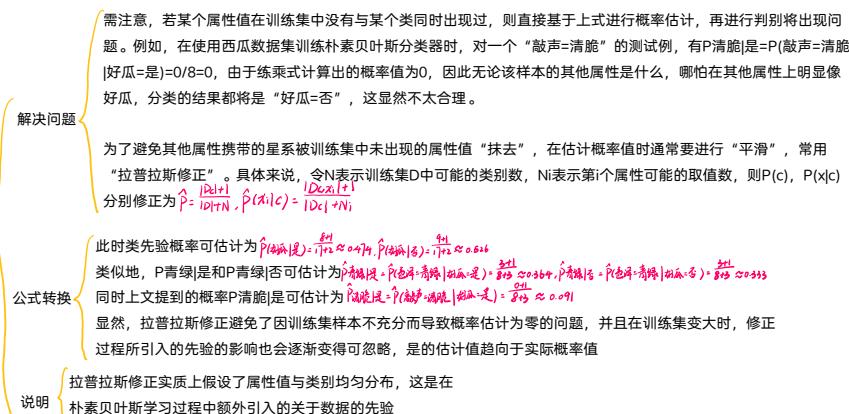
编号	色泽	根蒂	敲声	纹理	脐部	触感	密度	含糖量	好瓜
1	青绿	蜷缩	清脆	清晰	凹陷	硬滑	0.691	0.460	是
2	乌黑		沉闷						
3	...								
4	淡白								否

目标: 根据属性判断是否是好瓜

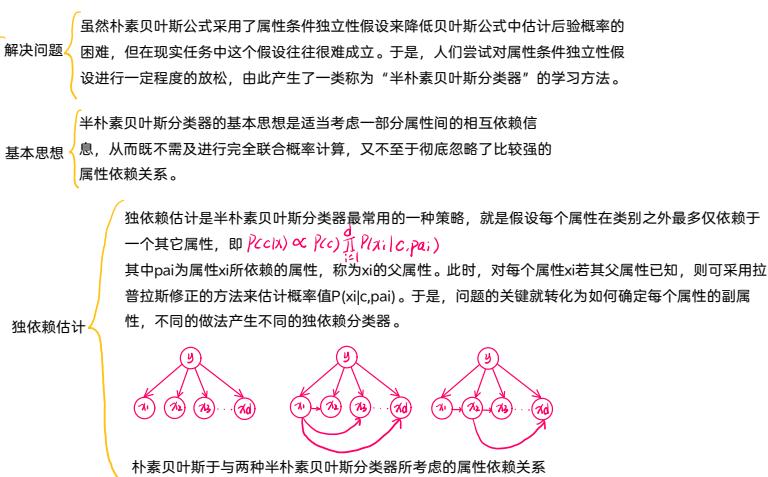
叶斯分类器



拉普拉斯修正



半朴素贝叶斯分类器



继续以上例中的硬币为例，假设现在我们拥有两枚这样特殊的硬币，每次抛掷后得到是花或字的概率可能并不为0.5。记这两枚硬币（编号1和2）抛掷后是花的概率分别为 P_1, P_2 。为了估计这两个概率，现做实验，每次取其中一枚硬币，连掷5下并记录下结果。

硬币 结果 统计 可以很容易地估计出 P_1 和 P_2 :

$$P_1 = (3+1+2)/15 = 0.4 \quad (\text{硬币1出现花的次数}/\text{硬币1总抛掷次数})$$

$$P_2 = (2+3)/10 = 0.5 \quad (\text{硬币2出现花的次数}/\text{硬币2总抛掷次数})$$

在上面的问题中，抹去每轮投掷时使用的硬币标记，如下：

硬币 结果 统计 现在目标仍未变，还是估计 P_1 和 P_2 ，该如何实现？

EM算法

1 字花花花字 2花-3字

在上面的问题中，抹去每轮投掷时使用的硬币标记，如下：

硬币 结果 统计 现在目标仍未变，还是估计P1和P2，该如何实现？

Nan	花花花花字	3花-2字
Nan	字字花花字	2花-3字
Nan	花字字字字	1花-4字
Nan	花字字字花	3花-2字
Nan	字花花字字	2花-3字

加入隐变量z 显然此时我们多了一个隐变量z，可以把它认为是一个5维的向量
(z1, z2, z3, z4, z5)，代表每次投掷时所使用的硬币，比如z1代表第一轮使用的硬币是1还是2。但是这个变量z并不知道，就无法去估计P1, P2，因此我们必须先估计出z，才能进一步估计P1和P2

