

1 对有限态混合总体, 以  $k=3$  为例, 即

$$X \stackrel{d}{=} p_1 N(\mu_1, \sigma^2) + p_2 N(\mu_2, \sigma^2) + p_3 N(\mu_3, \sigma^2)$$

其中  $p_1 < p_2 < p_3$ . 生成  $X$  的一组样本,

(1) 使用 EM 算法估计参数  $(p_1, p_2, \mu_1, \mu_2, \mu_3, \sigma^2)$ .

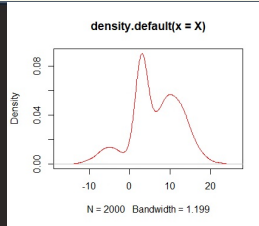
(2) 试用方差估计的某个方法, 估计(1)中各估计的标准差的估计.

分析过程: 这是一个正态混合参数估计问题,  $\theta = (p, \mu, \sigma^2)$  为待估参数, 其中

$p = (p_1, p_2, p_3)$ ,  $\mu = (\mu_1, \mu_2, \mu_3)$ . 假设  $p_1 = 0.1$ ,  $p_2 = 0.3$ ,  $p_3 = 0.6$ ,  $\mu_1 = 5$ ,  $\sigma_1^2 = 3$ ,  $\mu_2 = 3$ ,  $\sigma_2^2 = 1$

$\mu_3 = 10$ ,  $\sigma_3^2 = 4$ . 生成  $X$  的一组样本  $X = (x_1, \dots, x_N)$ ,  $N = 2000$ , 初始化参数, 利用 EM 算法迭代求解

```
1 set.seed(101)
2 N<-1000      ##生成X的一组样本
3 n1<-rbinom(1,size=N,prob=0.1)
4 n2<-rbinom(1,size=N,prob=0.3)
5 n3<-N-n1-n2
6 X1<-rnorm(n1,-5,3)
7 X2<-rnorm(n2,3,1)
8 X3<-rnorm(n3,10,4)
9 X0<-c(X1,X2,X3)
10 X0<-X[sample(1:N)]
11 plot(density(X),col="red")
12
```



```
13 EM<-function(X){
14   ##initial value规定初始值##
15   p10<-0.1
16   p20<-0.3
17   p30<-0.6
18   mu10<-5.2
19   sigma10<-3.1
20   mu20<-3.3
21   sigma20<-1.2
22   mu30<-9.3
23   sigma30<-4.4
24   para<-c(p10,p20,p30,mu10,sigma10,mu20,sigma20,mu30,sigma30)
25   tol<-1e-8
26   para.old<-para+1
27   ##
28   for(j in 1:1000){
29     #count<-count+1
30     vp10<-(para[1]*dnorm(X,para[4],para[5]))/(para[1]*dnorm(X,para[4],para[5])
31       +para[2]*dnorm(X,para[6],para[7])
32       +para[3]*dnorm(X,para[8],para[9]))
33     vp20<-(para[2]*dnorm(X,para[6],para[7]))/(para[1]*dnorm(X,para[4],para[5])
34       +para[2]*dnorm(X,para[6],para[7])
35       +para[3]*dnorm(X,para[8],para[9]))
36     vp30<-(para[3]*dnorm(X,para[8],para[9]))/(para[1]*dnorm(X,para[4],para[5])
37       +para[2]*dnorm(X,para[6],para[7])
38       +para[3]*dnorm(X,para[8],para[9]))
39
40     phi1<-sum(vp10)
41     phi2<-sum(vp20)
42     phi3<-sum(vp30)
43     p11<-phi1/N
44     p21<-phi2/N
45     p31<-phi3/N
46     phix1<-sum(X*vp10)
47     mu11<-phix1/phi1
48     phix2<-sum(X*vp20)
49     mu21<-phix2/phi2
50     phix3<-sum(X*vp30)
51     mu31<-phix3/phi3
52     phixmu1<-sum((X-mu11)^2*vp10)
53     sigma11<-sqrt(phixmu1/phi1)
54     phixmu2<-sum((X-mu21)^2*vp20)
55     sigma21<-sqrt(phixmu2/phi2)
56     phixmu3<-sum((X-mu31)^2*vp30)
57     sigma31<-sqrt(phixmu3/phi3)
58
59     para<-c(p11,p21,p31,mu11,sigma11,mu21,sigma21,mu31,sigma31)
60
61     if (sqrt(sum((para-para.old)^2))/sqrt(sum(para.old^2))<tol)
62       {break}
63     para.old<-para
64   }
65   return(list(estimate=para,iter=j,tol=tol))
66 }
```

```
67 #print(list(estimate=para,iter=j,tol=tol))
68
69 B <- 50
70 n <- N
71 theta.b <- matrix(0,nrow=B,ncol=9)
72 theta.b[1,]<-EM(X0)$estimate
73
```

```
69 B <- 50
70 n <- N
71 theta.b <- matrix(0,nrow=B,ncol=9)
72 theta.b[1,]<-EM(X0)$estimate
73 for (b in 2:B) {
74   i <- sample(1:n, size = n, replace = TRUE)
75   Xnew <- X0[i]
76   theta.b[b,]<- EM(Xnew)$estimate
77 }
78 print(theta.b)
79 sqrt(diag(cov(theta.b)))      #计算样本协方差矩阵
```

