

线性模型和广义线性模型中的一种 强影响点的显著性检验法 Influential Point Test for Linear Model and Generalized Linear Model

龙蓓 林路*
Long Bei Lin Lu

(桂林医学院数学教研室 桂林市乐群路 56 号 541001)

(Teaching and Research Section of Math., Guilin Medical College, 56 Lequnlu, Guilin, Guangxi, 541001)

摘要 根据强影响点的实际意义, 提出一种强影响点的显著性检验模型。解决了线性模型和广义线性模型的强影响点的显著性检验问题, 其中的检验统计量分别是 F 检验统计量和 Score 检验统计量, 实例表明此法较好。

关键词 统计诊断 线性模型 广义线性模型 强影响点

中图法分类号 O 212.1

Abstract The influential point test for the linear model and the generalized linear model are discussed. We give the influential point test model and F statistic for the linear model and score statistic for the generalized linear model. Two numerical examples are given to illustrate our results.

Key words statistical diagnostic, linear model, generalized linear model, influential point

影响分析是统计诊断中十分活跃的分支, 其初期的最有实用价值的内容就是研究某些特定点对统计分析(如参数估计)的影响。虽然线性模型中强影响点的显著性检验已有文章作过探讨, 但对广义线性模型的强影响点的显著性检验问题却研究甚少。作者根据强影响点的实际意义, 提出了一种强影响点的显著性检验模型, 解决了线性模型和广义线性模型的强影响点的显著性检验问题, 其中的检验统计量分别是 F 检验统计量和 Score 检验统计量, 实例表明此法较好。

1 检验模型

本文考察的线性模型为

$$Y = X\beta + \epsilon, \quad \epsilon \sim N(0, \sigma^2 I) \quad (1)$$

其中 $Y = (y_1, y_2, \dots, y_n)$ 是观测向量, $X = (x_1, x_2, \dots, x_n)^T$ 是 $n \times p$ 列满秩设计阵, $\beta = (\beta_1, \beta_2, \dots, \beta_p)^T$ 为 $p \times 1$ 未知参数, $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^T$ 为 $n \times 1$ 随机误差向量, β 在模型(1)中的最小二乘估计(最

大似然估计)为 $\hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p)^T = (X^T X)^{-1} X^T Y$,

在模型(1)中, 删除 $J = \{i_1, i_2, \dots, i_k\}$ 中数据, $J \subset \{1, 2, \dots, n\}$, 得到数据删除模型。

$$Y(J) = X(J)\beta + \epsilon(J), \quad \epsilon(J) \sim N(0, \sigma^2 I(J)) \quad (2)$$

其中 $Y(J), X(J), \epsilon(J), I(J)$ 分别是 Y, X, ϵ, I 删除第 i_1 行, 第 i_2 行, \dots , 第 i_k 行后得到的向量或矩阵, β 在模型(2)中的最小二乘估计(最大似然估计)为:

$$\hat{\beta}(J) = [\hat{\beta}(J)_1, \hat{\beta}(J)_2, \dots, \hat{\beta}(J)_p]^T = (X^T(J)X(J))^{-1} X^T(J)Y(J).$$

本文还研究实用上较常用的典则联系的广义线性模型。

$$E(Y) = \mu, \quad \theta = G(\mu) = X\beta \quad (3)$$

其中 Y, X 和 β 与模型(1)中的表示式一致, $\mu = (\mu_1, \mu_2, \mu_3, \dots, \mu_n)^T, \theta = (\theta_1, \theta_2, \dots, \theta_n)^T = [g(\mu_1), g(\mu_2), \dots, g(\mu_n)]^T$, y_i 服从指数族分布, 其分布密度为:

$$p(y_i, \theta_i) = \exp\left\{\frac{\theta_i y_i - \psi(\theta_i)}{a(\varphi)} + c(y_i, \varphi)\right\},$$

