

基于 C_p 统计量的自变量选择原则

The Independent Variable Selection Rule Based on the C_p Statistics

周婉枝

Zhou Wanzi

陈宇丹

Chen Yudan

(广西大学数学与信息科学系,
南宁市西乡塘路 530004)(Dept. of Math. & Inf's Sci. Guangxi University,
10 Xixiangtang Road, Nanning, Guangxi, 530004)

(广西医科大学基础部)

南宁市滨湖路6号 530021)

(Basics Dept., Guangxi Medical University,
6 Binhu Road, Nanning, 530021)

摘要 对回归模型 $y = \beta_0 + \beta_1 x_1 + \dots + \beta_t x_t$ (其中 y 是 $q \times 1$ 随机向量, β_i 为 $q \times 1$ 维参数向量), 提出在 $q \geq 1$ 的情况下, 基于 C_p 统计量的自变量的选择原则: 选择自变量是子集 P , 使其相对应的 C_p 值满足条件 $C_p \leq \gamma_a(t, n - t - 1, q)$.

关键词 选择原则 C_p 统计量 置信椭圆

Abstract For the regression model $y = \beta_0 + \beta_1 x_1 + \dots + \beta_t x_t$, where y is $q \times 1$ random vector and β_i is $q \times 1$ parameter vector, the independent variable selection rule based on the C_p statistics under the condition of $q \geq 1$: The selection of the independent variable subset P can be done through the calculation of C_p corresponding to P , which satisfied $C_p \leq \gamma_a(t, n - t - 1, q)$.

Key words selection rule, C_p statistics, confidence ellipsoid

1 C_p 统计量

对回归模型 $y = \beta_0 + \beta_1 x_1 + \dots + \beta_t x_t$, 1964 年 C. L Mallows 在 $q = 1$ 的条件下提出了 C_p 统计量及基于 C_p 统计量的自变量选择原则^[1], 本文讨论在 $q \geq 1$ 的条件下自变量的选择原则.

1.1 模型

$$\text{原模型} \begin{cases} Y = XB + e \\ e = (e_{(1)} \dots e_{(n)})' \\ e_{(1)} \dots e_{(n)} \text{ 相互独立且同分布 } Ng(0, \Sigma) \end{cases} \quad (1)$$

这里 Y 为 $n \times q$ 随机观察阵, X 为 $n \times (t + 1)$ 设计矩阵, B 为 $(t + 1) \times q$ 参数矩阵, e 为 $n \times q$ 误差阵选择自变量子集 P , 相应于 P 的资料阵 X .

$$\text{选模型} \begin{cases} Y = X_P B_P + e \\ e = (e_{(1)} \dots e_{(n)})' \\ e_{(1)} \dots e_{(n)} \text{ 相互独立同分布 } Ng(0, \Sigma) \end{cases} \quad (2)$$

1.2 C_p 统计量

在模型(1)真下, 选用模型(2), 文献[2]在损失

函数

$$J_p = \text{tr}[\Sigma^{-1}(Y_P - XB)'(Y_P - XB)]$$

导出了统计量

$$C_p = (2p - t - 1)q + (n - t - q - 2)\text{tr}[(RSS)^{-1}(RSS_P - RSS)]$$

2 主要结果

定义: P 是自变量子集, X 分割为 (X_1, X_R) , 其中 X_1 相应于子集 P 的资料阵. B 相应地剖成 B_1, B_R , 若 $B_R = 0$, 则称 P 为合宜的.

利用此定义, 若 P 合宜, 选用模型(2)是正确的.

2.1 模型中心化

对模型(1), 把常数项单独列出改写为

$$\begin{cases} Y_i = 1\beta_0' + XB_i + e & 1 \text{ 为元素均为 1 的向量} \\ e_{(1)} \dots e_{(n)} \text{ 独立同分布 } Ng(0, \Sigma) \end{cases}$$

令 $\bar{X} = X - (1/n) 11'X = (I - (1/n) 11')X$

令 $\bar{\beta}_0' = \beta_0' + (1/n) 1'XB_i$, 则(1)改写为

$$\begin{cases} Y = 1\bar{\beta}_0' + \bar{X}B_i + e \\ e_{(1)} \dots e_{(n)} \text{ 独立同分布 } Ng(0, \Sigma) \end{cases}$$

此时, 由于 $l' \bar{X} = 0$, 计算得

$$\begin{pmatrix} \hat{\beta}_0' \\ \hat{B}_i \end{pmatrix} = \begin{pmatrix} (1/n) l' Y \\ (\bar{X}' \bar{X})^{-1} \bar{X}' Y \end{pmatrix}.$$

$\hat{B}_i = (\bar{X}' \bar{X})^{-1} \bar{X}' Y$ 相当于从 $Y = \bar{X} B_i + e$ 进行最小二乘估计而得.

