

相关性度量与估计效率

刘兰亭 张信东

摘要 在这篇文章中我们给出了最小二乘拟合值与残差之间的相关性度量和最小二乘估计相对效率的关系

关键词 相关性度量 估计效率

一 引言

考虑线性模型

$$y = X\beta + e, e \sim N(0, \Gamma) \quad (1)$$

其中 y 是 $n \times 1$ 观测向量, X 是已知的 $n \times p$ 设计矩阵, β 是未知的 $p \times 1$ 参数向量, e 是 $n \times 1$ 随机误差向量, Γ 是 $n \times n$ 正定矩阵, 即 $\Gamma > 0$ 。假定 $rk(X) = p$, $n \geq 2p$ 。 β 的 BLU 估计(即最佳线性无偏估计)与 LS 估计(即最小二乘估计)分别为

$$\beta^* = (X^T \Gamma^{-1} X)^{-1} X^T \Gamma^{-1} y, \text{Cov}(\beta^*) = (X^T \Gamma^{-1} X)^{-1}$$

$$\hat{\beta} = (X^T X)^{-1} X^T y, \text{Cov}(\hat{\beta}) = (X^T X)^{-1} X^T \Gamma X (X^T X)^{-1}$$

熟知 $\text{Cov}(\beta^*) \leq \text{Cov}(\hat{\beta})$ 。但是, 在许多实际问题中, Γ 是未知的, 于是只能用 $\hat{\beta}$ 来代替 β^* , 这就要蒙受一些损失。为此, 许多作者利用

$$e(\hat{\beta}) = \frac{|\text{Cov}(\beta^*)|}{|\text{Cov}(\hat{\beta})|} \quad (2)$$

作为 LS 估计的相对效率。显然, $0 \leq e(\hat{\beta}) \leq 1$, 而且其值愈大, 表示估价 $\hat{\beta}$ 与 β^* 愈接近, 当然所蒙受的损失就愈小。

许多专家和学者提出和研究了度量两个随机向量 x 和 z 相关程度的数量指标。例如相关系数, Hotelling 和张尧庭提出的广义相关系数和 Lindley (见[1]) 提出的相关性量度 L 。对一般线性模型, 王松桂 (见[2]) 研究了广义相关系数与估计效率的关系。本文讨论在模型(1)情形相关性度量 L 与估计效率的关系, 并据此得出相关性度量的界。

对随机向量 x 和 z , Lindley 相关性度量定义为

$$L(\dot{x}, z) = \iint f(s, t) \ln \frac{f(s, t)}{g(s) \cdot h(t)} ds \cdot dt \quad (3)$$

其中 $g(s)$, $h(t)$ 和 $f(s, t)$ 分别为 x 、 z 和 $(x^r, z^r)^r$ 的分布密度或联合分布密度。

二 主要结果

不失一般性假定 $X^r X = I_p$, 选 $n \times (n-p)$ 矩阵 N 使 $(X \mid N)^r (X \mid N) = I_n$ 。记 $\Sigma_1 = XX^r \Gamma XX^r$, $\Sigma_2 = NN^r \Gamma NN^r$, 和

$$\Sigma = \begin{pmatrix} XX^r \Gamma XX^r & XX^r \Gamma NN^r \\ NN^r \Gamma XX^r & NN^r \Gamma NN^r \end{pmatrix}$$

易知 $\text{rk}(\Sigma_1) = p$, $\text{rk}(\Sigma_2) = n-p$ 和 $\text{rk}(\Sigma) = n$

根据 Rao 介绍的奇异型正态分布 (见 [3]), 我们可以得到最小二乘拟合值 $XX^r y$ 、残差 $NN^r y$ 和 $\begin{pmatrix} XX^r y \\ NN^r y \end{pmatrix}$ 的分布密度分别为

$$g(s) = \frac{1}{(2\pi)^{\frac{p}{2}} \cdot (\det(X^r \Gamma X))^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(s - X\beta)^r \Sigma_1^{-1}(s - X\beta)\right\} \quad (4)$$

$$N_1^r s = N_1^r X\beta \quad \text{a.s.}$$

$$h(t) = \frac{1}{(2\pi)^{\frac{n-p}{2}} \cdot (\det(N^r \Gamma N))^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}t^r \Sigma_2^{-1}t\right\} \quad (5)$$

$$N_2^r t = 0 \quad \text{a.s.}$$

$$f(s, t) = \frac{1}{(2\pi)^{\frac{n}{2}} (\det \Gamma)^{\frac{1}{2}}} \exp\left\{-\frac{1}{2} \begin{pmatrix} s - X\beta \\ t \end{pmatrix}^r \Sigma^{-1} \begin{pmatrix} s - X\beta \\ t \end{pmatrix}\right\} \quad (6)$$

$$N_3^r \begin{pmatrix} s \\ t \end{pmatrix} = N_3^r \begin{pmatrix} X\beta \\ 0 \end{pmatrix}$$