$i = 1, 2, \dots, n$. 其中 φ 为多余参数, 因此假设 $a(\varphi) = 1$, g 为一元严格增函数, 即为 ψ 的反函数, 记模型(3)中 β 的最大似然估计为 $\hat{\beta}$, $\hat{\beta}$ 可由 $G-N$ 迭代公式求解。

1997-10-20 收稿。

*邵阳师范专科学校数学系, 湖南邵阳, 422000 (Department of Math., Shaoyang Teachers College, Shaoyang, Hunan, 422000)

在模型(3)中, 删除 $J = \{i_1, i_2, \dots, i_k\}$ 中数据, 得到其数据删除模型.

$$E(Y(J)) = \mu(J), \theta(J) = G(J)(\mu) = X(J)\beta \quad (4)$$

其中 $Y(J), \mu(J), \theta(J), G(J), X(J)$ 分别是 Y, μ, θ, G, X 删除第 i_1 行, 第 i_2 行, \dots , 第 i_k 行后得到的向量或矩阵, 模型(4)中 β 的最大似然估计记为 $\beta(J)$.

记 $Y_J = (y_{i_1}, y_{i_2}, \dots, y_{i_k})^T, X_J = (x_{i_1}, x_{i_2}, \dots, x_{i_k})^T$, 在影响分析中, 常用 β 与 $\beta(J)$ 的某种距离描述数据 (Y_J, X_J) 对回归分析(如参数估计)的影响. 例如, 对参数估计来说, 若 $\beta(J)$ 与 β 的距离很大, 则数据 (Y_J, X_J) 为强影响点; 若 $\beta(J)$ 与 β 相去不远, 则 (Y_J, X_J) 不是强影响点, 基于这种思想, 本文考察原模型(1)(或(3))和数据删除模型(2)(或(4))中未知参数 β 的真值之间的差异^[1], 若删除数据 (Y_J, X_J) 后, 模型(2)(或(4))中的 β 的真值与模型(1)(或(3))中 β 的真值差异很大, 则数据 (Y_J, X_J) 为强影响点, 否则, (Y_J, X_J) 不是强影响点. 在模型(1)(或(3))中, 用 β 近似地表示 β 的真值, 我们就有如下想法: 若模型(2)(或(4))中 β 的真值与 β 差异很大, 则 (Y_J, X_J) 是强影响点, 否则 (Y_J, X_J) 不是强影响点.

由上述思想, 我们在数据删除模型(2)和(4)中提出如下假设检验:

$$H: \beta = \beta \quad (5)$$

若原假设 H 被拒绝, 则认为模型(2)(或(4))中的 β 的真值与模型(1)(或(3))中的 β 有显著差异, 从而数据 (Y_J, X_J) 是强影响点.

2 检验统计量

先求线性模型中强影响点的检验统计量, 显然, 检验(5)是模型(2)中 β 的线性假设检验.

模型(2)的残差平方和为:

$$RSS(J) = \sum_{j \in J} (y_j - x_j^T \beta(J))^2 \quad (6)$$

模型(2)中, 在条件 $H: \beta = \beta$ 下的残差平方和为:

$$RSS(J)_H = \sum_{j \in J} (y_j - x_j^T \beta)^2 \quad (7)$$

记模型(1)的残差平方和为:

$$RSS = \sum_{j=1}^n (y_j - x_j^T \beta)^2 \quad (8)$$

且记:

$$\hat{e}_j = y_j - x_j^T \beta, \hat{e}_j = (\hat{e}_{i_1}, \hat{e}_{i_2}, \dots, \hat{e}_{i_k})^T, P_J = X_J(X_J^T X_J)^{-1} X_J^T$$

