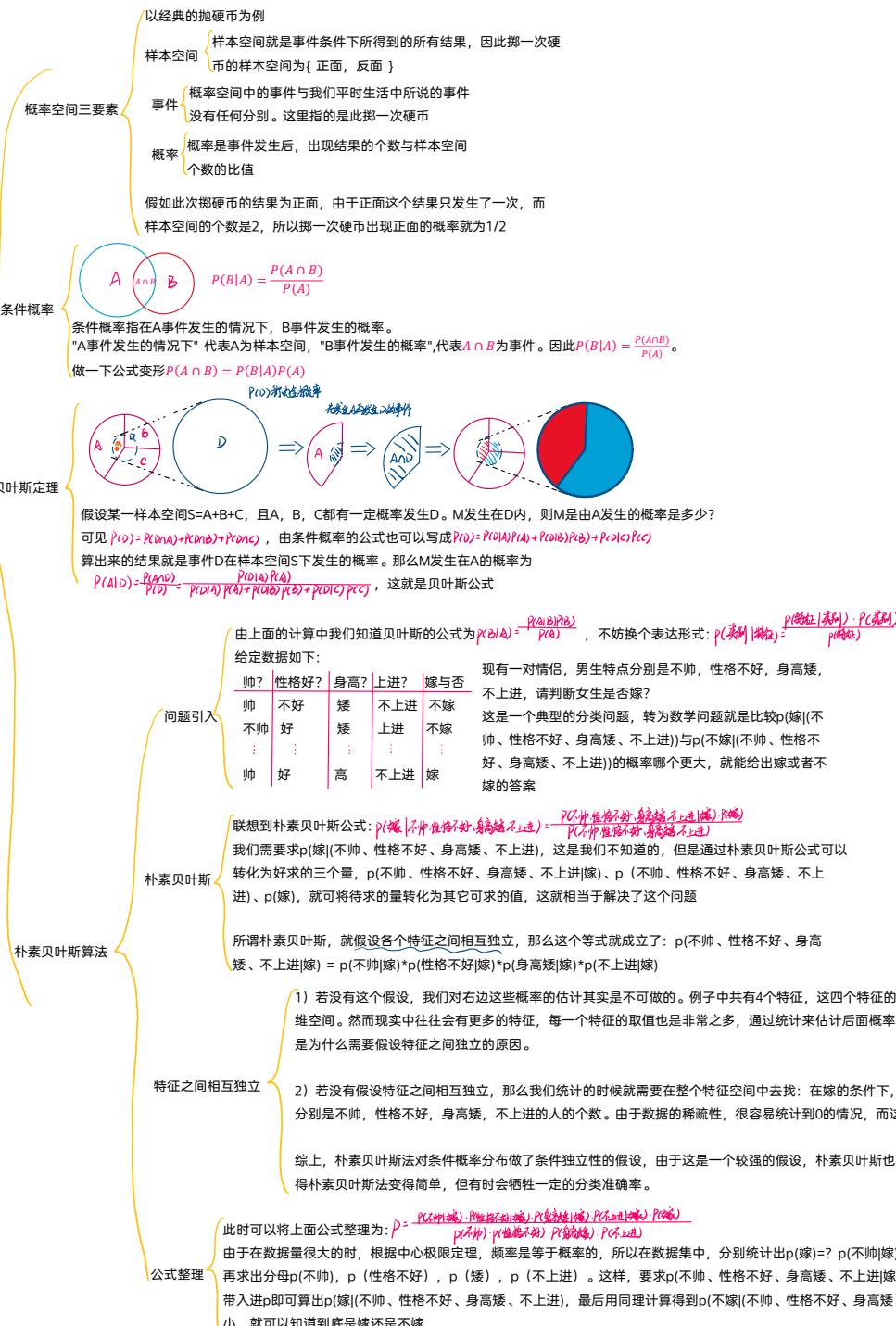
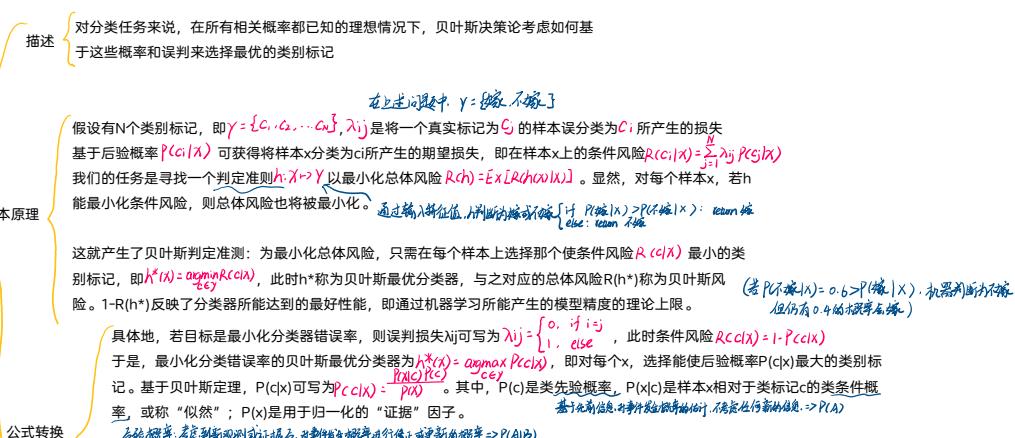


