

半参数估计的自然样条函数法

吴云¹ 孙海燕¹ 马学忠¹

(1 武汉大学测绘学院, 武汉市珞喻路129号, 430079)

摘要: 用补偿最小二乘原理, 得到了参数和非参数分量的惟一解, 并通过模拟计算, 对半参数回归模型和参数模型的计算结果进行了比较。结果表明, 半参数回归方法能较好地地将观测值中具有连续光滑特性的系统误差分离出来。

关键词: 半参数回归; 系统误差; 自然样条函数; 平滑参数; 补偿最小二乘原理

中图法分类号: P207.2

在数据处理中, 观测值中的系统误差或参数模型与实际情况的偏差是影响参数估计精度的主要因素。例如, 在GPS长基线测量中, 双差观测值中仍然存在电离层和对流层折射误差的影响^[1], 如果在参数模型中不顾及这些系统误差, 将会严重地影响到平差精度, 甚至会导致错误的结果^[2]。观测值中的系统误差往往随着某些因素的连续变化, 有时其函数特性十分复杂, 难以用少量参数准确地表达。自然样条函数^[3]属于无穷维函数空间的光滑曲线, 用来描述连续变化的模型误差有其独特的灵活性和简洁性。本文用自然样条函数来表示随时间连续变化的系统误差, 即用参数在表示观测值中可线性化部分的基础上, 增加用自然样条函数表示随时间连续变化的系统误差部分, 也就是非参数分量部分, 从而构成半参数模型。为了求得参数与半参数分量的惟一解, 引入了补偿最小二乘原理 (penalised least squares, PLS)^[4]。本文的算例验证了自然样条函数的半参数模型和补偿最小二乘原理在数据处理中的可行性。

1 数学模型及其解算

1.1 半参数回归模型中的自然样条插值函数

半参数回归模型^[5]为:

$$L = Bx + S + \Delta \quad (1)$$

式中, L 为 n 维观测向量; x 为 t 维参数向量, t 为必要观测数; Δ 为 n 维偶然误差向量; B 为列满秩设计矩阵; 非参数向量 $S = [s_1, s_2, \dots, s_n]^T$,

是描述系统误差的 n 维向量, 其中 s_i 是某些特定量 t_i 的函数。式(1)的误差方程为:

$$V = Bx + S - L \quad (2)$$

式中, $S = (\hat{s}_1, \hat{s}_2, \dots, \hat{s}_n)^T$, 即系统误差 S 的估值。

设 $s(t)$ 为区间 $[t_1, t_n]$ 上的自然样条插值函数, $t_i (i = 1, 2, \dots, n)$ 为节点, 且 $t_1 < t_i < t_n$ 。 $s(t)$ 满足插值条件:

$$s(t_i) = s_i \quad (3)$$

可以找到惟一的满足上述条件的自然样条插值函数^[3], 因此, 观测方程为:

$$L_i = B(t_i)x + s(t_i) + \Delta_i, i = 1, 2, \dots, n \quad (4)$$

1.2 补偿最小二乘原理及其解

补偿最小二乘原理是在最小二乘法的目标函数上增加一个包含非参数部分的补偿项, 即

$$\text{sum} = \sum_{i=1}^n (L_i - Bx - s(t_i))^2 + \alpha \int_{t_1}^{t_n} (s''(t))^2 dt = \min \quad (5)$$

式(5)的前一项是残差平方和, 后一项是补偿项。一个函数当其一阶导数较小时, 二阶导数与其曲率值很接近, 且曲率小, 在几何上理解为“平滑”, 而自然样条插值函数是最光滑的曲线插值函数^[3], 所以, $\int_{t_1}^{t_n} (s''(t))^2 dt$ 刻画了 $s(t)$ 的光滑程度; $\alpha > 0$ 称为光滑参数, 以平衡拟合和光滑程度。如果拟合程度要求较高, 则其光滑程度较差, 反之亦然, 因此, 对 α 要合理地选择。

补偿最小二乘原理的补偿项可以表达为^[4]:

$$\alpha \int_{t_1}^{t_n} (s''(t))^2 dt = \alpha S^T F G^{-1} F^T S \quad (6)$$

式中, $F = (f_{ij})$ 和 $G = (g_{ij})$ 分别是 $n \times (n-2)$ 和 $(n-2) \times (n-2)$ 阶矩阵, 其矩阵中元素的值由 t_i 之间的间隔决定。设 $h_i = t_{i+1} - t_i (i = 1, 2, \dots, n-1)$, 则

$$f_{ij} = \begin{cases} h_j^{-1}, & i = j \\ -(h_j^{-1} + h_{j+1}^{-1}), & i = j + 1 \\ h_{j+1}^{-1}, & i = j + 2 \\ 0, & \text{其他} \end{cases} \quad (7)$$

$$g_{ij} = \begin{cases} 1/3(h_{i-1} + h_i), & i = j, j = 2, \dots, n-1 \\ 1/6h_{i+1}, & i = j-1, j = 2, \dots, n-2 \\ 1/6h_i, & i = j+1, j = 1, \dots, n-3 \\ 0, & \text{其他} \end{cases} \quad (8)$$

