

基于数据分组与右删失情形下对数正态分布的参数估计 *

刘 欣 陈 惠 费鹤良

(上海师范大学数理信息学院, 上海, 200234)

摘要

本文研究了对数正态分布数据在分组与删失情形下参数的估计问题. 一是给出未知参数的极大似然估计存在且唯一的充要条件. 二是利用EM算法对参数值进行了估计.

关键词: 对数正态分布, 分组右删失, 极大似然估计, EM算法.

学科分类号: O213.2.

§1. 引言

对数正态分布是可靠性和寿命试验中广泛应用的模型之一, 已有很多著作讨论了对数正态分布在完全样本、Type-I和Type-II截尾样本下的参数估计, 可见Lawless (1983), Bain和Engelhardt (1991)等. 对于对数正态分布在分组数据下的MLE, 王静(2003)给出了MLE存在且唯一的充要条件, 以及证明了MLE具有相合性与收敛速度服从重对数律. 现在我们讨论比分组数据更为一般的情况, 就是在每个分组的右端再去掉一些未失效的产品, 这就使数据成了分组与右删失的情形. 对于数据分组和右删失情形下Weibull分布参数的MLE已有了一定的研究. Liu (2001)讨论了Weibull分布极大似然估计存在且唯一的充要条件. 本文我们讨论对数正态分布在分组与右删失情形下未知参数的MLE存在且唯一的充要条件, 并且用EM算法给出了参数的估计值.

设某产品的寿命服从对数正态分布, 其分布函数为

$$F(x; \mu, \sigma) = \int_0^x \frac{1}{\sqrt{2\pi}\sigma t} \exp\left[-\frac{(\ln t - \mu)^2}{2\sigma^2}\right] dt, \quad x > 0. \quad (1.1)$$

这里 μ, σ 都未知且 $\sigma > 0$, $\mu \in (-\infty, +\infty)$.

现取 n 个产品进行寿命试验, 获得数据如下: 将 $[0, +\infty)$, 分成 $N + 1$ 个区间, 前 N 个区间记作 $[T_{i-1}, T_i]$, 其中 $i = 1, \dots, N$, $0 = T_0 < T_1 < \dots < T_N < \infty$. d_i 为落入到第 i 个区间 $[T_{i-1}, T_i]$ 中的失效产品数, 并记在 T_i 时刻截尾产品数为 λ_i . 这样我们有 $\sum_{i=1}^N (d_i + \lambda_i) = n$.

*国家自然科学基金(10571057, 10671129)和教育部高校博士点专项科研基金项目(20060270002).

本文2005年4月29日收到, 2005年10月13日收到修改稿.

似然函数为

$$L(\mu, \sigma; d, \lambda) = \prod_{i=1}^N \left(\int_{T_{i-1}}^{T_i} \frac{1}{\sqrt{2\pi}\sigma t} \exp \left[-\frac{(\ln t - \mu)^2}{2\sigma^2} \right] dt \right)^{d_i} \left(1 - \int_0^{T_i} \frac{1}{\sqrt{2\pi}\sigma t} \exp \left[-\frac{(\ln t - \mu)^2}{2\sigma^2} \right] dt \right)^{\lambda_i},$$

其中 $d = (d_1, d_2, \dots, d_N)$, $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_N)$.

作变换 $x = (\ln t - \mu)/\sigma$, 令 $\alpha = 1/\sigma$, $\beta = \mu/\sigma$, 则 $L(\mu, \sigma) = l(\alpha, \beta)$,

$$l(\alpha, \beta) = \prod_{i=1}^N \left[\left(\int_{\alpha \ln T_{i-1} - \beta}^{\alpha \ln T_i - \beta} g(x) dx \right)^{d_i} \left(1 - \int_0^{\alpha \ln T_i - \beta} g(x) dx \right)^{\lambda_i} \right]. \quad (1.2)$$

对数似然函数为

$$h(\alpha, \beta) = \ln l(\alpha, \beta) = \sum_{i=1}^N \left[d_i \ln \int_{\alpha \ln T_{i-1} - \beta}^{\alpha \ln T_i - \beta} g(x) dx + \lambda_i \ln \left(1 - \int_0^{\alpha \ln T_i - \beta} g(x) dx \right) \right], \quad (1.3)$$

其中 $g(x) = (1/\sqrt{2\pi}) \cdot \exp\{-x^2/2\}$. 这样我们将求 $L(\mu, \sigma)$ 的最大值点存在且唯一的充要条件转化为了求 $h(\alpha, \beta)$ 的最大值点存在且唯一的充要条件, 两者是一致的.

§2. MLE存在且唯一的充要条件

如果在整个寿命试验过程中没有一个产品失效, 那么未知参数 (μ, σ) 的MLE是不存在的. 所以我们假设至少有一个产品失效. 类似于文献[4]的方法可将数据的所有情形分成9种情况讨论:

情形1: $d_1 > 0, d_2 = \dots = d_N = 0, d_1 + \lambda_1 < n$.