我们可得到假设检验(5)的 F 检验统计量^[3]:

$$\begin{aligned} F_H(J) &= \frac{1}{p} (RSS(J)_H - RSS(J)) / \\ & \frac{1}{n-p-k} RSS(J) \\ &= \frac{n-p-k}{p} \cdot \frac{RSS - \hat{e}_j^T \hat{e}_j - (RSS - \hat{e}_j^T (I - P_J)^{-1} \hat{e}_j)}{RSS - \hat{e}_j^T (I - P_J)^{-1} \hat{e}_j} \\ &= \frac{n-p-k}{p} \cdot \frac{\hat{e}_j^T ((I - P_J)^{-1} - I) \hat{e}_j}{RSS - \hat{e}_j^T (I - P_J)^{-1} \hat{e}_j} \\ &= \frac{n-p-k}{p} \cdot \frac{\hat{e}_j^T (I - P_J)^{-1} P_J \hat{e}_j}{RSS - \hat{e}_j^T (I - P_J)^{-1} \hat{e}_j}, \end{aligned} \quad (9)$$

于是, 我们有如下定理.

定理 1 在模型(2)中, 若 $H: \beta = \beta$ 为真, 则:

$$F_H(J) = \frac{n-p-k}{p} \cdot \frac{\hat{e}_j^T (I - P_J)^{-1} P_J \hat{e}_j}{RSS - \hat{e}_j^T (I - P_J)^{-1} \hat{e}_j} \sim F(p, n-p-k) \quad (10)$$

其中, $F(p, n-p-k)$ 为自由度为 $(p, n-p-k)$ 的 F 分布.

由定理 1 知, 给定显著性水平 $\alpha, 0 < \alpha < 1$, 记 $F_\alpha(p, n-p-k)$ 是自由度为 $(p, n-p-k)$ 的 F 分布的上侧 α 分位点, 当 $F_H(J) > F_\alpha(p, n-p-k)$ 时, 不定式(5), 即认为 (Y_J, X_J) 是强影响点; 不然就接受(5), 即认为 (Y_J, X_J) 不是强影响点.

由 $F_H(J)$ 的表达式(9)知, $F_H(J)$ 与 Cook 统计量和 $W-K$ 统计量等度量影响的统计量有相当的一致性^[1]. 在本文后面的实例中, 将再次证实这种一致性.

下面将定理 1 推广(或者说应用)到约束线性模型强影响点的显著性检验.

考察约束线性模型

$$\begin{cases} Y = X\beta + \epsilon \\ A^T \beta = 0 \end{cases} \quad \epsilon \sim N(0, \sigma^2 I) \quad (11)$$

其中 A 为 $p \times q$ 列满秩矩阵, $q < p$, 此模型的数据删除形式为

$$\begin{cases} Y(J) = X(J)\beta + \epsilon(J) \\ A^T \beta = 0 \end{cases} \quad \epsilon(J) \sim N(0, \sigma^2 I(J)) \quad (12)$$

在模型(11), (12)中, 记约束 $A^T \beta = 0$ 的解为 $\beta = Bt$, B 是 $p \times (p-g)$ 列满秩矩阵, 满足 $A^T Bt = 0$, 则分别得到与(11)和(12)等价的无约束线性模型^[4]

$$Y = (XB)t + \epsilon \quad \epsilon \sim N(0, \sigma^2 I) \quad (13)$$

和 $Y(J) = (X(J)B)t + \epsilon(J)$,

$$\epsilon(J) \sim N(0, \sigma^2 I(J)) \quad (14)$$

若记(11)中 β 的最小二乘估计为 $\hat{\beta}_c$, 残差为 $\hat{e}_j(c) = y_j - x_j^T \hat{\beta}_c$, (13)中 t 的最小二乘估计为 \hat{t} , 要检验数据 (Y_J, X_J) 是否为模型(11)即模型(13)的强影响点,

只要对模型(14) 检验如下假设:

$$H_c: t = t \quad (15)$$

由式(9) 知, 检验(15) 的 F 统计量为

$$F_{H_c}(J) = \frac{n-p+q-k}{p-q} \frac{\hat{e}_j(c)(I-P_J(c))^{-1}P_J(c)\hat{e}_j(c)}{RSS_c - \hat{e}_j(c)(I-P_J(c))^{-1}\hat{e}_j(c)} \quad (16)$$

