

线性模型下分布函数的估计

李斌

(盐城工学院学生处, 盐城 224003)

摘要 利用辅助信息, 针对线性模型 $y_i = \beta x_i + x_i^q \epsilon_i, \epsilon_i, i. i. d., E\epsilon_i = 0$, 给出了总体分布函数的估计量: $\hat{F}(t) = \frac{1}{N} [\sum_{j \in S} \Delta(t - y_j) + \sum_{i \in S} \frac{1}{n} \sum_{j \in S} \Delta(t - \hat{\beta} \cdot x_i - x_i^{q^*} u_j)]$, $u_j = (y_j - \hat{\beta} \cdot x_j) / x_j^{q^*}, j \in S$, 改进了 Chambers-Dunstan 的结果。

关键词 抽样调查 分布函数 辅助信息 有限总体

分类号 O51

随着社会的发展, 抽样调查的应用日益广泛。但大多数关于抽样调查的文献及教科书都着重讨论总体均值的估计方法。然而在实际工作中, 常要估计某一有限总体的分布函数。在简单随机样本的情况下, 我们可以用经验分布函数作为总体分布函数的估计, Gliveko-Cantelli 定理^[1]证明了其相合性。

估计某一总体 Y 的分布函数, 常有与之相关的总体 X (已知的辅助信息) 可以利用, 用 Y_1, \dots, Y_N 表示待估有限总体, X_1, \dots, X_N 表示辅助信息, $(x_1, y_1), \dots, (x_n, y_n)$ 为一样本, $F_Y(t)$ 表示 Y 的分布函数。问题是如何给出 $F_Y(t)$ 的估计量? 近几年来有不少文献对此进行了讨论。

1986年, Chambers, Dunstan 假设总体服从模型:

$$y_i = \beta x_i + x_i^q \epsilon_i, \quad \epsilon_i, i. i. d., E\epsilon_i = 0 \tag{1}$$

其中 $i=1, 2, \dots, N, 0 \leq q \leq 1$ 已知, β 为待估参数。在一个具体问题中, 取 $q = \frac{1}{2}$, 给出了估计量

$$\hat{F}_m(t) = \frac{1}{N} [\sum_{j \in S} \Delta(t - y_j) + \sum_{i \in S} \frac{1}{n} \sum_{j \in S} \Delta(t - \hat{\beta} x_i - \sqrt{x_i} u_j)] \tag{2}$$

其中, $\Delta(x) = \begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases}$

$$\hat{\beta} = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i}, \quad u = \frac{y_i - \hat{\beta} x_i}{\sqrt{x_i}}$$

并证明了 $\hat{F}_m(t)$ 为模型无偏估计量^[2]。

本文主要对 Chambers-Dunstan 方法进行改进, 讨论线性模型(1)下, 分布函数的估计。

对于(1), 如果 q 已知, 利用最小二乘法估计 β , 可得:

$$\hat{\beta} = \frac{\sum_{i=1}^n y_i x_i^{1-2q}}{\sum_{i=1}^n x_i^{1-2q}} \tag{3}$$

而 q 未知时, 可用 t 替换

$$\hat{\beta}(t) = \frac{\sum_{i=1}^n y_i x_i^{1-2t}}{\sum_{i=1}^n x_i^{1-2t}} \tag{4}$$

如此, 只要考虑 t 对估计量的影响。

定理 1: 设 $x_1 < x_2 < \dots < x_{n-1} < x_{2m}, K = \max\{i | x_i < x_m\}, L = \min\{i | x_i > x_{n-1}\}$, 如果 $\sum_{i=1}^L y_i^2 > 0$, 且 $\sum_{i=L}^m y_i^2 > 0$. 则方程

$$\sum_{i=1}^m \left(\frac{y_i}{x_i}\right)^2 = \sum_{i=m+1}^{2m} \left(\frac{y_i}{x_i}\right)^2 \quad (5)$$

存在唯一解 q .