情形2: $d_1 > 0, d_2 = \dots = d_N = 0, d_1 + \lambda_1 = n$.

情形3: $d_1 > 0, d_2 > 0, d_3 = \dots = d_N = \lambda_2 = \dots = \lambda_N = 0$.

情形4: $d_1 > 0, d_2 > 0, d_3 = \dots = d_N = 0$, 并且存在一个整数 $j \geq 2$, 使得 $\lambda_j > 0$.

情形5: $d_1 > 0$, 并且存在一个整数 $j > 2$, 使得 $d_j > 0$.

情形6: $d_1 = \dots = d_{N-1} = 0, d_N > 0$.

情形7: $d_1 = 0$, 存在整数 i_1, i_2 , $(1 < i_1 < i_2 \leq N)$, 使得 $d_{i_1} > 0, \lambda_{i_2} > 0$.

情形8: $d_1 = 0$, 存在整数 i_1, i_2 , $(1 < i_1 < i_1 + 1 < i_2 \leq N)$, 使得 $d_{i_1} > 0, d_{i_2} > 0$, 但是 $\lambda_{i_1+1} = \dots = \lambda_N = 0$.

情形9: $d_1 = 0$ 存在整数 i_1 , $(1 < i_1 < N)$, 使得 $d_{i_1} > 0, d_{i_1+1} > 0$, 而且 $d_{i_1+2} = \dots = d_N = \lambda_{i_1+1} = \dots = \lambda_N = 0$.

2.1 若干引理的介绍

引理 2.1^[5] 设 $g(x) = (1/\sqrt{2\pi}) \cdot \exp\{-x^2/2\}$, $G \triangleq G(u, v) = \int_u^v g(x)dx$, ($-\infty < u < v < +\infty$), 则 $\ln G$ 的 Hessian 矩阵

$$H(u, v) = \begin{pmatrix} \frac{\partial^2 \ln G}{\partial u^2} & \frac{\partial^2 \ln G}{\partial u \partial v} \\ \frac{\partial^2 \ln G}{\partial u \partial v} & \frac{\partial^2 \ln G}{\partial v^2} \end{pmatrix}$$

是负定的.

引理 2.2^[5] 设 $g(x) = (1/\sqrt{2\pi}) \cdot \exp\{-x^2/2\}$, $-\infty < c < d < +\infty$, $\alpha \in (0, +\infty)$, $\beta \in (-\infty, +\infty)$,

$$\begin{aligned} A_1 &= A_1(\alpha, \beta) = \ln \int_{\alpha c - \beta}^{\alpha d - \beta} g(x)dx, \\ A_2 &= A_2(\alpha, \beta) = \ln \int_{-\infty}^{\alpha d - \beta} g(x)dx, \\ A_3 &= A_3(\alpha, \beta) = \ln \int_{\alpha c - \beta}^{+\infty} g(x)dx, \end{aligned}$$

则 A_1 的 Hessian 矩阵 $H_1(\alpha, \beta)$ 是负定的, A_2 和 A_3 的 Hessian 矩阵 $H_2(\alpha, \beta)$ 和 $H_3(\alpha, \beta)$ 是半负定的.

引理 2.3 设 $h(\alpha, \beta)$ 由(1.3)确定, d_i, λ_i , $1 \leq i \leq N$ 如前定义, 若下列条件:

条件1: $d_1 > 0$, $d_2 > 0$, $d_3 = \dots = d_N = 0$, 并且存在一个整数 $j \geq 2$, 使得 $\lambda_j > 0$.

条件2: $d_1 > 0$, 并且存在一个整数 $j > 2$, 使得 $d_j > 0$.

条件3: $d_1 = 0$, 存在整数 i_1, i_2 , $(1 < i_1 < i_2 \leq N)$, 使得 $d_{i_1} > 0$, $\lambda_{i_2} > 0$.

条件4: $d_1 = 0$, 存在整数 i_1, i_2 , $(1 < i_1 < i_1 + 1 < i_2 \leq N)$, 使得 $d_{i_1} > 0$, $d_{i_2} > 0$, 但是 $\lambda_{i_1+1} = \dots = \lambda_N = 0$.

任何一条满足, 则 $h(\alpha, \beta)$ 的 Hessian 矩阵 $\tilde{H}(\alpha, \beta)$ 在 $(0, +\infty) \times (-\infty, +\infty)$ 上处处负定.