其中 RSS_c 为(11) 即(13) 的残差平方和 $\sum_{j=1}^n (y_j - x_j^T \beta)^2$ $\hat{e}_j(c) = (\hat{e}_{j1}(c), \hat{e}_{j2}(c), \dots, \hat{e}_{jk}(c))^T$, 由文献[3] 知:

$$P_J(c) = X_J B (B^T (X^T X) B)^{-1} B X_J^T = X_J ((X^T X)^{-1} - (X^T X)^{-1} M (X^T X)^{-1}) X_J^T$$

其中 $M = A (A^T (X^T X)^{-1} A)^{-1} A$, 由此得到如下定理.

定理 2 模型(12)(即模型(14) 中, 若式(15) 为真, 则:

$$F_{H_c}(J) = \frac{n-p+q-k}{p-q} \frac{\hat{e}_j(c)(I-P_J(c))^{-1}P_J(c)\hat{e}_j(c)}{RSS_c - \hat{e}_j(c)(I-P_J(c))^{-1}\hat{e}_j(c)} \sim F(p-q, n-p+q-k) \quad (17)$$

定理 2 表明, 若 $F_{H_c}(J) > F_{\alpha}(p-q, n-p+q-k)$, 则数据 (Y_J, X_J) 是模型(11) 的强影响点; 不然, (Y_J, X_J) 不是(11) 的强影响点.

下面求广义线性模型中强影响点的检验统计量.

广义线性模型(4)(略去与参数无关的项后) 的对数似数函数为:

$$L(\beta) = \sum_{j \in J} (y_j x_j^T \beta - \psi(x_j^T \beta)) \quad (18)$$

$$\text{且: } L(\beta) = X^T(J) S(J) \quad (19)$$

其中 $S(J)$ 是由 $S_j, j \in \{1, 2, \dots, n\} - J$ 组成的列向量, $S_j = y_j - \psi(\theta_j) = y_j - \mu_j$

$$\dot{L}(\beta) = -X^T(J) W(J) X(J) \quad (20)$$

其中 $W(J)$ 是 $W_j, j \in \{1, 2, \dots, n\} - J$ 构成的对角阵, $W_j = \text{Var}(y_j) = \psi''(\theta_j)$.

以下, 若无特殊说明, $S(J), W(J)$ 及 $S = (s_1, s_2, \dots, s_n)^T, S_J = (s_{j1}, s_{j2}, \dots, s_{jk})^T, W = \text{diag}(w_1, w_2, \dots, w_n), W_J = \text{diag}(w_{j1}, w_{j2}, \dots, w_{jk})$ 等均表示在 β 处计值. 我们知道^[3], 模型(4) 对于假设检验(5) 的 Score 检验统计量为:

$$SC_J = -(\dot{L}(\beta))^T \dot{L}^{-1}(\beta) \dot{L}(\beta) \rightarrow \chi^2(p) \quad (21)$$

由式(19), (20) 有:

$$SC_J = S^T(J) X(J) (X^T(J) W(J) X(J))^{-1} X^T(J) S(J) \quad (22)$$

由于在 β 处有 $X^T S = 0$ (即 β 满足模型(3) 的似然方程), 故

$$X^T(J) S(J) = X^T S - X_J^T S_J = -X_J^T S_J \quad (23)$$

又由于

$$X^T(J) W(J) X(J) = X^T W X - X_J^T W_J X_J \quad (24)$$

由式(22), (23), (24) 及矩阵和式求逆公式, 得到

$$SC_J = S_J^T X_J (X^T W X)^{-1} X_J^T (I - W_J X_J (X^T W X)^{-1} X_J^T)^{-1} S_J \quad (25)$$

所以, 有

定理 3 在模型(4) 中, 若 $H: \beta = \beta$ 为真, 则 Score 检验统计量

$$SC_J = S_J^T X_J (X^T W X)^{-1} X_J^T (I - W_J X_J (X^T W X)^{-1} X_J^T)^{-1} S_J \rightarrow \chi^2(p) \quad (26)$$

定理 3 表明, 若 $SC_J > \chi_{\alpha}^2(p)$, 则认为数据 (Y_J, X_J) 是模型(3) 中的强影响点; 否则, 数据 (Y_J, X_J) 不是模型(3) 中的强影响点.