# 贝叶斯分类器

## 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

于是，最小化分类错误率的贝叶斯最优分类器为  $\hat{y}(x) = \arg\max_{c \in \{A, B\}} P(c|x)$ ，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。  
似然概率率，若要判断测试时证据后，对事件发生概率进行修正的频率  $\Rightarrow P(A|B)$

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。这种通过事实反过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。通过比对设置不同参数下得到的概率，概率越大的可能是正确的。  
用 0.5 作为参数，概率为  $C(10)_0.5^6(1-0.5)^4 \approx 0.21$ ，再用 0.6 作参数： $C(10)_0.6^6(1-0.6)^4 \approx 0.25$ ，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。

解决问题

设置硬币的参数为  $\theta$ ，可以得到似然函数：

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。通过投入更多的实验结果，参数会越来越明确。

极大似然估计

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  $L(\theta) = f(x_1|\theta) f(x_2|\theta) \dots f(x_n|\theta)$ ，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。

### 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  $\theta_c$  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： $P(D_c|\theta_c) = \prod_{i \in D_c} P(x_i|\theta_c)$   
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。为了防止上式练习操作造成下溢，通常使用对数似然： $LL(\theta_c) = \log P(D_c|\theta_c) = \sum_{i \in D_c} \log P(x_i|\theta_c)$   
此时参数  $\theta_c$  的极大似然估计  $\hat{\theta}_c = \arg\max_{\theta_c} LL(\theta_c)$ 。  
例如，在连续属性情形下，假设概率密度函数  $P(x|\theta_c) \sim N(\mu_{\theta_c}, \sigma_{\theta_c}^2)$ ，则参数  $\mu_{\theta_c}$  和  $\sigma_{\theta_c}^2$  的极大似然估计为  $\hat{\mu}_{\theta_c} = \frac{1}{|D_c|} \sum_{x \in D_c} x$ ,  $\hat{\sigma}_{\theta_c}^2 = \frac{1}{|D_c|} \sum_{x \in D_c} (x - \hat{\mu}_{\theta_c})^2$   
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  $(x - \hat{\mu}_{\theta_c})(x - \hat{\mu}_{\theta_c})^T$  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

朴素贝叶斯分类器表达式

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  $P(x_i|c)$ 。  
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  $P(x_i|c) = \frac{|D_c \cap x_i|}{|D_c|}$ 。  
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  $\mu_{ci}$  和  $\sigma_{ci}^2$  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  $P(x_i|c) = \frac{1}{\sqrt{2\pi}\sigma_{ci}} \exp\left(-\frac{(x_i - \mu_{ci})^2}{2\sigma_{ci}^2}\right)$ 。