2.2 回归系数 B_i 的置信椭圆

设 $A_n \sim Wq(n, I)$, $D \sim Wq(m, I)$ A_n 与 D 相互独立

$(n-q-1) \text{tr} A_n^{-1} D$ 的分布为 $r(m, n, q)$

由文献 [2], 可得 $r(m, n, q)$ 的极限分布为 $x_{mq}^2 (n \rightarrow \infty)$.

定理: \hat{B}_i 是模型 (1) 下 B_i 的最小二乘估计, RSS 为相互残差阵, 令 $Q = \bar{X}' \bar{X}$, 则有

$$(n-t-q-2) \text{tr} [(RSS)^{-1} (\hat{B}_i - B_i) Q (\hat{B}_i - B_i)] \sim r(t, n-t-1, q)$$

证: 由 2.1 的结论,

$$\begin{aligned} \hat{B}_i &= (\bar{X}' \bar{X})^{-1} \bar{X}' Y \\ &= B_i + (\bar{X}' \bar{X})^{-1} \bar{X}' e, \text{rank}(\bar{X}) = t \end{aligned}$$

$$\text{故 } (\hat{B}_i - B_i)' Q (\hat{B}_i - B_i) = e' P_X e \sim Wq(t, \Sigma)$$

而 $RSS \sim Wq(n-t-1, \Sigma)$ 且因为 $(I - (1/n) l l' - P_X) P_X = 0$

有 RSS 与 $(\hat{B}_i - B_i)' Q (\hat{B}_i - B_i)$ 独立

$$\text{从而 } (n-t-q-2) \text{tr} [(RSS)^{-1} (\hat{B}_i - B_i)' Q (\hat{B}_i - B_i)] \sim r(t, n-t-1, q)$$

B_i 的置信椭圆

定义: 令 $v = n-t-q-2$

$$S_\alpha = \{B_i: v \cdot \text{tr}[(RSS)^{-1} (\hat{B}_i - B_i)' Q (\hat{B}_i - B_i)] \leq r_\alpha(t, n-t-1, q)\}$$

$$P\{r(t, n-t-1, q) \leq r_\alpha(t, n-t-1, q)\} = 1 - \alpha$$

由 $r(t, n-t-1, q)$ 的渐近分布, n 足够大, $r_\alpha(t, n-t-1, q)$ 可由 x_{α}^2 的百分位点得到.

2.3 自变量子集 P 合宜时与 Cp 的关系

定理: 以下三点等价 对子集 P

(i) P 合宜且 $B_i \in S_\alpha$

(ii) $v \cdot \text{tr} (RSS)^{-1} (RSS_s - RSS) \leq r_\alpha(t, n-t-1, q)$

(iii) $Cp \leq (2p-t-1) q + r_\alpha(t, n-t-1, q)$

先证

$$(a) \text{RSS}_s - \text{RSS} = (B_s^{(0)} - \hat{B}_i)' Q (B_s^{(0)} - \hat{B}_i),$$

其中 $B_s^{(0)} = \begin{pmatrix} \tilde{B}_s^{(-0)} \\ 0 \end{pmatrix}$, $\tilde{B}_s^{(-0)}$ 是指在选模型下 \tilde{B}_s 的最小二乘估计且除去常数项, 对模型中心化后所求的 \tilde{B}_s . 正是没有中心化前的 $\tilde{B}_s^{(-0)}$, 为方便起见, 用中心化后求得的 \tilde{B}_s 代替 $\tilde{B}_s^{(-0)}$, 中心化后 \bar{X} 的仍用 \bar{X} 表