3 算例

下面的两个算例是分别应用定理 1 和定理 3 检验线性模型及广义线性模型中单个数据是否是强影响点, 若逐个考察单个数据 (y_j, x_j^T) 的强影响点判别问题, 式(9) 为

$$F_H(j) = \frac{n-p-1}{p} \frac{P_{jj} e_j^2 / (1-P_{jj})}{RSS - e_j^2 / (1-P_{jj})}$$

其中, $P_{jj} = x_j^T (X^T X)^{-1} x_j$, 为便于应用, 将 $F_H(j)$ 用 Cook 统计量和 $W-K$ 统计量表示:

$$F_H(j) = \frac{(1-P_{jj}) \sigma^2}{\sigma^2(j)} D_j \quad (27)$$

$$F_H(j) = \frac{1-P_{jj}}{p} (W_{k_j})^2 \quad (28)$$

其中 $D_j = \frac{P_{jj} e_j^2}{P \sigma^2 (1-P_{jj})^2}$ 是 (y_j, x_j^T) 的 Cook 统计量, $W-k_j = \frac{\sqrt{P_{jj}} \cdot e_j^2}{\sigma^2(j) (1-P_{jj})}$ 是 (y_j, x_j^T) 的 $W-K$ 统计量, $\sigma^2 = \frac{RSS}{n-p}, \sigma^2(j) = \frac{RSS(j)}{n-p-1}$,

对于单个数据 (y_j, x_j^T) , 记 $h_{jj} = W_j x_j^T (X^T W X)^{-1} x_j$, 式(25) 为:

$$SC_J = \frac{h_{jj} s_j^2}{W_j (1-h_{jj})} = \frac{h_{jj}}{1-h_{jj}} \cdot r_{pj} \quad (29)$$

其中 r_{pj} 为 (y_j, x_j^T) 的 Pearson 残差: $r_{pj} = s_j w_j^{-1/2}$,

例 1 在 BOQ 数据^[5] 中, 已知 $n = 25$, 对 25 组数据进行线性回归分析 ($p = 8$), 表 1 给出每个数据 (y_j, x_j^T) 的 Cook 统计量, $W-K$ 统计量, F 统计量的值是由式(28) 算得的, 该表说明, 第 23 号点的 D_j 和 $W-k_j$ 都特别大, 若用 D_j 和 $W-k_j$ 为标准, 可以怀疑第 23 号点是强影响点^[1], 若用 F 统计量检验, 只有第 23 号点的 F 统计量的值 $F_H(j) = 3.3839 > F_{0.05}(8, 16) = 2.59$, 即强影响点的显著性检验表明, 只有第 23 号点是强影响点, 检