(1).

```
> print(list(estimate=para,iter=j,tol=tol))
$estimate
[1] 0.1086194 0.2966098 0.5947708
[4] -4.9767161 3.0172355 2.9069963
[7] 0.9630215 10.2146407 3.9936090

$iter
[1] 92

$tol
[1] 1e-08
```

结论: 图中 estimate 数据即为各参数的估计

(2).

```
> sqrt(diag(cov(theta.b)))
[1] 0.00977055 0.02499198 0.02629029 0.39873862
[5] 0.33572892 0.08494352 0.07279765 0.24781455
[9] 0.15521946
```

10.2

2. 假设一种灯泡的寿命服从某个参数为  $\theta$  的指数分布. 为了估计  $\theta$ , 我们测试了  $n$  个灯泡直至他们失效, 记失效的时间是  $y_1, \dots, y_n$ ; 另外还有一组独立的实验测试了  $m$  个灯泡, 但是单个灯泡的失效时间没有被记录, 只记录下来在时刻  $t$  时共有  $r$  个灯泡失效. 试使用 EM 算法求  $\theta$  的极大似然估计的迭代公式.

提示: 可观测数据:  $y_1, \dots, y_n, r$

完全数据:  $y_1, \dots, y_n, z_1, \dots, z_m$ , 其中  $z_1, \dots, z_m$  为另一组实验中  $m$  个灯泡的失效时间.

$r$  个在  $t$  时失效

$m-r$  个在  $t$  时未失效.

如上述自己设定与数据生成数据, 并给出用 EM 算法求得的  $\hat{\theta}_{MLE}$ .

分析: (1).  $E(X) = \frac{1}{\theta}$ ,  $f(x) = \theta e^{-\theta x}$   $x > 0$ ,  $F(x) = \begin{cases} 1 - e^{-\theta x}, & x > 0 \\ 0, & x \leq 0 \end{cases} \rightarrow 1 - e^{-\theta x} = 1 - \frac{1}{e^{\theta x}} = \frac{e^{\theta x} - 1}{e^{\theta x}}$

$$L(\theta|X) = \theta^{n+m} e^{-\theta (\sum_{i=1}^n y_i + \sum_{i=1}^m z_i)}$$

其中,  $x = (n, z)$ .

$$\Rightarrow \lg L(\theta|X) = (n+m) \lg \theta - \theta (n\bar{y} + \sum_{i=1}^m z_i)$$

① E步:  $Q(\theta|\theta^{(t)}) = E_x \{ \lg L(\theta|X) | y, \theta_t \}$

$$= (n+m) \lg \theta - \theta [n\bar{y} + \underbrace{(m-r)(t + \frac{1}{\theta^{(t)}})}_{\text{估计失效时间}} + r(\frac{1}{\theta^{(t)}} - t_{ht})] \quad (*)$$