证明: 令 $a_i = \ln T_i$, 则 $-\infty \equiv a_0 < a_1 \dots < a_N < +\infty$. 设

$$\begin{aligned} h_d^{(i)}(\alpha, \beta) &= \ln \int_{\alpha \ln T_{i-1} - \beta}^{\alpha \ln T_i - \beta} g(x)dx = \ln \int_{\alpha a_{i-1} - \beta}^{\alpha a_i - \beta} g(x)dx, \\ h_\lambda^{(i)}(\alpha, \beta) &= \ln \left\{ 1 - \int_{\alpha \ln T_{i-1} - \beta}^{\alpha \ln T_i - \beta} g(x)dx \right\} = \ln \left\{ 1 - \int_{\alpha a_{i-1} - \beta}^{\alpha a_i - \beta} g(x)dx \right\}, \end{aligned}$$

$i = 1, 2, \dots, N$. 分别记 $h_d^{(i)}(\alpha, \beta)$ 与 $h_\lambda^{(i)}(\alpha, \beta)$ 的 Hessian 矩阵为 $H_d^{(i)}(\alpha, \beta)$ 与 $H_\lambda^{(i)}(\alpha, \beta)$. 由(1.3)式知

$$\tilde{H}(\alpha, \beta) = \sum_{i=1}^N (d_i H_d^{(i)}(\alpha, \beta) + \lambda_i H_\lambda^{(i)}(\alpha, \beta)).$$

由引理2.2知, $H_d^{(i)}(\alpha, \beta)$, $i = 2, \dots, N$ 都是负定的; $H_d^{(1)}(\alpha, \beta)$ 与 $H_\lambda^{(i)}(\alpha, \beta)$, $i = 1, \dots, N$ 都是半负定的, 从而只要至少存在一个 i , $2 \leq i \leq N$ 使 $d_i > 0$, $\tilde{H}(\alpha, \beta)$ 就是负定的.

若条件1成立, 有 $d_2 > 0$, 所以 $\tilde{H}(\alpha, \beta)$ 是负定的.

若条件2成立, 有 $j > 2$, 使 $d_j > 0$, 所以 $\tilde{H}(\alpha, \beta)$ 是负定的.

若条件3成立, 有 $1 < i_1 < i_2 \leq N$, 使 $d_{i_1} > 0$, $\lambda_{i_2} > 0$, 所以 $\tilde{H}(\alpha, \beta)$ 是负定的.

若条件4成立, 有 $1 < i_1 < i_2 \leq N$, 使 $d_{i_1} > 0$, $d_{i_2} > 0$, 所以 $\tilde{H}(\alpha, \beta)$ 是负定的.

从而上述条件中任何一条满足, $\tilde{H}(\alpha, \beta)$ 就是负定的. \square

引理 2.4 设 $f(\theta)$ 是 Θ 上有二阶连续偏导数的实值函数, Θ 是 R^k 中的凸开集. 又设梯度向量 ∇f 在 Θ 中至少有一处为零, $f(\theta)$ 的 Hessian 矩阵 $H(\theta)$ 在 Θ 中处处负定, 则 $f(\theta)$ 在 Θ 上严格凹, 在 Θ 中有唯一的最大值点, 而且没有其它的极值点和稳定点.

2.2 MLE存在且唯一的充要条件

定理 2.1 在分组与右删失数据情形下, 为了对数正态分布参数的最大似然估计存在且唯一, 必须且只需下列条件之一满足.

条件1: $d_1 > 0$, $d_2 > 0$, $d_3 = \dots = d_N = 0$, 并且存在一个整数 $j \geq 2$, 使得 $\lambda_j > 0$.

条件2: $d_1 > 0$, 并且存在一个整数 $j > 2$, 使得 $d_j > 0$.

条件3: $d_1 = 0$, 存在整数 i_1, i_2 , $(1 < i_1 < i_2 \leq N)$, 使得 $d_{i_1} > 0$, $\lambda_{i_2} > 0$.

条件4: $d_1 = 0$, 存在整数 i_1, i_2 , $(1 < i_1 < i_1 + 1 < i_2 \leq N)$, 使得 $d_{i_1} > 0$, $d_{i_2} > 0$, 但是 $\lambda_{i_1+1} = \dots = \lambda_N = 0$.

条件1, 2, 3, 4 分别是上述9种情形中的4, 5, 7, 8.

证明: 充分性: 由[3]知, 存在有界闭区域 $[\alpha_1, \alpha_2] \times [\beta_1, \beta_2]$, 使 $h(\alpha, \beta)$ 在其上存在最大值点 $(\hat{\alpha}, \hat{\beta})$ 且 $(\hat{\alpha}, \hat{\beta})$ 也是 $h(\alpha, \beta)$ 在 $D = (0, +\infty) \times (-\infty, +\infty)$ 上的最大值点. 由于 $(\hat{\alpha}, \hat{\beta})$ 是 D 之内点, 故梯度向量 ∇h 在 $(\hat{\alpha}, \hat{\beta})$ 为零. 根据引理2.3, $h(\alpha, \beta)$ 的 Hessian 矩阵是负定的. 再利用引理2.4 知 $h(\alpha, \beta)$ 的最大值点存在且唯一. 充分性证毕.

下证必要性, 只要分别对情况1, 3, 6, 9, 2 证明其MLE不存在或不唯一即可.