解决思路 要估计z，我们又得知道P1和P2，这样我们才能用最大似然概率法则去估计z，但是我们的目标就是估计P1和P2，现在又需要让它作为条件，陷入了循环中。

突破问题的关键，就是先随机初始化一个P1和P2，用它来估计z，然后基于z，还是按照最大似然反过来去估计新的P1和P2。

当新的P1和P2与我们的初始值一样时，说明了我们初始的P1和P2很可能是真实的值；当新估计出来的P1与P2与初始化差别很大时，继续用新的P1, P2迭代，直至收敛。

不妨给P1, P2赋初始值：P1=0.2, P2=0.7，看看第一轮最可能是哪个硬币：

- 1) 如果是硬币1，得出“花花花花字”的概率为 $0.2 * 0.2 * 0.2 * 0.8 * 0.8 = 0.00512$
- 2) 如果是硬币2，得出“花花花花字”的概率为 $0.7 * 0.7 * 0.7 * 0.3 * 0.3 = 0.03087$

一次求出其它4轮中的相应概率，做成表格如下：按照极大似然法则：

轮数	若是硬币1	若是硬币2	第1轮中最有可能的是硬币2
1	0.00512	0.03087	第2轮中最有可能的是硬币1
2	0.02048	0.01323	第3轮中最有可能的是硬币1
3	0.08192	0.00567	第4轮中最有可能的是硬币2
4	0.00512	0.03087	第5轮中最有可能的是硬币1
5	0.02048	0.01323	

把上面的值作为z的估计值，然后按照极大似然概率法则来估计新的P1, P2：

$$P1 = (2+1+2)/15 = 0.33, P2 = (3+3)/10 = 0.6$$

现在打开上帝视角，可以知道P1, P2的极大似然估计其实就是上面例子中的0.4和0.5（下文中

将这两个值称为P1, P2的真实值）。对比我们初始化的P1, P2和新估计出的P1, P2：

初始化的P1	估计出的P1	真实的P1	初始化的P2	估计出的P2	真实的P2
0.2	0.33	0.4	0.7	0.6	0.5

不难看出，我们估计的P1, P2相比于它们的初始值，更接近它们的真实值了。可以期待，我们继续按照上面的思路，用估计出的P1和P2再来估计z，再用z来估计新的P1和P2，反复迭代下去，就可以最终得到P1=0.4, P2=0.5。此时无论怎样迭代，P1和P2的值都会保持0.4和0.5不变，这样我们就找到了P1和P2的极大似然估计。

- 注意事项
- 1) 数学证明，新估计出的P1和P2一定会更接近真实的P1和P2
 - 2) 迭代不一定会收敛于真实的P1, P2，这取决于P1, P2的初始值

EM进阶版 我们用极大似然估计法则估计出z的值，然后再用z值按照极大似然概率估计法则估计出新的P1和P2。也就是说，我们使用了一个最可能的z值，而不是所有可能的z值。

如果考虑所有可能的z值，对每一个z值都估计出一个新的P1和P2，将每一个z值概率大小作为权重，将所有新的P1和P2分别加权相加，这样的P1和P2应该会更好。

在本例中，z值共有 $2^5 = 32$ 种，而我们实际上并不需要算这么多次，只需要用期望来简化运算。

轮数	若是硬币1	若是硬币2	利用这张表，我们可以计算出每轮投掷使用硬币1或硬币2的概率。比如第1轮，使用硬币1的概率是： $0.00512 / (0.00512 + 0.03087) = 0.14$ ，使用硬币2的概率是 $1 - 0.14 = 0.86$ 。可以依次计算出其它4轮的概率：
1	0.00512	0.03087	
2	0.02048	0.01323	
3	0.08192	0.00567	
4	0.00512	0.03087	
5	0.02048	0.01323	