证明: 不妨令

$$S(t) = \frac{\sum_{i=1}^m \left(\frac{y_i}{x_i}\right)^2}{\sum_{i=m+1}^{2m} \left(\frac{y_i}{x_i}\right)^2} \quad (6)$$

$\forall \Delta > 0, S(t) > 0$

$$\begin{aligned} S(t + \Delta) &= \frac{\sum_{i=1}^m \left(\frac{y_i}{x_i + \Delta}\right)^2}{\sum_{i=m+1}^{2m} \left(\frac{y_i}{x_i + \Delta}\right)^2} = \\ &= \frac{\sum_{i=1}^m \left(\frac{y_i}{x_i}\right)^2 \cdot \left(\frac{1}{x_i + \Delta}\right)^2}{\sum_{i=m+1}^{2m} \left(\frac{y_i}{x_i}\right)^2 \cdot \left(\frac{1}{x_i + \Delta}\right)^2} = \\ &= \frac{\sum_{i=1}^m \left(\frac{y_i}{x_i}\right)^2 \cdot \left(\frac{x_m}{x_i}\right)^{2\Delta}}{\sum_{i=m+1}^{2m} \left(\frac{y_i}{x_i}\right)^2 \cdot \left(\frac{x_m}{x_i}\right)^{2\Delta}} \quad (7) \end{aligned}$$

$$\sum_{i=1}^K \left(\frac{y_i}{x_i}\right)^2 \cdot \left(\frac{x_m}{x_i}\right)^{2\Delta} > \sum_{i=1}^K \left(\frac{y_i}{x_i}\right)^2 > 0$$

$$\sum_{i=K+1}^m \left(\frac{y_i}{x_i}\right)^2 \cdot \left(\frac{x_m}{x_i}\right)^{2\Delta} = \sum_{i=K+1}^m \left(\frac{y_i}{x_i}\right)^2$$

因此 $\sum_{i=1}^m \left(\frac{y_i}{x_i}\right)^2 \cdot \left(\frac{x_m}{x_i}\right)^{2\Delta} > \sum_{i=1}^m \left(\frac{y_i}{x_i}\right)^2 > 0$ (8)

同理 $0 < \sum_{i=m+1}^{2m} \left(\frac{y_i}{x_i}\right)^2 \cdot \left(\frac{x_m}{x_i}\right)^{2\Delta} < \sum_{i=m+1}^{2m} \left(\frac{y_i}{x_i}\right)^2$ (9)

从(7)(8)(9)知, $S(t + \Delta) > S(t)$, 即 $S(t)$ 严格单调增加.

又 $S(+\infty) = \lim_{t \rightarrow +\infty} \left[\frac{\sum_{i=1}^m \left(\frac{y_i}{x_i}\right)^2}{\sum_{i=m+1}^{2m} \left(\frac{y_i}{x_i}\right)^2} \right] = \lim_{t \rightarrow +\infty} \left\{ \frac{[\sum_{i=1}^m y_i^2 \left(\frac{x_i}{x_i}\right)^{2t}]}{[\sum_{i=m+1}^{2m} y_i^2 \left(\frac{x_i}{x_i}\right)^{2t}]} \right\} = +\infty$ (10)

$\left(\frac{x_K}{x_i}\right)^{2t} > 0, \lim_{t \rightarrow +\infty} \sum_{i=m+1}^{2m} y_i^2 \left(\frac{x_i}{x_i}\right)^{2t} = 0$

同理 $S(-\infty) = 0$, 由介值定理知 $S(t) = 1$ 有唯一解, 即(5)有唯一解 q .

定理 2: 设 $2 < X < M, x_1 \leq x_2 \leq \dots \leq x_n$ 且等号不全成立, 则下述结论成立: ① 任意实数 $t, \hat{\beta}(t)$ 是 β 无偏估计量. ② 任意实数 $t, \hat{\beta}(t)$ 是 β 的弱相合估计. ③ n 固定时, 任意 $t \neq q$, 有 $Var(\hat{\beta}(q)) < Var(\hat{\beta}(t))$

证明: ① $\hat{\beta}(t) = \frac{\sum_{i=1}^n [(\beta x_i + x_i^q \epsilon_i) x_i^{1-2t}]}{\sum_{i=1}^n x_i^{2-2t}} = \frac{\sum_{i=1}^n \beta x_i^{2-2t}}{\sum_{i=1}^n x_i^{2-2t}} + \frac{\sum_{i=1}^n x_i^{1+q-2t} \epsilon_i}{\sum_{i=1}^n x_i^{2-2t}} = \beta + \frac{\sum_{i=1}^n x_i^{1+q-2t} \epsilon_i}{\sum_{i=1}^n x_i^{2-2t}}$

$$E(\hat{\beta}(t)) = \beta + \frac{\sum_{i=1}^n x_i^{1+q-2t} E\epsilon_i}{\sum_{i=1}^n x_i^{2-2t}} = \beta$$