情形1:

$$\begin{aligned} l(\alpha, \beta) &= [\Phi(t_1)]^{d_1} \prod_{i=1}^N [1 - \Phi(t_i)]^{\lambda_i} < [\Phi(t_1)]^{d_1} \prod_{i=1}^N [1 - \Phi(t_1)]^{\lambda_i} \\ &= [\Phi(t_1)]^{d_1} [1 - \Phi(t_1)]^{n-d_1} \leq \left(\frac{d_1}{n}\right)^{d_1} \left(\frac{n-d_1}{n}\right)^{n-d_1}. \end{aligned}$$

其中, $\Phi(t)$ 为标准正态分布的分布函数. $t_i = \alpha \ln T_i - \beta$. 取 $\beta = \alpha \ln T_1 - P_{d_1/n}$, $P_{d_1/n}$ 为标准正态分布的 d_1/n 分位数, 即 $\Phi(P_{d_1/n}) = d_1/n$. 则 $\alpha \ln T_i - \beta = \alpha \ln T_i - (\alpha \ln T_1 - P_{d_1/n}) =$

$\alpha \ln(T_i/T_1) + P_{d_1/n} \longrightarrow P_{d_1/n}$, ($\alpha \rightarrow 0^+$). 此时

$$\begin{aligned}\lim_{\alpha \rightarrow 0^+} l(\alpha, \beta) &= \lim_{\alpha \rightarrow 0^+} \left\{ [\Phi(t_1)]^{d_1} \prod_{i=1}^N [1 - \Phi(t_i)]^{\lambda_i} \right\} \\ &= [\Phi(P_{d_1/n})]^{d_1} \lim_{\alpha \rightarrow 0^+} \prod_{i=1}^N \left[1 - \Phi\left(\alpha \ln \frac{T_i}{T_1} + P_{d_1/n}\right) \right]^{\lambda_i} \\ &= \left(\frac{d_1}{n} \right)^{d_1} \prod_{i=1}^N [1 - \Phi(P_{d_1/n})]^{\lambda_i} = \left(\frac{d_1}{n} \right)^{d_1} \left(\frac{n - d_1}{n} \right)^{n - d_1}.\end{aligned}$$

所以MLE不存在.

情形3: 易知 $d_1 + d_2 + \lambda_1 = n$,

$$\begin{aligned}l(\alpha, \beta) &= [\Phi(t_1)]^{d_1} [\Phi(t_2) - \Phi(t_1)]^{d_2} [1 - \Phi(t_1)]^{\lambda_1} \\ &< [\Phi(t_1)]^{d_1} [1 - \Phi(t_1)]^{d_2 + \lambda_1} \leq \left(\frac{d_1}{n} \right)^{d_1} \left(\frac{n - d_1}{n} \right)^{n - d_1}.\end{aligned}$$

取 $\beta = \alpha \ln T_1 - P_{d_1/n}$, 则 $\alpha \ln T_1 - \beta = P_{d_1/n}$, $\alpha \ln T_2 - \beta = \alpha \ln(T_2/T_1) + P_{d_1/n} \longrightarrow +\infty$, ($\alpha \rightarrow +\infty$). 此时

$$\begin{aligned}\lim_{\alpha \rightarrow +\infty} l(\alpha, \beta) &= \lim_{\alpha \rightarrow +\infty} \left\{ [\Phi(P_{d_1/n})]^{d_1} \left[\Phi\left(\alpha \ln \frac{T_2}{T_1} + P_{d_1/n}\right) - \Phi(P_{d_1/n}) \right]^{d_2} [1 - \Phi(P_{d_1/n})]^{\lambda_1} \right\} \\ &= \left[\frac{d_1}{n} \right]^{d_1} \left[1 - \Phi(P_{d_1/n}) \right]^{d_1} = \left(\frac{d_1}{n} \right)^{d_1} \left(\frac{n - d_1}{n} \right)^{n - d_1}.\end{aligned}$$

所以MLE不存在.

情形6: 若 $\lambda_N \neq 0$, 则

$$\begin{aligned}l(\alpha, \beta) &= [\Phi(t_N) - \Phi(t_{N-1})]^{d_N} \prod_{i=1}^N [1 - \Phi(t_i)]^{\lambda_i} \\ &< [\Phi(t_N)]^{d_N} [1 - \Phi(t_N)]^{\lambda_N} \\ &= \frac{d_N^{d_N} \lambda_N^{\lambda_N}}{(d_N + \lambda_N)^{d_N + \lambda_N}}.\end{aligned}$$