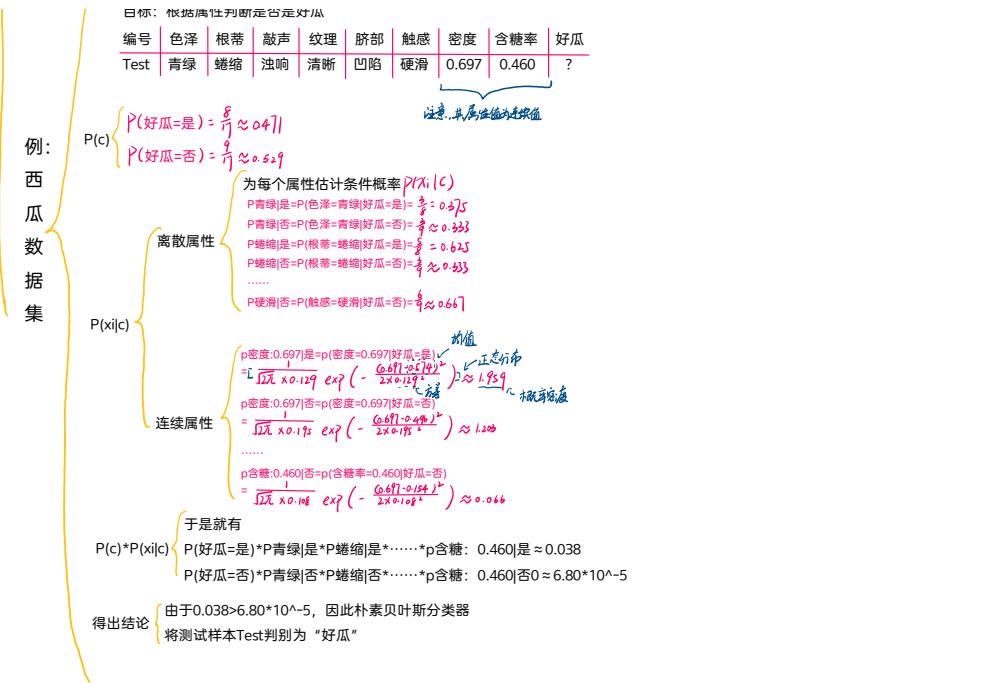
西瓜数据集 (1-8号为好瓜, 9-17号为坏瓜)

| 编号 | 色泽 | 根蒂 | 敲声 | 纹理 | 脐部 | 触感 | 密度    | 含糖率   | 好瓜 |
|----|----|----|----|----|----|----|-------|-------|----|
| 1  | 青绿 | 蜷缩 | 浊响 | 清晰 | 凹陷 | 硬滑 | 0.691 | 0.460 | 是  |
| 2  | 乌黑 |    |    |    |    |    |       |       |    |
| 3  |    |    |    |    |    |    |       |       |    |
| 4  | 淡白 |    |    |    |    |    |       |       |    |
| 5  |    |    |    |    |    |    |       |       |    |
| 6  |    |    |    |    |    |    |       |       |    |
| 7  |    |    |    |    |    |    |       |       |    |
| 8  |    |    |    |    |    |    |       |       |    |
| 9  |    |    |    |    |    |    |       |       |    |
| 10 |    |    |    |    |    |    |       |       |    |
| 11 |    |    |    |    |    |    |       |       |    |
| 12 |    |    |    |    |    |    |       |       |    |
| 13 |    |    |    |    |    |    |       |       |    |
| 14 |    |    |    |    |    |    |       |       |    |
| 15 |    |    |    |    |    |    |       |       |    |
| 16 |    |    |    |    |    |    |       |       |    |
| 17 |    |    |    |    |    |    |       |       |    |

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 | ?  |

# 斯分类器



# 拉普拉斯修正

需要注意，若某个属性值在训练集中没有与某个类同时出现过，则直接基于上式进行概率估计，再进行判别将出现问题。例如，在使用西瓜数据集训练朴素贝叶斯分类器时，对一个“敲声=清脆”的测试例，有 $P(\text{清脆}|\text{是}) = P(\text{敲声=清脆}|\text{好瓜=是}) = 0/8 = 0$ ，由于训练乘式计算出的概率值为0，因此无论该样本的其他属性是什么，哪怕在其他属性上明显像好瓜，分类的结果都将是“好瓜=否”，这显然不太合理。

解决问题  
为了避免其他属性携带的星系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”。具体来说，令 $N$ 表示训练集D中可能的类别数， $N_i$ 表示第*i*个属性可能的取值数，则 $P(c)$ ,  $P(x_i|c)$  分别修正为  $\hat{P} = \frac{|D|+1}{|D|+N}$ ,  $\hat{P}(x_i|c) = \frac{|D_{ci}|+1}{|D_c|+N_i}$