轮数	zi=硬币1	zi=硬币2	该表中的左右两列表示期望值。第一行当中，0.86表示从期望的角度看，这轮抛掷使用硬币2的概率是0.86。相比于前面的方法，这种说法更加严谨。这样我们在估计P1或P2时，就可以用上全部的数据，而不是部分的数据。
1	0.14	0.86	
2	0.61	0.39	
3	0.94	0.06	
4	0.14	0.86	
5	0.61	0.39	在这一步中，我们实际上是估计出了z的概率分布，这步被称作E步。

结合下表，我们按照期望最大似然估计概率的法则来估计新的P1和P2：

硬币	结果	统计
Nan	花花花花字	3花-2字
Nan	字字花花字	2花-3字
Nan	花字字字字	1花-4字
Nan	花字字字花	3花-2字
Nan	字花花字字	2花-3字

以P1估计为例，这里我们要做的是在之前推算出来的使用硬币1的概率下，去计算出

现在这种抛掷情况下的可能性，就是计算一下期望看看有多少次等价的硬币1抛出来为

正，多少次抛出来为负。以第一轮为例，硬币1在第一轮中被使用的概率是0.14，现

在出现了3花2字的情况，相当于： $0.14 * 3 = 0.42$ 花， $0.14 * 2 = 0.28$ 字。

依次算出其它4轮，列表如下：

$\text{花} \quad \text{字}$ $\text{总计} \quad 4.22 \quad 7.98$

计算 $P1 = 4.22 / (4.22 + 7.98) = 0.35$

可以看到，改变了z值的估计方法后，新估计出的P1更加接近真实值0.4。原因就是我

们使用了所有抛掷的数据，而不是之前只使用了部分的数据。

这步中，我们根据E步中求出的z的概率分布，根据极大似然概率法则估计P1和P2，被

称作M步。

在上述例子中，我们抹去了硬币的标记，这种“未观测变量”叫做“隐变量”。令X表示已观测变量集，Z表示隐变量集， Θ 表示模型参数。若对 Θ 做极大似然估计，则应最大化对数似然： $LL(\Theta|X, Z) = \ln P(X, Z|\Theta)$

然而由于Z是隐变量，上式无法直接求解，此时我们可通过对Z计算期望，来最大化已观测数据的对数“边际似然”： $LL(\Theta|X) = \ln P(X|\Theta) = \ln \sum_z P(X, Z|\Theta)$

EM算法是常用的估计参数隐变量的利器，其基本思想是：若参数 Θ 已知，则可根据训练数据推出最优隐变量Z的值(E步)；反之，若Z的值已知，则可方便地对参数 Θ 做极大似然估计(M步)。

EM算法

Expectation-Maximization

在“EM初级版”中，其实是以初始值 Θ^0 为起点，对 $LL(\Theta|X)$ ，可迭代执行以下步骤直至收敛：

1) 基于 Θ^t 判断隐变量Z的期望，记为 Z_t

2) 基于已观测变量X和 Z_t 对参数 Θ 做极大似然估计，记为 Θ^{t+1}

EM算法

Expectation-Maximization

在“EM初级版”中，其实就是以初始值 Θ^0 为起点，对 $LL(\Theta|X)$ ，可迭代执行以下步骤直至收敛：

- 1) 基于 Θ^t 判断隐变量Z的期望，记为 Z_t
- 2) 基于已观测变量X和 Z_t 对参数 Θ 做极大似然估计，记为 Θ^{t+1}

进一步，在“EM进阶版”当中，我们不是取Z的期望，而是基于 Θ^t 计算隐变量Z的概率分布 $P(Z|X,\Theta^t)$ ，则EM算法的两个步骤是：

- 1) 以当前参数 Θ^t 推断隐变量分布 $P(Z|X,\Theta^t)$ ，并计算对数似然 $LL(\Theta|X,Z)$ 关于Z的期望： $Q(\Theta|\Theta^{t+1}) = E_{Z|X,\Theta^t} LL(\Theta|X,Z)$
- 2) 寻找参数最大化期望似然，即 $\Theta^{t+1} = \underset{\Theta}{\operatorname{argmax}} Q(\Theta|\Theta^t)$

EM算法使用两个步骤交替计算，新得到的参数值重新被用于E步，直至收敛到局部最优解。