取 $\beta = \alpha \ln T_N - P_{d_N/(d_N + \lambda_N)}$, 则 $\alpha \ln T_N - \beta = P_{d_N/(d_N + \lambda_N)}$, $\alpha \ln T_i - \beta = \alpha \ln(T_i/T_N) + P_{d_N/(d_N + \lambda_N)} \longrightarrow -\infty$, ($\alpha \rightarrow +\infty$). 此时

$$\begin{aligned}\lim_{\alpha \rightarrow +\infty} l(\alpha, \beta) &= \lim_{\alpha \rightarrow +\infty} \left\{ \left[\Phi(P_{d_N/(d_N + \lambda_N)}) - \Phi\left(\alpha \ln \frac{T_{N-1}}{T_N} + P_{d_N/(d_N + \lambda_N)}\right) \right]^{d_N} \right. \\ &\quad \times \left. \prod_{i=1}^N \left[1 - \Phi\left(\alpha \ln \frac{T_i}{T_N} + P_{d_N/(d_N + \lambda_N)}\right) \right]^{\lambda_i} \right\} \\ &= [\Phi(P_{d_N/(d_N + \lambda_N)})]^{d_N} [1 - \Phi(P_{d_N/(d_N + \lambda_N)})]^{\lambda_N} \\ &= \frac{d_N^{d_N} \lambda_N^{\lambda_N}}{(d_N + \lambda_N)^{d_N + \lambda_N}}.\end{aligned}$$

若 $\lambda_N = 0$, 则 $l(\alpha, \beta) = [\Phi(t_N) - \Phi(t_{N-1})]^{d_N} \prod_{i=1}^N [1 - \Phi(t_i)]^{\lambda_i} < 1$. 取 $\beta = \alpha \ln T^*$, $T_{N-1} < T^* < T_N$, 则 $t_i = \alpha \ln T_i - \beta = \alpha \ln(T_i/T^*) \rightarrow -\infty$ ($i \leq N-1$) 或 $\rightarrow +\infty$ ($i = N$) ($\alpha \rightarrow +\infty$). 此时

$$\begin{aligned}\lim_{\alpha \rightarrow +\infty} l(\alpha, \beta) &= \lim_{\alpha \rightarrow +\infty} \left\{ \left[\Phi\left(\alpha \ln \frac{T_N}{T^*}\right) - \Phi\left(\alpha \ln \frac{T_{N-1}}{T^*}\right) \right]^{d_N} \times \prod_{i=1}^{N-1} \left[1 - \Phi\left(\alpha \ln \frac{T_i}{T^*}\right) \right]^{\lambda_i} \right\} \\ &= 1.\end{aligned}$$

所以MLE不存在.

情形9:

$$\begin{aligned}l(\alpha, \beta) &= \prod_{j=1}^{i_1} [1 - \Phi(t_i)]^{\lambda_j} \prod_{j=i_1}^{i_1+1} [\Phi(t_i) - \Phi(t_{i-1})]^{d_j} \\ &< [1 - \Phi(t_{i_1})]^{\lambda_{i_1}} [\Phi(t_{i_1})]^{d_{i_1}} [1 - \Phi(t_{i_1})]^{d_{i_1+1}} \\ &= [\Phi(t_{i_1})]^{d_{i_1}} [1 - \Phi(t_{i_1})]^{\lambda_{i_1} + d_{i_1+1}} \\ &\leq \frac{d_{i_1}^{d_{i_1}} (\lambda_{i_1} + d_{i_1+1})^{\lambda_{i_1} + d_{i_1+1}}}{(d_{i_1} + \lambda_{i_1} + d_{i_1+1})^{d_{i_1} + \lambda_{i_1} + d_{i_1+1}}}.\end{aligned}$$

取 $\beta = \alpha \ln T_{i_1} - P_{(\lambda_{i_1} + d_{i_1+1})/(d_{i_1} + \lambda_{i_1} + d_{i_1+1})}$, 则 $\alpha \ln(T_j/T_{i_1}) + P_{(\lambda_{i_1} + d_{i_1+1})/(d_{i_1} + \lambda_{i_1} + d_{i_1+1})} \rightarrow -\infty$ ($1 \leq j \leq i_1 - 1$) 或 $\rightarrow +\infty$ ($j = i_1 + 1$) ($\alpha \rightarrow +\infty$, 此时