示. 故 $Q = X' X$

$$\text{相应于子集 } P \text{ 记 } Q = \begin{pmatrix} B & C \\ C' & D \end{pmatrix} \quad Q^{-1} = \begin{pmatrix} B_1 & C_1 \\ C_1' & D_1 \end{pmatrix}$$

$$\begin{aligned} \text{由于 } \text{RSS}_s - \text{RSS} &= Y' (P_X - P_{X_1}) Y \\ &= \hat{B}'_R (D - C' B^{-1} C) \hat{B}_R \end{aligned}$$

因 $B_s^{(0)} = (\tilde{B}'_s, 0')'$

$$\begin{aligned} (B_s^{(0)} - \hat{B}_i)' Q (B_s^{(0)} - \hat{B}_i) &= \tilde{B}'_s B \tilde{B}_s + \hat{B}'_R Q \hat{B}_R - \tilde{B}'_s B \hat{B}_R - \tilde{B}'_s C \hat{B}_R \\ &\quad - \hat{B}'_R B \tilde{B}_s - \hat{B}'_R C' \tilde{B}_s \end{aligned}$$

由正规方程 $X = (X_1, X_R)$

$$\begin{cases} B \tilde{B}_s + C \hat{B}_R = X_1 Y \\ C' \tilde{B}_s + D \hat{B}_R = X_R Y \\ B \tilde{B}_s = X_1 Y \end{cases} \Rightarrow \begin{cases} \tilde{B}_s = B^{-1} X_1 Y \\ \hat{B}_R = B^{-1} [X_1' - (D - C' B^{-1} C)^{-1} (X_1' - C' B^{-1} X_1)] Y \\ \hat{B}_R = (D - C' B^{-1} C)^{-1} (X_1' - C' B^{-1} X_1) Y \end{cases}$$

代入 $(B_s^{(0)} - \hat{B}_i)' Q (B_s^{(0)} - \hat{B}_i)$ 各项得

$$(B_s^{(0)} - \hat{B}_i)' Q (B_s^{(0)} - \hat{B}_i) = \hat{B}'_R (D - C' B^{-1} C) \hat{B}_R$$

从而 (a) 成立

定理的证明: (ii) \rightarrow (i), 由 (a),

$$\text{RSS}_s - \text{RSS} = (B_s^{(0)} - \hat{B}_i)' Q (B_s^{(0)} - \hat{B}_i)$$

可知 $B_s^{(0)} \in S_\alpha$, 且 $B_s^{(0)} = (\tilde{B}'_s, 0')'$, 故 P 合宜, 故 (i) 成立.

(i) \Rightarrow (ii), 若 (i) 成立, 则存在 H , $H = (H_s' H_R')'$, $H_R = 0$

且 $v \cdot (n-t-q-2) \text{tr} [(RSS)^{-1} (H - \hat{B}_i)' Q (H - \hat{B}_i)] \leq r_\alpha(t, n-t-1, q)$

因 $B_s^{(0)} = (\tilde{B}'_s, 0')'$, 故

$$\hat{B}'_i X' X B_s^{(0)} = Y' (X_1 \quad X_R) \begin{pmatrix} \tilde{B}_s \\ 0 \end{pmatrix} = Y' X_1 \tilde{B}_s$$

$$B_s^{(0)'} X' X B_s^{(0)}$$

$$= (Y' X_1 B^{-1}, 0) \begin{pmatrix} B & C \\ C' & D \end{pmatrix} \begin{pmatrix} B^{-1} X_1' Y \\ 0 \end{pmatrix}$$

$$= Y' X_1 \tilde{B}_s$$

所以 $\hat{B}'_i X' X B_s^{(0)} = B_s^{(0)'} X' X B_s^{(0)}$

$$\hat{B}'_i X' X H = Y' X_1 H$$

$$B_s^{(0)'} X' X H = Y' X_1 H$$

从而 $(\hat{B}_i - B_s^{(0)})' Q (B_s^{(0)} - H) = 0$

$$(\hat{B}_i - H)' Q (B_i - H)$$

$$= (\hat{B}_i - B_s^{(0)})' Q (\hat{B}_i - B_s^{(0)})$$

$$+ (B_s^{(0)} - H)' Q (B_s^{(0)} - H)$$

(下转第 68 页 Continue on page 68)