其中 N_1 , N_2 和 N_3 分别是 $n \times (n-p)$, $n \times p$ 和 $2n \times n$ 矩阵, 且分别满足

$$\mu(N_1) = \mu^\perp(\Sigma_1), \quad \mu(N_2) = \mu^\perp(\Sigma_2), \quad \mu(N_3) = \mu^\perp(\Sigma)$$

$\mu(A)$ 表示由 A 的列向量张成的线性子空间, “ \perp ” 表示线性子空间的正交补空间, Σ_1^- , Σ_2^- 和 Σ^- 分别为 Σ_1 , Σ_2 和 Σ 的广义逆。由于密度函数与广义逆的选取无关, 故可用 Σ_1 和 Σ_2 的 “+” 号逆 Σ_1^+ 和 Σ_2^+ 分别代替 Σ_1^- 和 Σ_2^- , 至于 Σ^- 我们亦可用 Σ 的 “+” 号逆代替, 亦可任选一个广义逆, 譬如

$$\Sigma^- = \begin{pmatrix} XX^r \Gamma^{-1} XX^r & XX^r \Gamma^{-1} - X(X^r \Gamma X)^{-1} X^r \\ \Gamma^{-1} XX^r - X(X^r \Gamma X)^{-1} X^r & \Gamma^{-1} \end{pmatrix} \quad (7)$$

另外我们不难得到

$$\Sigma_1^+ = X(X'GX)^{-1}X', \quad \Sigma_2^+ = N(N'GN)^{-1}N' \quad (8)$$

为了证明简洁, 定义集合 $M_1 = \left\{ \begin{pmatrix} s \\ t \end{pmatrix} : N_1^+s = 0, N_2^+t = 0 \right\}$, $M_2 = \left\{ \begin{pmatrix} s \\ t \end{pmatrix} : N_3^+ \begin{pmatrix} s \\ t \end{pmatrix} = 0 \right\}$ 和 $M = \left\{ \begin{pmatrix} s \\ t \end{pmatrix} : N_1^+s = 0, N_2^+t = 0, N_3^+ \begin{pmatrix} s \\ t \end{pmatrix} = 0 \right\}$. 则 $M_2 \subset M_1$, 从而

$$M = M_2 \quad (9)$$

为给出主要结论首先证明如下引理

引理 1 设 $\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_p^2$ 是 $XX'y$ 和 $NN'y$ 之间的典型相关系数, 则

$$\prod_{i=1}^p (1 - \rho_i^2) = \det(I_p - (X'GX)^{-1/2} X'GN(N'GN)^{-1} N'GX(X'GX)^{-1/2})$$

证明 因为

$$\text{Cov} \begin{pmatrix} XX'y \\ NN'y \end{pmatrix} = \Sigma \triangleq \begin{pmatrix} \Sigma_1 & \Sigma_{12} \\ \Sigma_{21} & \Sigma_2 \end{pmatrix}$$

其中 $\Sigma_{12} = XX'GN$, $\Sigma_{21} = \Sigma_{12}'$. 所以 $XX'y$ 和 $NN'y$ 之间的典型相关系数 ρ^2 应满足

$$\det(\rho^2 I_p - \Sigma_1^{-1} \Sigma_{12} \Sigma_2^{-1} \Sigma_{21}) = 0$$

而它又等价于

$$\det(\rho^2 I_p - X(X'GX)^{-1} X'GN(N'GN)^{-1} N'GX(X'GX)^{-1}) = 0$$

由于 $X(X'GX)^{-1} X'GN(N'GN)^{-1} N'GX(X'GX)^{-1}$ 与 $(X'GX)^{-1/2} X'GN(N'GN)^{-1} \times N'GX(X'GX)^{-1/2}$ 具有完全相同的非零特征根, 因此 $1 - \rho_i^2$ ($i = 1, 2, \dots, p$) 是 $I_p - (X'GX)^{-1/2} X'GN(N'GN)^{-1} N'GX(X'GX)^{-1/2}$ 的全部特征根. 由此即得结论.

定理 1 对于模型(1), 最小二乘拟合值 $XX'y$ 和残差 $NN'y$ 之间的相关性度量与 LS 估计相对效率具有关系

$$L(XX'y, NN'y) = -\frac{1}{2} \ln e(\hat{\beta}) \quad (10)$$

证明 将(4)–(6)式代入(3)式, 并利用(9)式

$$\begin{aligned} L(XX'y, NN'y) &= \iint_{\substack{N_1^+s = N_1^+X\beta \\ N_2^+t = 0 \\ N_3^+ \begin{pmatrix} s \\ t \end{pmatrix} = N_3^+ \begin{pmatrix} s \\ t \end{pmatrix}}} f(s, t) \cdot \ln \frac{f(s, t)}{g(s)h(t)} ds dt \\ &= \iint_M f(s + X\beta, t) \cdot \ln \frac{f(s + X\beta, t)}{g(s + X\beta)h(t)} ds dt \\ &= \iint_M -\frac{1}{2} \ln \frac{\det(\Gamma)}{\det(X'GX) \cdot \det(N'GN)} f(s + X\beta, t) ds dt \\ &\quad - \frac{1}{2} \iint_{M_2} \begin{pmatrix} s \\ t \end{pmatrix}' \left[\Sigma^- - \begin{pmatrix} \Sigma_1^+ & 0 \\ 0 & \Sigma_2^+ \end{pmatrix} \right] \begin{pmatrix} s \\ t \end{pmatrix} f(s + X\beta, t) ds dt \end{aligned}$$

$$= -\frac{1}{2} \ln \frac{\det(\Gamma)}{\det(X' \Gamma X) \cdot \det(N' \Gamma N)} - E \begin{pmatrix} s \\ t \end{pmatrix}' \left[\Sigma^- - \begin{pmatrix} \Sigma_1^+ & 0 \\ 0 & \Sigma_2^+ \end{pmatrix} \right] \begin{pmatrix} s \\ t \end{pmatrix}$$