② M步: 由  $\frac{\partial Q}{\partial \theta} = 0 \Rightarrow \frac{n+m}{\theta} - n\bar{y} + (m-r)(t + \frac{1}{\theta^{(t)}}) + r(\frac{1}{\theta^{(t)}} - t_{ht}) = 0$

$$\Rightarrow \theta^{(t+1)} = (n+m) \cdot [n\bar{y} + (m-r)(t + \frac{1}{\theta^{(t)}}) + r(\frac{1}{\theta^{(t)}} - t_{ht})]$$

得到 EM 再次迭代公式

(\*) 式:  $E(\sum_{i=1}^m z_i | r, \theta^{(t)}, x) = \sum_{i=1}^r E(z_i | r, \theta^{(t)}, x) + \sum_{i=r+1}^m E(z_i | z_i \geq t, x, \theta^{(t)})$

$$= r \cdot \frac{\int_0^t z f(z, \theta^{(t)}) dz}{\int_0^t f(z, \theta^{(t)}) dz} + (m-r) \cdot \frac{\int_t^{+\infty} z f(z, \theta^{(t)}) dz}{\int_t^{+\infty} f(z, \theta^{(t)}) dz}$$

A =  $\frac{\int_0^t z \theta^{(t)} e^{-\theta^{(t)} z} dz}{\int_0^t \theta^{(t)} e^{-\theta^{(t)} z} dz} = \frac{\int_0^t z de^{-\theta^{(t)} z}}{\int_0^t de^{-\theta^{(t)} z} \cdot (-\frac{1}{\theta^{(t)}})} = \frac{te^{-\theta^{(t)} t} - \int_t^{+\infty} e^{-\theta^{(t)} z} dz}{e^{-\theta^{(t)} t} - 1}$

$\left( \int_0^t e^{-\theta^{(t)} z} dz = -\frac{1}{\theta^{(t)}} (e^{-\theta^{(t)} t} - 1) \right) = \frac{te^{-\theta^{(t)} t}}{e^{-\theta^{(t)} t} - 1} + \frac{1}{\theta^{(t)}}$

B =  $\frac{\int_t^{+\infty} z \theta^{(t)} e^{-\theta^{(t)} z} dz}{\int_t^{+\infty} \theta^{(t)} e^{-\theta^{(t)} z} dz} = \frac{\int_t^{+\infty} z de^{-\theta^{(t)} z}}{\int_t^{+\infty} de^{-\theta^{(t)} z}} = \frac{ze^{-\theta^{(t)} z} \Big|_t^{+\infty} - \int_t^{+\infty} e^{-\theta^{(t)} z} dz}{e^{-\theta^{(t)} t}}$

$$= \frac{0 - te^{-\theta^{(t)} t} + \frac{1}{\theta^{(t)}} (0 - e^{-\theta^{(t)} t})}{0 - e^{-\theta^{(t)} t}} = t + \frac{1}{\theta^{(t)}}$$

2) R code :

```
85 ##2
86 set.seed(101)
87 N<-1000
88 n<-600
89 m<-400
90 t<-5
91 r<-160
92 X1<-rexp(n,1/10) #生成样本 #设定样本服从theta=10的指数分布
93 X2<-rexp(m,1/10)
94 X0<-c(X1,X2)
95 X0<-X0[sample(1:N)]
96 plot(density(X0))
97
98 theta0<-9
99 para<-theta0
100 tol<-1e-8
101 para.old<-para+1
102 ###
103 for(j in 1:1000){ #进行迭代
104   ht<-exp(-t*para)/(1-exp(-t*para))
105   theta1<-(m+n)/(sum(X1)+(m-r)*(t+1/para)+r*(1/para-t*ht))
106   para<-theta1
107
108   if (sqrt(sum((para-para.old)^2))/sqrt(sum(para.old^2))<tol)
109     {break}
110   para.old<-para
111 }
112 list(estimate=para,iter=j,tol=tol)
```

```
> list(estimate=para,iter=j,tol=tol)
$estimate
[1] 0.1034652

$iter
[1] 14

$tol
[1] 1e-08
```

结论: estimate 即为得到的  
样本 $\theta$ 估计值.