② $Var(\hat{\beta}(t)) = \frac{\sum_{i=1}^n x_i^{2+2q-4t} Var(\epsilon_i)}{\left(\sum_{i=1}^n x_i^{2-2t}\right)^2} = \sigma_\epsilon^2 \cdot \frac{\sum_{i=1}^n x_i^{2+2q-4t}}{\left(\sum_{i=1}^n x_i^{2-2t}\right)^2}$

由于 $2 < X < M$, 任取 n, i, j 有: $0 < \frac{x_i^{1+q-2t}}{x_i^{2-2t}} < M < +\infty$

即: $0 < x_i^{1+q-2t} < M x_i^{2-2t}$ 对任意 n, i, j

令 $x_k^{2-2t} = \min\{x_i^{2-2t} | i = 1, 2, \dots, n\}$

$x_L^{1+q-2t} = \max\{x_i^{1+q-2t} | i = 1, 2, \dots, n\}$

则 $Var(\hat{\beta}(t)) \leq \frac{\sum_{i=1}^n (x_i^{1+q-2t})^2}{\left(\sum_{i=1}^n x_k^{2-2t}\right)^2} \cdot \sigma_\epsilon^2 = \frac{n(x_L^{1+q-2t})^2}{n^2(x_k^{2-2t})^2} \cdot \sigma_\epsilon^2 \leq \frac{(M^t)^2}{n} \cdot \sigma_\epsilon^2 \rightarrow 0 (n \rightarrow +\infty \text{ 时})$

又 $\hat{\beta}(t)$ 无偏, 故 $\hat{\beta}(t)$ 为 β 的弱相合估计。

$$\text{③} \quad \text{令 } V(t) \triangleq \text{Var}(\hat{\beta}(t)) / \sigma_0^2$$

$$\text{则 } V(t) = \frac{\sum_{i=1}^n (x_i^{1+q-2t})^2 / (\sum_{i=1}^n x_i^{2-2t})^2}{(\sum_{i=1}^n x_i^{1-q-2t}) (\sum_{i=1}^n x_i^{1-2t})^{-1}}$$

$$\frac{dV(t)}{dt} = \left(\sum_{i=1}^n x_i^{2+2q-4t} \cdot \ln x_i \cdot (-4) \right) \left(\sum_{i=1}^n x_i^{1-2t} \right)^{-2} - \left(\sum_{i=1}^n x_i^{1-q-2t} \right) (-2) \left(\sum_{i=1}^n x_i^{1-2t} \right)^{-3} \left(\sum_{i=1}^n x_i^{1-2t} \ln x_i \cdot (-2) \right) =$$

$$4 \left[\sum_{i=1}^n x_i^{2+2q-4t} \right]^{-1} \left\{ \left(\sum_{i=1}^n x_i^{2+2q-4t} \right) \left(\sum_{i=1}^n x_i^{1-2t} \ln x_i \right) - \left(\sum_{i=1}^n x_i^{1-q-2t} \ln x_i \right) \left(\sum_{i=1}^n x_i^{1-2t} \right) \right\} \triangleq$$

$$4 \left[\sum_{i=1}^n x_i^{2+2q-4t} \right]^{-1} W(t)$$

$$4 \left[\sum_{i=1}^n x_i^{2+2q-4t} \right]^{-1} > 0, W(t), \frac{dV(t)}{dt}, \frac{d\text{Var}(\hat{\beta}(t))}{dt} \text{ 同号}$$

$$W(t) = \sum_{i,j=1}^n [x_i^{2+2q-4t} x_j^{-2t} \ln x_j] - \sum_{i,j=1}^n [x_j^{1+2q-4t} x_i^{-t} \ln x_j] = \sum_{i,j=1}^n x_i^{-t} x_j^{1-2t} \ln x_j [x_i^{2q-2t} - x_j^{2q-2t}] =$$

$$\sum_{i,j} \{ x_i^{2-2t} x_j^{2-2t} [x_i^{q-t} + x_j^{q-t}] \} \ln x_j [x_i^{q-t} - x_j^{q-t}] = \sum_{i \neq j} D(i, j) \ln x_j [x_i^{q-t} - x_j^{q-t}] =$$