$$\begin{aligned}\lim_{\alpha \rightarrow +\infty} l(\alpha, \beta) &= \lim_{\alpha \rightarrow +\infty} \left\{ \prod_{j=1}^{i_1} \left[1 - \Phi\left(\alpha \ln \frac{T_i}{T_{i_1}} + P_{(\lambda_{i_1} + d_{i_1+1})/(d_{i_1} + \lambda_{i_1} + d_{i_1+1})}\right) \right]^{\lambda_j} \right. \\ &\quad \times \prod_{j=i_1}^{i_1+1} \left[\Phi\left(\alpha \ln \frac{T_i}{T_{i_1}} + P_{(\lambda_{i_1} + d_{i_1+1})/(d_{i_1} + \lambda_{i_1} + d_{i_1+1})}\right) \right. \\ &\quad \left. - \Phi\left(\alpha \ln \frac{T_{i-1}}{T_{i_1}} + P_{(\lambda_{i_1} + d_{i_1+1})/(d_{i_1} + \lambda_{i_1} + d_{i_1+1})}\right) \right]^{d_j} \Big\} \\ &= \lim_{\alpha \rightarrow +\infty} \left\{ \prod_{j=1}^{i_1-1} \left[1 - \Phi\left(\alpha \ln \frac{T_i}{T_{i_1}} + P_{(\lambda_{i_1} + d_{i_1+1})/(d_{i_1} + \lambda_{i_1} + d_{i_1+1})}\right) \right]^{\lambda_j} \right. \\ &\quad \times \left[1 - \Phi(P_{(\lambda_{i_1} + d_{i_1+1})/(d_{i_1} + \lambda_{i_1} + d_{i_1+1})}) \right]^{\lambda_{i_1}} \left[\Phi(P_{(\lambda_{i_1} + d_{i_1+1})/(d_{i_1} + \lambda_{i_1} + d_{i_1+1})}) \right. \\ &\quad \left. - \Phi\left(\alpha \ln \frac{T_{i-1}}{T_{i_1}} + P_{(\lambda_{i_1} + d_{i_1+1})/(d_{i_1} + \lambda_{i_1} + d_{i_1+1})}\right) \right]^{d_{i_1}} \\ &\quad \times \left[\Phi\left(\alpha \ln \frac{T_{i_1+1}}{T_{i_1}} + P_{(\lambda_{i_1} + d_{i_1+1})/(d_{i_1} + \lambda_{i_1} + d_{i_1+1})}\right) \right. \\ &\quad \left. - \Phi(P_{(\lambda_{i_1} + d_{i_1+1})/(d_{i_1} + \lambda_{i_1} + d_{i_1+1})}) \right]^{d_{i_1+1}} \Big\} \\ &= \frac{d_{i_1}^{d_{i_1}} (\lambda_{i_1} + d_{i_1+1})^{\lambda_{i_1} + d_{i_1+1}}}{(d_{i_1} + \lambda_{i_1} + d_{i_1+1})^{d_{i_1} + \lambda_{i_1} + d_{i_1+1}}}.\end{aligned}$$

所以MLE不存在.

情形2:

$$l(\alpha, \beta) = [\Phi(t_1)]^{d_1} [1 - \Phi(t_1)]^{\lambda_1} \leq \left(\frac{d_1}{d_1 + \lambda_1} \right)^{d_1} \left(1 - \frac{d_1}{d_1 + \lambda_1} \right)^{\lambda_1} = \frac{d_1^{d_1} \lambda_1^{\lambda_1}}{(d_1 + \lambda_1)^{d_1 + \lambda_1}}.$$

此时考虑选取 α, β 使 $\Phi(t_1) = \int_{-\infty}^{\alpha \ln T_1} g(x) dx = d_1 / (d_1 + \lambda_1)$, 就有

$$l(\alpha, \beta) = \frac{d_1^{d_1} \lambda_1^{\lambda_1}}{(d_1 + \lambda_1)^{d_1 + \lambda_1}}$$

成立, 而这样 α, β 的有无穷多组, 所以此时 MLE 存在但不唯一. \square

§3. EM算法

我们在第2节中讨论了MLE存在且唯一的充要条件. 虽然似然方程有可能通过直接计算(如用牛顿法)求解, 但EM算法更方便高效, 所以本节介绍用EM算法求得参数的估计. 记 $X'_j, j = 1, \dots, n$ 为独立同分布随机变量, 其分布函数为(1.1). 作变换 $X = \ln X'$, 记 $a_i = \ln T_i, i = 1, \dots, N$. 则 $X_j, i = 1, \dots, n$ 独立同分布于正态分布 $N(\mu, \sigma^2)$. 它们分别落入区间 $[a_{i-1}, a_i]$ 或在 a_i 时刻被截尾, 即落入区间 $[a_i, +\infty)$, 我们只能观测到落在区间 $[a_{i-1}, a_i]$ 的 X_j 的个数 d_i 及在 a_i 被截尾的个数 λ_i , 其中 $i = 1, \dots, N, -\infty = a_0 < a_1 < \dots < a_N < +\infty$. 记随机变量 X_j 的全体为 X , 其总数为 n , 观测结果为 Y , X_{ih}, X_{il} 分别为落入区间 $[a_{i-1}, a_i]$ 与在 a_i 被截尾的随机变量.

E步: 注意到这样一个事实, X 实际上已经包含了 Y 的所有信息, 所以有 $p(\mu, \sigma^2 | X, Y) = p(\mu, \sigma^2 | X)$, 由正态分布的密度函数可以得到

$$\begin{aligned} \log p(\mu, \sigma^2 | X) &= \sum_{i=1}^N d_i \log \left[\frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(X_{ih} - \mu)^2}{2\sigma^2} \right] \right] \\ &\quad + \sum_{i=1}^N \lambda_i \log \left[\frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(X_{il} - \mu)^2}{2\sigma^2} \right] \right] \\ &= \sum_{i=1}^N (d_i + \lambda_i) (-\log \sqrt{2\pi}\sigma) - \sum_{i=1}^N d_i \frac{(X_{ih} - \mu)^2}{2\sigma^2} \\ &\quad - \sum_{i=1}^N \lambda_i \frac{(X_{il} - \mu)^2}{2\sigma^2}, \end{aligned}$$