这是由于 $f(s + X\beta, t)$ 是正态分布 $N(0, \Sigma)$ 在超平面 M_2 上的概率分布密度; 又由于

$$\begin{aligned} & E \begin{pmatrix} s \\ t \end{pmatrix}' \left[\Sigma^- - \begin{pmatrix} \Sigma_1^+ & 0 \\ 0 & \Sigma_2^+ \end{pmatrix} \right] \begin{pmatrix} s \\ t \end{pmatrix} \\ &= \text{tr} \left(\left[\Sigma^- - \begin{pmatrix} X(X' \Gamma X)^{-1} X' & 0 \\ 0 & N(N' \Gamma N)^{-1} N' \end{pmatrix} \right] \Sigma \right) \\ &= \text{tr}(\Sigma^- \Sigma) - \text{tr} \begin{pmatrix} X X' & X(X' \Gamma X)^{-1} X' \Gamma N N' \\ N(N' \Gamma N)^{-1} N' \Gamma X X' & N N' \end{pmatrix} \\ &= \text{rk}(\Sigma) - n = 0 \end{aligned}$$

于是, 利用引理 1

$$\begin{aligned} L(XX'y, NN'y) &= -\frac{1}{2} \ln \frac{\det(\Gamma)}{\det(X' \Gamma X) \cdot \det(N' \Gamma N)} \\ &= -\frac{1}{2} \ln [\det(I_p - (X' \Gamma X)^{-\frac{1}{2}} X' \Gamma N (N' \Gamma N)^{-1} N' \Gamma X (X \Gamma X)^{-\frac{1}{2}})] \\ &= -\frac{1}{2} \ln \prod_{i=1}^p (1 - \rho_i^2) \end{aligned}$$

而 Bartmann 和 Bloomfield 已证明 (见 [4])

$$e(\hat{\beta}) = \prod_{i=1}^p (1 - \rho_i^2) \quad (11)$$

故

$$L(XX'y, NN'y) = -\frac{1}{2} \ln e(\hat{\beta}) \quad \text{定理证毕}$$

根据 [5] ((2.8) 式) 与定理 1 不难得到

推论 1 在定理 1 条件之下, 有

$$0 \leq L(XX'y, NN'y) \leq -\frac{1}{2} \sum_{i=1}^p \ln \frac{4\lambda_i \lambda_{n-i+1}}{(\lambda_i + \lambda_{n-i+1})^2} \quad (12)$$

其中 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$ 是 Γ 的特征根。

根据 [1] ((9.7.16) 式) 给出的两个非奇异正态变量之间的相关性度量, 不难得到

$$L(X'y, N'y) = -\frac{1}{2} \ln \prod_{i=1}^p (1 - \gamma_i^2) \quad (13)$$

其中 $\gamma_1^2 \geq \gamma_2^2 \geq \dots \geq \gamma_p^2$ 是 $X'y$ 和 $N'y$ 之间的典型相关系数。而它们又与 $XX'y$ 和 $NN'y$ 之间的典型相关系数是一致的。故可得

定理 2 对于线性模型 (1), $N'y$ 和 $X'y$ 之间的相关性度量与估计效率的关系为

$$L(X'y, N'y) = -\frac{1}{2} \ln e(\hat{\beta}) \quad (14)$$

类似于推论1有

推论2 在定理2条件之下, 有

$$0 \leq L(X'y, N'y) \leq -\frac{1}{2} \sum_{i=1}^p \ln \frac{4\lambda_i \lambda_{n-i+1}}{(\lambda_i + \lambda_{n-i+1})^2} \quad (15)$$

参 考 文 献

- [1] Kullback S. Information theory and statistics. New York, John Wiley, 1959, 8
- [2] 王松桂. 广义相关系数与估计效率. 科学通报, 1985, 30(19), 1521~1529
- [3] Rao C R. Linear statistical inference and its applications. New York, John Wiley, 1973, 527
- [4] Bartmann F C, Bloomfield P. Inefficiency and correlation. Biometrika, 1981, 68(1), 67~11
- [5] Bloomfield P, Watson G S. The inefficiency of least squares. Biometrika, 1975, 62(1), 121~128

CORRELATION MEASURE AND ESTIMATION EFFICIENCY

Liu Lanting Zhang Xindong

Abstract

In this paper, We give out the correlation measure relationship between the least square fitted values and the residuals to the relative efficiency of the least square estimation.

Key words Correlation measure Estimation efficiency