表 6 龙州站洪峰作业预报登记表

Table 6 Record of flood peak forecast from Longzhou station

发布时间 Broadcast time	预报值 Forecast		实测值 Measure		预见期内 水位变幅 Stage fluctuation during forecast period (m)	洪水 总涨幅 Total fluctuation of flood (m)	发预报 时水位 Stage at broad cast time (m)	预 报 准确率 Forecast accuracy (%)	有 效 预见期 Useful forecast period (h)	起涨 Start to rise	
	出现时间 Peak time	水 位 Stage	出现时间 Peak time	水 位 Stage						时间 time	水 位 Stage
	d-h	(m)	d-h : min	(m)						d-h	(m)
1971-07-23-20	24-23	119.50	25-07 : 00	120.32	11.08	13.42	108.42	82.6	35.0	23-09	106.9
1971-08-29-20	30-12	116.00	30-16 : 0	116.68	6.68	9.20	109.32	73.9	20.0	29-08	107.4
1975-09-31-20	02-02	115.80	02-01 : 01	116.24	6.58	8.82	109.28	74.6	39.0	31-06	107.4
1978-05-17-05	18-08	116.00	18-16 : 00	116.83	7.24	9.41	108.76	76.9	35.0	16-11	106.0
1980-07-24-05	25-08	120.00	25-12 : 00	120.21	7.24	9.91	112.76	73.1	31.0	23-20	110.30
1985-09-11-22	12-20	117.50	13-03 : 00	118.17	5.25	7.88	112.25	66.6	29.0	11-4	110.2
1986-07-23-05	24-12	122.00	25-01 : 30	125.89	11.23	16.93	110.77	66.3	44.5	22-17	108.9
1986-07-23-20	24-13	125.00	25-01 : 30	125.89	5.55	6.44	119.45	86.2	29.5	22-17	108.9

的预报难题。以上预报方案的创立，开创了探讨上游无情报提供的洪水预报科研的良好开端，特别是中小流域，除了可借用毗邻流域的水文，雨量资料进行本河段预报外，尚可以直接应用本站的实测水位、流量、雨量资料进行本站洪峰的预报尝试，解决上游无情报的困难。

《上游无情报的水文预报技术》被江苏扬州水利专科学校李慧珑教授称为“教科书没有的、行之有效

的方法”

参考文献

- 1 扬振怀等. 中国水利百科全书. 北京: 水利电力出版社, 1990. 12.
- 2 华东水利学院水文系编. 水文预报. 中国工业出版社, 1962. 8.

(责任编辑: 邓大玉 蒋汉明)

(上接第 18 页 Continue from page 18)

$$\begin{aligned}
 & v \cdot \text{tr} [(RSS)^{-1} (RSS_0 - RSS)] \\
 &= v \cdot \text{tr} (RSS)^{-1} (\hat{B}_i - B_0^{(0)})' Q (\hat{B}_i - B_0^{(0)}) \\
 &\leq v \cdot \text{tr} (RSS)^{-1} (\hat{B}_0 - B_0^{(0)})' Q (\hat{B}_i - B_0^{(0)}) \\
 &\quad + (B_0^{(0)} - H)' Q (B_0 - H) \\
 &= v \cdot \text{tr} (RSS)^{-1} (\hat{B}_i - H)' Q (\hat{B}_i - H) \\
 &\leq r_\alpha(t, n-t-1, q)
 \end{aligned}$$

故 (ii) 成立

(ii) ⇔ (iii) 由 C_p 的定义直接得到 P 合宜与 C_p 的关系

定理的意义在于 (i) ⇔ (ii)，因为 (i) 的含义在于“在置信椭圆 S_α 中包含了一个形如 $(B_0' \quad 0)'$

的矩阵”，表明从参数 B_t 的区域估计的观点看，“弃后 r 个自变量”与数据是相容（当然在所给的置信系数下），(i) 与 (iii) 的等价性正是把这个事实归于 C_p 的计算，因此，选择 P ，首先应注意满足 $C_p \leq (2p-t-1)q + r_\alpha(t, n-t-1, q)$ 的子集。

参考文献

- 1 Mallows C L. Some Commerts on C_p . Thehnometrics, 1964, 15: 661~675.
- 2 周婉枝. C_p 统计量. 广西大学学报, 1995, 20 (3): 291~294.

(责任编辑: 莫鼎新 邓大玉)