公式转换  
此时类先验概率可估计为  $\hat{P}(\text{好瓜=是}) = \frac{1}{17} \approx 0.059$ ,  $\hat{P}(\text{好瓜=否}) = \frac{16}{17} \approx 0.941$   
类似地， $P(\text{青绿}|\text{是})$  和  $P(\text{青绿}|\text{否})$  可估计为  $\hat{P}(\text{青绿}|\text{是}) = \frac{3+1}{17+1} \approx 0.235$ ,  $\hat{P}(\text{青绿}|\text{否}) = \frac{3+1}{17+1} \approx 0.235$   
同时上文提到的概率 $P(\text{清脆}|\text{是})$  可估计为  $\hat{P}(\text{清脆}|\text{是}) = \frac{3+1}{17+1} \approx 0.235$   
显然，拉普拉斯修正避免了因训练集样本不充分而导致概率估计为零的问题，并且在训练集变大时，修正过程所引入的先验的影响也会逐渐变得可忽略，是的估计值趋向于实际概率值

说明  
拉普拉斯修正实质上假设了属性值与类别均匀分布，这是在朴素贝叶斯学习过程中额外引入的关于数据的先验

# 半朴素贝叶斯分类器

解决问题  
虽然朴素贝叶斯公式采用了属性条件独立性假设来降低贝叶斯公式中估计后验概率的困难，但在现实任务中这个假设往往很难成立。于是，人们尝试对属性条件独立性假设进行一定程度的放松，由此产生了一类称为“半朴素贝叶斯分类器”的学习方法。

基本思想  
半朴素贝叶斯分类器的基本思想是适当考虑一部分属性间的相互依赖信息，从而既不需进行完全联合概率计算，又不至于彻底忽略了比较强的属性依赖关系。

独依赖估计  
独依赖估计是半朴素贝叶斯分类器最常用的一种策略，就是假设每个属性在类别之外最多仅依赖于一个其它属性，即  $P(c|x) \propto P(c) \prod_{i=1}^d P(x_i|c, pa_i)$   
其中 $pa_i$ 为属性 $x_i$ 所依赖的属性，称为 $x_i$ 的父属性。此时，对每个属性 $x_i$ 若其父属性已知，则可采用拉普拉斯修正的方法来估计概率值 $P(x_i|c, pa_i)$ 。于是，问题的关键就转化为如何确定每个属性的父属性，不同的做法产生不同的独依赖分类器。



朴素贝叶斯与两种半朴素贝叶斯分类器所考虑的属性依赖关系

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

可以很容易地估计出 $P_1$ 和 $P_2$ ： $P_1 = (3+1)/15 = 0.4$  (硬币1出现花的次数/硬币1总抛掷次数)

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

$P_2 = (2+3)/10 = 0.5$

| 硬币  | 结果    | 统计    |
|-----|-------|-------|
| Nan | 花字字字字 | 1花-4字 |

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

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

$(z_1, z_2, z_3, z_4, z_5)$ ，代表每次投掷时所使用的硬币，比如 $z_1$ 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 $P_1, P_2$ ，因此我们必须先

估计出 $z$ ，才能进一步估计 $P_1$ 和 $P_2$

# EM算法

加入隐变量z  
 Nan 子子化化子 2化-3子 既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  
 Nan 花字字字字 1花-4字  $(z_1, z_2, z_3, z_4, z_5)$ , 代表每次投掷时所使用的硬币，比如 $z_1$ 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 $P_1, P_2$ ，因此我们必须先估计出z，才能进一步估计 $P_1$ 和 $P_2$

解决思路  
 要估计z，我们又得知道 $P_1$ 和 $P_2$ ，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 $P_1$ 和 $P_2$ ，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 $P_1$ 和 $P_2$ ，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 $P_1$ 和 $P_2$ 。

当新的 $P_1$ 和 $P_2$ 与我们的初始值一样时，说明了我们初始的 $P_1$ 和 $P_2$ 很可能是真实的值；当新估计出来的 $P_1$ 与 $P_2$ 与初始化差别很大时，继续用新的 $P_1, P_2$ 迭代，直至收敛。