$$\sum_{i > j} D(i, j) \ln x_i [x_i^{q-t} - x_j^{q-t}] + \sum_{i < j} D(i, j) \ln x_j [x_i^{q-t} - x_j^{q-t}] =$$

$$\sum_{i > j} D(j, i) \ln x_i [x_i^{q-t} - x_j^{q-t}] + \sum_{i > j} D(i, j) \ln x_j [x_i^{q-t} - x_j^{q-t}] =$$

$$\sum_{i > j} D(i, j) [\ln x_j - \ln x_i] [x_i^{q-t} - x_j^{q-t}]$$

其中 $D(i, j) = x_i^{2-2t} x_j^{2-2t} [x_i^{q-t} + x_j^{q-t}] = D(j, i) > 0$

由此可得, 当 $t=q$ 时, $W(t)=0$; $t>q$ 时, $q-t<0, x_i^{q-t}-x_j^{q-t} \geq 0 (i>j)$, 等号成立的充要条件是 $x_i = x_j$, 于是当 x_1, \dots, x_n 不全相等时, $W(t)>0$ 。同理, $t<q$ 时, $W(t)<0$ 。又 $W(t), d(\text{Var}(\hat{\beta}(t)))/dt$ 同号, 故 $\text{Var}(\hat{\beta}(t))$ 在 $t=q$ 处取最小值。

综上, 如果有 \hat{q} , 可用 $\hat{\beta}(\hat{q})$ 作为 β 的估计量。而对给定的 $\hat{\beta}$, 对 $y_i - \hat{\beta}x_i, x_i$ 用定理 1 解方程 $\sum_{i=1}^m \frac{(y_i - \hat{\beta}x_i)^2}{x_i^2} = \sum_{i=m+1}^{2m} \frac{(y_i - \hat{\beta}x_i)^2}{x_i^2}$ (10) 可得解 $\hat{q}(\hat{\beta})$ 。

任取 $t \in [0, 1]$ 作为迭代初值, 由定理 2 知 $\hat{\beta}(t) \triangleq \hat{\beta}_1$ 是 β 的无偏、相合估计, 只是估计量的方差较 $\text{Var}(\hat{\beta}(\hat{q}))$ 大些, 再利用方程(10)得到 $\hat{q}_1 = \hat{q}(\hat{\beta}(t))$, 依此类推反复迭代, 对 $\hat{\beta}_k, \hat{q}_k$, 令 $\hat{\beta}_{k+1} = \hat{\beta}(\hat{q}_k), \hat{q}_{k+1} = \hat{q}(\hat{\beta}_{k+1}), \dots$ 设 $\hat{\beta}_k \rightarrow \hat{\beta}^*, \hat{q}_k \rightarrow \hat{q}^*, \hat{\beta}^*, \hat{q}^*$ 就是最终相应参数的估计量, 据此给出有限总体分布函数的估计量:

$$\hat{F}(t) = \frac{1}{N} \left[\sum_{j \in S} \Delta(t - y_j) \right] + \sum_{i=1}^s \frac{1}{n} \sum_{j=1}^n \Delta(t - \hat{\beta}^* x_j - x_j^{\hat{q}^*} u_i) \tag{11}$$

其中 $u_i = (y - \hat{\beta}^* x_j) / x_j^{\hat{q}^*}, j \in S$ 。

参 考 文 献

- 1 陈家鼎, 李东风. 数理统计学讲义. 北京: 高等教育出版社, 1993
- 2 Chambers, R. L & Dunstan, R. (1986) Estimating distribution functions from survey data. Biometrika, Vol. 73

The Estimation of Scattering Functions Under the Lined Model

Lu Bin

(The Students Department of Yancheng Institute of Technology, Yancheng 224003, PRC)

Abstract Using assisted message, the writing gives the estimating quality of the total scattering functions: $\hat{F}(t) = \frac{1}{N} \left[\sum_{j \in S} \Delta(t - y_j) + \sum_{i=1}^s \frac{1}{n} \sum_{j=1}^n \Delta(t - \hat{\beta}^* x_j - x_j^{\hat{q}^*} u_i) \right], u_i = (y_i - \hat{\beta}^* x_i) / x_i^{\hat{q}^*}, j \in S$, under the lined model: $y_i = \beta x_i + x_i^q \varepsilon_i, \varepsilon_i \sim i.i.d., E\varepsilon_i = 0$ and improves the result of Chambers-Dunstan.

Keywords selective examination; scattering functions; assisted message; limited entity