验结果与文献 [2] 的影响分析一致。

表 1 BOQ 数据的影响度量

Table 1 BOQ data influence measurings

No.	P_{jj}	$W - k_j$	D_j	$F_H(j)$
1	0.257 3	-0.043 3	0.000	0.000 2
2	0.160 9	-0.031 8	0.000	0.000 1
3	0.161 4	-0.201 6	0.005	0.004 3
4	0.163 1	-0.077 8	0.001	0.000 6
5	0.147 5	-0.174 5	0.004	0.003 2
6	0.158 9	-0.247 9	0.008	0.006 5
7	0.182 9	0.356 3	0.016	0.013 0
8	0.359 1	-0.353 3	0.016	0.010 0
9	0.288 0	0.202 9	0.005	0.003 7
10	0.129 5	0.035 3	0.000	0.000 1
11	0.124 1	0.584 9	0.039	0.037 5
12	0.202 4	0.223 0	0.007	0.005 0
13	0.080 2	-0.053 7	0.000	0.000 3
14	0.096 9	-0.007 3	0.000	0.000
15	0.557 6	-2.828 2	0.761	0.442 3
16	0.402 4	-0.044 4	0.000	0.000 1
17	0.368 2	-1.016 2	0.123	0.081 6
18	0.446 5	1.128 6	0.154	0.088 1
19	0.088 6	0.310 6	0.012	0.011 0
20	0.366 3	-2.178 7	0.417	0.376 0
21	0.070 4	0.565 2	0.034	0.037 1
22	0.785 4	-3.071 5	1.079	0.253 1
23	0.988 5	-48.517 9	115.041	3.383 9
24	0.876 2	8.537 3	5.889	1.127 9
25	0.546 7	0.472 2	0.029	0.012 6

例 2 表 2 是一组人造的 Logistic 回归数据^[4], 并给出 Pearson 残差和 Score 检验统计量, 结果显示,

只有第 11 号数据的 $S_{C_j} = 2.061 3 > \chi_{0.20}^2(1) = 1.642$, 即强影响点的显著性检验表明, 只有第 11 号数据是强影响点, 这结果与文献 [2] 中用其他统计量分析的结果是一致的。

表 2 Logistic 回归数据影响度量

Table 2 Influence measurings of Logistic regression data

No.	x_j	Logit (y_j/n_j)	rp_j	h_{jj}	S_{C_j}
1	1	0.5	-0.506	0.201 8	0.064 7
2	2	0.5	-0.615	0.211 5	0.101 5
3	3	0.8	-0.247	0.159 1	0.011 5
4	4	1.0	-0.055	0.122 7	0.000 4
5	5	1.3	0.248	0.101 1	0.006 9
6	6	1.3	0.157	0.093 0	0.002 5
7	7	1.6	0.417	0.097 0	0.018 7
8	8	1.8	0.535	0.111 8	0.036 0
9	9	2.1	0.715	0.136 0	0.080 4
10	10	2.1	0.642	0.168 4	0.083 5
11	17	1.0	-1.383	0.518 7	2.061 3

参考文献

- 1 林 路. 若干有偏估计的强影响点的显著性检验. 数学的实践与认识, 1997, 27 (3).
- 2 韦博成, 鲁国斌, 史建清. 统计诊断引论. 南京: 东南大学出版社, 1990.
- 3 陈希孺, 王松桂. 近代回归分析. 合肥: 安徽教育出版社, 1987.
- 4 虞克明, 王静龙. 约束线性模型异常值检验. 高校应用数学学报, 1994, 9 (4).
- 5 Myers R H. Classial and modern regression with application. Boston: Duxbury press, 1986.
- 6 Cox D R, Hinkley D V. Theoretical statistics London: Chapman and Hall, 1974.

(责任编辑: 黎贞崇 邓大玉)

德国环境、生态、气候的优先研究领域

(1) 实施新的能源研究计划——“第四 能源研究计划”。研究重点: 提高煤和其他矿物能源的利用率并减少对环境的污染, 可再生能源的利用, 核能源研究以及能源的信息处理和分析系统。计划实施可使德国的二氧化碳排放量减少 35%。

(2) 制定新的环境研究计划。研究重点: 生态研究、环境技术研究和大气研究。计划实施可持续发展生态保护系统和集约型生产的环保技术, 并建立相应的示范样板。

(3) 建立一个通过人工智能并且实现了广泛联网的、具有最大灵活性和最佳效果的交通系统, 减少资源消耗和环境负荷。

(摘自中国科学院《科学发展报告》1997. P51)