

一、问题描述

一个袋子中三种硬币的混合比例为： s_1, s_2 与 $1 - s_1 - s_2$ ($0 \leq s_i \leq 1$)，三种硬币掷出正面的概率分别为： p, q, r 。

(1)自己指定系数 s_1, s_2, p, q, r ，生成 N 个投掷硬币的结果（由01构成的序列，其中1为正面，0为反面）。

(2)利用EM算法来对参数进行估计并与预先假定的参数进行比较。

二、实验原理

该问题中，参数为 $\theta = (s_1, s_2, p, q, r)$ ，依次为选择第一个硬币的概率，选择第二个硬币的概率，第一个硬币，第二个硬币，第三个硬币抛掷时出现正面的概率。序列 $x = \{x_1, x_2, \dots, x_N\}$ ， x_i 为单次投硬币的结果，1为正面，0为反面。则有：

$$p(x_i|\theta) = s_1 \cdot p^{x_i} \cdot (1-p)^{1-x_i} + s_2 \cdot q^{x_i} \cdot (1-q)^{1-x_i} + (1-s_1-s_2) \cdot r^{x_i} \cdot (1-r)^{1-x_i}$$

$$\begin{aligned} p(x|\theta) &= \prod_{i=1}^N p(x_i|\theta) \\ &= \prod_{i=1}^N [s_1 \cdot p^{x_i} \cdot (1-p)^{1-x_i} + s_2 \cdot q^{x_i} \cdot (1-q)^{1-x_i} + (1-s_1-s_2) \cdot r^{x_i} \cdot (1-r)^{1-x_i}] \end{aligned}$$

设 x_i 来自硬币 A, B, C 的概率分别为 u_{1i}, u_{2i}, u_{3i} ，即：

$$\begin{aligned} u_{1i} &= \frac{p(x_i, A|\theta)}{p(x_i|\theta)} \\ &= \frac{1 \cdot p^{x_i} \cdot (1-p)^{1-x_i}}{s_1 \cdot p^{x_i} \cdot (1-p)^{1-x_i} + s_2 \cdot q^{x_i} \cdot (1-q)^{1-x_i} + (1-s_1-s_2) \cdot r^{x_i} \cdot (1-r)^{1-x_i}} \\ u_{2i} &= \frac{p(x_i, B|\theta)}{p(x_i|\theta)} \\ &= \frac{s_2 \cdot q^{x_i} \cdot (1-q)^{1-x_i}}{s_1 \cdot p^{x_i} \cdot (1-p)^{1-x_i} + s_2 \cdot q^{x_i} \cdot (1-q)^{1-x_i} + (1-s_1-s_2) \cdot r^{x_i} \cdot (1-r)^{1-x_i}} \\ u_{3i} &= \frac{p(x_i, C|\theta)}{p(x_i|\theta)} \\ &= \frac{(1-s_1-s_2) \cdot r^{x_i} \cdot (1-r)^{1-x_i}}{s_1 \cdot p^{x_i} \cdot (1-p)^{1-x_i} + s_2 \cdot q^{x_i} \cdot (1-q)^{1-x_i} + (1-s_1-s_2) \cdot r^{x_i} \cdot (1-r)^{1-x_i}} \end{aligned}$$

则有：

$$\begin{aligned} p(x|\theta) &= \prod_{i=1}^N \{ [s_1 \cdot p^{x_i} \cdot (1-p)^{1-x_i}]^{u_{1i}} + [s_2 \cdot q^{x_i} \cdot (1-q)^{1-x_i}]^{u_{2i}} \\ &\quad + [(1-s_1-s_2) \cdot r^{x_i} \cdot (1-r)^{1-x_i}]^{u_{3i}} \} \end{aligned}$$

实验 4	s1	s2	s3	p	q	r	正面概率
真实值	0.2000	0.3000	0.5000	0.3000	0.4000	0.8000	0.5800
迭代初值	0.4000	0.3000	0.3000	0.5000	0.4000	0.3000	0.4100
迭代终值	0.4244	0.2980	0.2776	0.6597	0.5637	0.4537	0.5739
生成的数据中正面个数占比							0.5739

实验 5	s1	s2	s3	p	q	r	正面概率
真实值	0.4000	0.2000	0.4000	0.1000	0.3000	0.5000	0.3000
迭代初值	0.1000	0.2000	0.7000	0.3000	0.4000	0.5000	0.4600
迭代终值	0.1106	0.2080	0.6814	0.1741	0.2470	0.3297	0.2953
生成的数据中正面个数占比							0.2953