$$\begin{aligned} Q(\mu, \sigma^2 | \mu^{(m)}, \sigma^{2(m)}, Y) &\cong E[\log p(\mu, \sigma^2 | X) | \mu^{(m)}, \sigma^{2(m)}, Y] \\ &= -n \log \sqrt{2\pi}\sigma - \sum_{i=1}^N d_i E \left[\frac{(X_{ih} - \mu)^2}{2\sigma^2} \middle| \mu^{(m)}, \sigma^{2(m)}, Y \right] \\ &\quad - \sum_{i=1}^N \lambda_i E \left[\frac{(X_{il} - \mu)^2}{2\sigma^2} \middle| \mu^{(m)}, \sigma^{2(m)}, Y \right]. \end{aligned}$$

X_{ih} 的条件密度

$$f(t|\mu^{(m)}, \sigma^{2(m)}, Y) = \frac{\frac{1}{\sqrt{2\pi}\sigma^{(m)}} \exp\left[-\frac{(t-\mu^{(m)})^2}{2\sigma^{2(m)}}\right]}{\int_{a_{i-1}}^{a_i} \frac{1}{\sqrt{2\pi}\sigma^{(m)}} \exp\left[-\frac{(x-\mu^{(m)})^2}{2\sigma^{2(m)}}\right] dx},$$

$$j = 1, 2, \dots, N, t \in [a_{i_1}, a_i).$$

X_{il} 的条件密度

$$f(t|\mu^{(m)}, \sigma^{2(m)}, Y) = \frac{\frac{1}{\sqrt{2\pi}\sigma^{(m)}} \exp\left[-\frac{(t-\mu^{(m)})^2}{2\sigma^{2(m)}}\right]}{1 - \int_{-\infty}^{a_i} \frac{1}{\sqrt{2\pi}\sigma^{(m)}} \exp\left[-\frac{(x-\mu^{(m)})^2}{2\sigma^{2(m)}}\right] dx},$$

$$j = 1, 2, \dots, N, t \in [a_i, +\infty).$$

为方便起见将 $f(x_{ih}) = f(t|\mu^{(m)}, \sigma^{2(m)}, Y)$, $f(x_{il}) = f(t|\mu^{(m)}, \sigma^{2(m)}, Y)$, 分别记为 $P_{ih}(t)$, $P_{il}(t)$. 则有

$$Q(\mu, \sigma^2 | \mu^{(m)}, \sigma^{2(m)}, Y) = -\frac{n}{2} \log 2\pi\sigma^2 - \sum_{i=1}^N d_i \int_{a_{i-1}}^{a_i} \frac{(x-\mu)^2}{2\sigma^2} P_{ih}(x) dx - \sum_{i=1}^N \lambda_i \int_{a_i}^{+\infty} \frac{(x-\mu)^2}{2\sigma^2} P_{il}(x) dx.$$

M步: 将 $Q(\mu, \sigma^2 | \mu^{(m)}, \sigma^{2(m)}, Y)$ 分别对 μ, σ^2 求导, 求出使 $Q(\mu, \sigma^2 | \mu^{(m)}, \sigma^{2(m)}, Y)$ 极大的点 $(\mu^{(m+1)}, \sigma^{2(m+1)})$.

首先对 μ 求导

$$\frac{\partial Q}{\partial \mu} = \sum_{i=1}^N d_i \int_{a_{i-1}}^{a_i} \frac{x-\mu}{\sigma^2} P_{ih}(x) dx + \sum_{i=1}^N \lambda_i \int_{a_i}^{+\infty} \frac{x-\mu}{\sigma^2} P_{il}(x) dx.$$

令 $\partial Q / \partial \mu = 0$, 有

$$\mu = \frac{\sum_{i=1}^N d_i \int_{a_{i-1}}^{a_i} x P_{ih}(x) dx + \sum_{i=1}^N \lambda_i \int_{a_i}^{+\infty} x P_{ih}(x) dx}{\sum_{i=1}^N d_i \int_{a_{i-1}}^{a_i} P_{ih}(x) dx + \sum_{i=1}^N \lambda_i \int_{a_i}^{+\infty} P_{ih}(x) dx}. \quad (3.1)$$

对 σ^2 求导

$$\frac{\partial Q}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \sum_{i=1}^N d_i \int_{a_{i-1}}^{a_i} \frac{(x-\mu)^2}{2\sigma^4} P_{ih}(x) dx + \sum_{i=1}^N \lambda_i \int_{a_i}^{+\infty} \frac{(x-\mu)^2}{2\sigma^4} P_{il}(x) dx.$$

令 $\partial Q / \partial \sigma^2 = 0$, 有

$$\sigma^2 = \frac{2}{n} \left\{ \sum_{i=1}^N d_i \int_{a_{i-1}}^{a_i} \frac{(x-\mu)^2}{2} P_{ih}(x) dx + \sum_{i=1}^N \lambda_i \int_{a_i}^{+\infty} \frac{(x-\mu)^2}{2} P_{il}(x) dx \right\}. \quad (3.2)$$

利用(3.1), (3.2)所得到的 μ, σ^2 就是我们所要寻找的 $\mu^{(m+1)}, \sigma^{2(m+1)}$. 这样就完成了一次迭代: $(\mu^{(m)}, \sigma^{2(m)}) \rightarrow (\mu^{(m+1)}, \sigma^{2(m+1)})$.