不妨给 $P_1, P_2$ 赋初始值： $P_1=0.2, P_2=0.7$ ，看看第一轮最可能是哪个硬币：  
 1) 如果是硬币1，得出“花花字花字”的概率为 $0.2*0.2*0.2*0.8*0.8=0.0512$   
 2) 如果是硬币2，得出“花花字花字”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 $P_1, P_2$ ：

$$P_1=(2+1+2)/15=0.33, P_2=(3+3)/10=0.6$$

现在打开上帝视角，可以知道 $P_1, P_2$ 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 $P_1, P_2$ 的真实值）。对比我们初始化的 $P_1, P_2$ 和新估计出的 $P_1, P_2$ ：

| 初始化的 $P_1$ | 估计出的 $P_1$ | 真实的 $P_1$ | 初始化的 $P_2$ | 估计出的 $P_2$ | 真实的 $P_2$ |
|------------|------------|-----------|------------|------------|-----------|
| 0.2        | 0.33       | 0.4       | 0.7        | 0.6        | 0.5       |

不难看出，我们估计的 $P_1, P_2$ 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 $P_1$ 和 $P_2$ 再来估计z，再用z来估计新的 $P_1$ 和 $P_2$ ，反复迭代下去，就可以最终得到 $P_1=0.4, P_2=0.5$ 。此时无论怎样迭代， $P_1$ 和 $P_2$ 的值都会保持0.4和0.5不变，这样我们就找到了 $P_1$ 和 $P_2$ 的极大似然估计。

注意事项  
 1) 数学证明，新估计出的 $P_1$ 和 $P_2$ 一定会更接近真实的 $P_1$ 和 $P_2$   
 2) 迭代不一定会收敛于真实的 $P_1, P_2$ ，这取决于 $P_1, P_2$ 的初始值

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 $P_1$ 和 $P_2$ 。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 $P_1$ 和 $P_2$ ，将每一个z值概率大小作为权重，将所有新的 $P_1$ 和 $P_2$ 分别加权相加，这样的 $P_1$ 和 $P_2$ 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 |   |

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 $P_1$ 或 $P_2$ 时，就可以用上全部的数据，而不是部分的数据。

在这一步骤中，我们实际上是估计出了z的概率分布，这步被称作E步。

EM进阶版

| 轮数 | zi=硬币1 | zi=硬币2 |
|----|--------|--------|
| 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   |

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 $P_1$ 和 $P_2$ ：

| 硬币  | 结果    | 统计    | 以 $P_1$ 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。 |
|-----|-------|-------|---|
| Nan | 花花字字字 | 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 | 花字字花花 | 3花-2字 | 依次算出其它四轮，列表如下： <u>待续</u> <small>如图所示，此处省略了其余四轮的结果</small>   |
| Nan | 字花花字字 | 2花-3字 |   |

| 轮数 | 花    | 字    | 计算 $P_1=4.22/(4.22+7.98)=0.35$  |
|----|------|------|---|
| 1  | 0.42 | 0.28 | 可以看到，改变了z值的估计方法后，新估计出的 $P_1$ 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，而不再是之前只使用了部分的数据。 |
| 2  | 1.22 | 1.83 |   |
| 3  | 0.94 | 3.76 |   |
| 4  | 0.42 | 0.28 | 这一步中，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 $P_1$ 和 $P_2$ ，被称作M步。                  |
| 5  | 1.22 | 1.83 |   |
| 总计 | 4.22 | 7.98 |   |

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， $\Theta$ 表示模型参数。若对 $\Theta$ 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 $\Theta$ 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 $\Theta$ 做极大似然估计(M步)。

EM算法

Expectation – Maximization

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

1) 基于 $\Theta^t$ 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 $\Theta$ 做极大似然估计，记为 $\Theta^{t+1}$

进一步，在“EM进阶版”当中，我们不是取Z的期望，而是基于 $\Theta^t$ 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 $\Theta^t$ 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： $Q(\Theta|\Theta^t) = E_{Z|X, \Theta^t} LL(\Theta|X, Z)$

2) 寻找参数最大化期望似然，即  $\Theta^{t+1} = \arg \max_{\Theta} Q(\Theta|\Theta^t)$

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。