实验 6	s1	s2	s3	p	q	r	正面概率
真实值	0.4000	0.2000	0.4000	0.1000	0.3000	0.5000	0.3000
迭代初值	0.4000	0.3000	0.3000	0.5000	0.4000	0.3000	0.4100
迭代终值	0.3839	0.3013	0.3148	0.3835	0.2931	0.2105	0.3018
生成的数据中正面个数占比							0.3018

四、结果分析

由于算法的输入只有一个由01构成的序列，而实际上信息量只有序列中0或1的占比。因此，迭代终值确定的综合正面概率 $s_1p + s_2q + s_3r$ ，能够收敛到生成的序列中正面个数占比的数值。

但在对于 s_1, s_2, p, q, r 这5个参数的估计方面，不论真实值与初值如何选择，EM算法效果都不是很好。因为对于这个问题，能使正面概率与真实值确定的正面概率相同的参数有很多组，而根据初值的选择，EM算法会收敛到其中一组参数上。亦即，EM算法确定的只是局部最优解。

五、附录（MATLAB 程序）

```
clear;clc;
%% 设置参数
[S1,S2]=deal(0.4,0.2);
[P,Q,R]=deal(0.1,0.3,0.5);
S3=1-S1-S2;%S3=0.2;
S=[S1,S2,S3,P,Q,R,S1*P+S2*Q+S3*R];
% 由初始参数确定的混合比例，正面概率，以及综合正面概率
%% 生成数据
N=10000; % 数据个数
X=zeros(N,1);
PI=[P,Q,R];
for i=1:N
```

```

n=randsrc(1,1,[1,2,3;S1,S2,S3]); % 选择一枚硬币
X(i)=randsrc(1,1,[1,0;PI(n),1-PI(n)]); % 投掷这枚硬币
end
%% 设置初值
[s1,s2]=deal(0.4,0.3);
[p,q,r]=deal(0.5,0.4,0.3);
s3=1-s1-s2;
s0=[s1,s2,s3,p,q,r,s1*p+s2*q+s3*r];
% 由迭代初值确定的混合比例，正面概率，以及综合正面概率
%% EM
u1=zeros(N,1);
u2=zeros(N,1);
u3=zeros(N,1);
%  $p(x=1)=s1*p+s2*q+s3*r$ 
%  $p(x=0)=s1*(1-p)+s2*(1-q)+s3*(1-r)$ 
%  $p(x)=s1*p^x*(1-p)^(1-x)+s2*q^x*(1-q)^(1-x)+s3*r^x*(1-r)^(1-x)$ 
M=10; % 迭代步数
s=zeros(M,length(s0));
for j=1:M
    % E-Step
    for i=1:N
        x=X(i);
        u1(i)=(s1*p^x*(1-p)^(1-x))/(s1*p^x*(1-p)^(1-x)+s2*q^x*(1-q)^(1-x)+s3*r^x*(1-r)^(1-x));
        u2(i)=(s2*q^x*(1-q)^(1-x))/(s1*p^x*(1-p)^(1-x)+s2*q^x*(1-q)^(1-x)+s3*r^x*(1-r)^(1-x));
        u3(i)=(s3*r^x*(1-r)^(1-x))/(s1*p^x*(1-p)^(1-x)+s2*q^x*(1-q)^(1-x)+s3*r^x*(1-r)^(1-x));
    end
    % M-Step
    s1=sum(u1)/N;
    s2=sum(u2)/N;
    s3=sum(u3)/N;
    p=sum(u1.*X)/sum(u1);
    q=sum(u2.*X)/sum(u2);
    r=sum(u3.*X)/sum(u3);
    s(j,:)= [s1,s2,s3,p,q,r,s1*p+s2*q+s3*r];
    % 迭代过程中得到的混合比例，正面概率，以及综合正面概率
end
%% Result
disp("由初始参数确定的混合比例，正面概率，以及综合正面概率"),disp(S)
disp("由迭代初值确定的混合比例，正面概率，以及综合正面概率"),disp(s0)
% disp("迭代过程中得到的混合比例，正面概率，以及综合正面概率"),disp(s)
disp("由迭代终值确定的混合比例，正面概率，以及综合正面概率"),disp(s(end,:))

```

```
disp("生成的数据中正面占比，根据迭代终值确定的正面概率")  
disp([sum(X)/N,s1*p+s2*q+s3*r])
```