§4. 模拟

记 $X'_i, i = 1, \dots, n$ 为独立同分布于对数正态分布 $LN(\mu, \sigma^2)$ 的随机变量, 取 $T_0 = 0, T_1 = 0.3, T_2 = 3, T_3 = 30, T_4 = 300, T_5 = 3000, T_6 = 30000, T_7 = 300000$, 误差精度0.01. 对 $j \leq 6$, 产品在 T_j 被截尾的概率取为 $j/7$, 而到 T_7 所有未失效的产品都被截尾. 分别对 $\mu = 5, \sigma = 2$ 与 $\mu = 8, \sigma = 1.5$ 进行迭代, 模拟结果见下表. (其中 μ', σ' 为参数初值)

表1 取 $\mu = 5, \sigma = 2, \mu' = 6, \sigma' = 1.5$ 的
模拟结果

样本量	模拟次数	$\bar{\mu}$	$\bar{\sigma}$
100	10	4.8659	1.9675
	30	4.8732	1.9644
	50	5.0945	1.9634
300	10	5.0294	2.0967
	30	4.9711	2.0899
	50	5.0295	2.0866
800	10	4.9774	2.1063
	30	5.0238	2.1057
	50	5.0297	2.1049
2000	10	5.0338	2.1077
	30	5.0197	2.1067
	50	4.9826	2.1054

表2 取 $\mu = 8, \sigma = 1.5, \mu' = 6, \sigma' = 2$ 的
模拟结果

样本量	模拟次数	$\bar{\mu}$	$\bar{\sigma}$
100	10	7.9359	1.4733
	30	7.9488	1.4655
	50	8.1127	1.4653
1-4	10	8.0687	1.4890
	30	8.0846	1.4881
	50	7.9711	1.4879
800	10	8.0354	1.5237
	30	7.9822	1.5211
	50	7.9823	1.5196
1-4	10	8.0346	1.5568
	30	8.0344	1.5509
	50	8.0288	1.5490

模拟结果表明随着样本量的增大, 估计值与参数真值的偏差有明显的减小. 当样本量为300时, 两个参数的平均偏差都在5%之内, 效果不错.

§5. 小结

本文针对数据分组与右删失情形下对数正态分布的参数估计问题进行了讨论, 主要给出未知参数的极大似然估计存在且唯一的充要条件; 并且介绍了如何利用EM算法对参数值进行估计. 从模拟结果中, 不难看出对于服从对数正态分布的分组右删失数据, 特别是样本较大的情况, 利用由EM算法得出的迭代公式可以得到对参数 μ, σ 的比较满意的估计 $\hat{\mu}, \hat{\sigma}$.

相应地, 也可以讨论这种情形下参数的区间估计问题, 如可以考虑极大似然估计是否具有渐近正态性, 从而研究区间估计或利用Bootstrap方法. 本文不做深入讨论.

参考文献

- [1] Lawless, J.F, *Statistical Models and Methods for Lifetime Data*, John wilay & sons, 1983.
- [2] Bain, L.J, Engelhardt, M., *Statistical Analysis of Reliability and Life Testing Models, Theory and Methods*, 2nd ed, New York: Marcel Dekker, 1991.
- [3] Cheng, K.F., Chen, C.H., Estimation of the Weibull parameters with grouped data, *Commu. Statist.-Theor.*, **17(2)**(1988), 325–341.
- [4] Liu, L.P., Existence of MLE for Weibull distribution with grouped and censored data, *Chinese Journal of Applied Probability and Statistics*, **17(2)**(2001), 133–138.
- [5] 王静, 分组数据情形下对数正态分布参数的最大似然估计, *应用数学学报*, **26(4)**(2003), 737–744.
- [6] 郑明, 杨艺, 郑宇, 基于分组数据的Weibull分布的参数估计, *高校应用数学学报A辑*, **18(3)**(2003), 303–310.

Estimation of the Parameters in the Lognormal Distribution with Grouped and Right-Censored Data

LIU XIN CHEN HUI FEI HELIANG

(College of Mathematics and Sciences, Shanghai Normal University, Shanghai, 200234)

In this paper we consider the estimation problem in the situation where the life of some products follows a lognormal distribution, and the observed data are grouped and right-censored. A necessary and sufficient condition for the existence and uniqueness of the Maximum Likelihood Estimate (MLE) of unknown parameters is given. Furthermore, we have estimated the unknown parameters by EM algorithm.

Keywords: Lognormal distribution, grouping and right-censoring, MLE, EM algorithm.

AMS Subject Classification: 62N01, 62N02.