式中, $j = 1, 2, \dots, n-2; i = 1, 2, \dots, n$ 。由于矩阵 G 为严格的对角占优矩阵, 即正定矩阵, 设 $K = FG^{-1}F^T$, 因此, K 是 $n \times n$ 阶的半正定矩阵。

根据式(5)和式(6), 按照求条件极值的拉格朗日乘数法, 构造如下函数:

$$\Phi = V^T P V + \alpha S^T K S + 2K_r^T (B\hat{x} + S - L - V) \quad (9)$$

式中, P 为对称正定方阵, 是观测值 L 的权; K_r 是拉格朗日常数。分别令 $\frac{\partial \Phi}{\partial V} = 0, \frac{\partial \Phi}{\partial S} = 0$ 及 $\frac{\partial \Phi}{\partial \hat{x}} = 0$, 可得:

$$K_r = P V \quad (10)$$

$$K_r = -\alpha K S \quad (11)$$

$$B^T K_r = 0 \quad (12)$$

将式(10)代入式(12), 并顾及式(2), 得:

$$B^T P B \hat{x} + B^T P S - B^T P L = 0 \quad (13)$$

将式(10)代入式(11), 顾及式(2), 得:

$$P B \hat{x} + (P + \alpha K) S = P L \quad (14)$$

由式(13)和式(14), 得法方程为:

$$\begin{bmatrix} B^T P B & B^T P \\ P B & P + \alpha K \end{bmatrix} \begin{bmatrix} \hat{x} \\ S \end{bmatrix} = \begin{bmatrix} B^T P L \\ P L \end{bmatrix} \quad (15)$$

在满足 $\text{rank}(F^T B) = t$ 的条件下, 法方程系数矩阵可逆, 未知量 \hat{x} 和 S 有惟一解。下面证明在满足 $\text{rank}(F^T B) = t$ 时, 法方程系数矩阵满秩。

构建二次型:

$$f = [x^T S^T] \begin{bmatrix} B^T P B & B^T P \\ P B & P + \alpha K \end{bmatrix} \begin{bmatrix} x \\ S \end{bmatrix} = (B\hat{x} + S)^T P (B\hat{x} + S) + \alpha S^T F G^{-1} F^T S \quad (16)$$

显然, $f \geq 0$, 由于 P 和 G^{-1} 为正定矩阵, 所以, 当 $B\hat{x} + S = 0$ (17)

且

$$F^T S = 0 \quad (18)$$

时, 二次型 $f = 0$ 。现用反证法证明式(17)和式(18)不能同时成立, 即只有 $f > 0$ 。

假设式(17)和式(18)同时成立, 在式(17)等号两边同乘 F^T , 兼顾式(18), 即得到 $F^T B \hat{x} = 0$ 。若 $\text{rank}(F^T B) = t$, 则对于任意向量 $x \neq 0$, 得 $F^T B \hat{x} \neq 0$, 所以在 $\text{rank}(F^T B) = t$ 的条件下, 此结论与假设矛盾, 式(17)和式(18)不能同时成立, 即二次型 $f \neq 0$, 只可能 $f > 0$, 这时法方程系数矩阵可逆, 方程有惟一解。

由法方程式(15)可以直接解得:

$$S = (P + \alpha K)^{-1} (P L - P B \hat{x}) \quad (19)$$

$$\hat{x} = (B^T P B)^{-1} (B^T P L - B^T P S) \quad (20)$$

将式(19)代入式(20), 并令 $M = (P + \alpha K)^{-1} P$, 则

$$(B^T P (I - M) B) \hat{x} = B^T P (I - M) L \quad (21)$$

同样, 由矩阵的逆变^[6]也可以证明, 只要满足 $\text{rank}(F^T B) = t$, 即 $(B P (I - M) B)$ 可逆, 则式(21)有惟一解:

$$\hat{x} = (B^T P (I - M) B)^{-1} (B^T P (I - M) L) \quad (22)$$

由式(19)和误差方程(2)可得到 S 和 V 。

2 算例分析

这里构造一个模拟的平差问题^[2], 分别建立参数模型和非参数模型并求解。

设有线性模型 $Y_1 = Bx$, 取 $x = [2, 3]^T$ 。 $B = (b_{i,j})$ 为 100×2 阶矩阵, $b_{i,1} = i/20, b_{i,2} = (i/20)^2, i = 1, 2, \dots, 100$ 。为了说明半参数估计的自然样条函数法, 假设观测值中含有随时间非周期连续变化的系统误差 $Y_2 = [y_1, y_2, \dots, y_{100}]^T, y_i = 2t_i + 3\sin(2t_i) \circ \sin(6t_i), t_i = \frac{2(i-1)\pi}{100}, i = 1, 2, \dots, 100$, 模拟的系统误差图见图 1。观测值的真值为 $L = Y_1 + Y_2$ 。观测偶然误差 Δ 是由 100 个服从 $N(0, 1)$ 分布的正态随机数组成的列向量, 于是观测值为 $L = Y_1 + Y_2 + \Delta$ 。平差时, 取观测值的权阵 P 为单位阵。

如果不顾及观测值中的系统误差 Y_2 , 建立参数模型并由最小二乘法平差, 其残差图见图 2。从图 2 可以看出, 残差不具备偶然误差的特性, 明显地含有系统误差。参数的估值为:

$$\hat{x} = (B^T PB)^{-1} B^T PL = [1.745, 3.516]^T$$

下面采用自然样条函数半参数模型,按照补偿最小二乘法平差,并取 $\alpha = 9 \times 10^{-6}$,其结果见图3和图4。

从图3可以看出,建立半参数模型并通过补偿

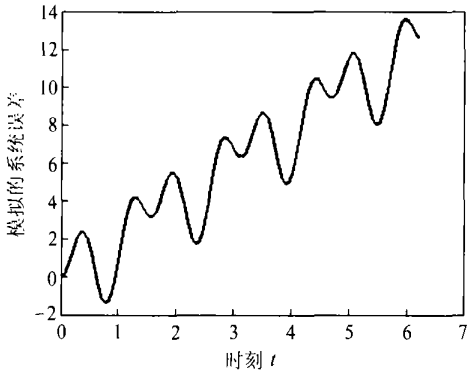


图1 模拟的观测值中的系统误差

Fig. 1 Simulated System Errors

最小二乘原理方法得到的系统误差,正确地反映了系统误差随时间的变化规律,为了解系统误差的性质提供了有价值的信息。图4是半参数模型下的残差,它具有随机性,并接近正态分布。参数估值为 $\hat{x} = [1.821, 3.154]^T$,与真实值也较接近。

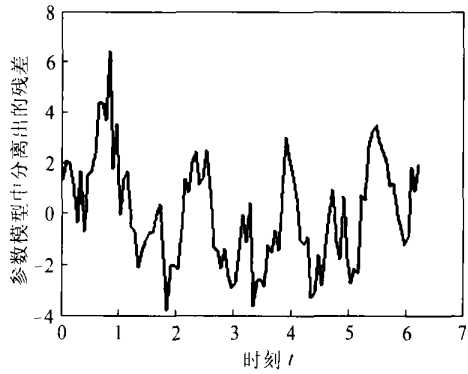


图2 由参数模型解算的残差

Fig. 2 Residual Estimated from Parametric Model

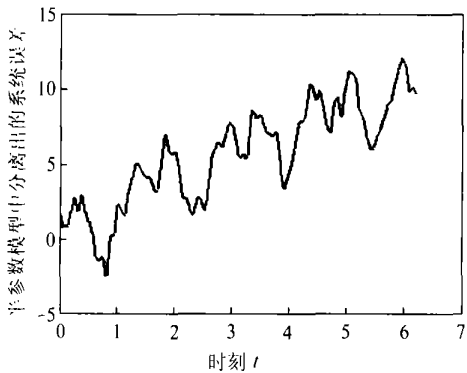


图3 半参数模型中分离出的系统误差

Fig. 3 Estimated System Errors from Semiparametric Model

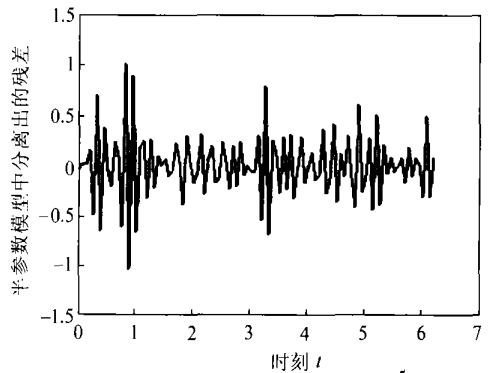


图4 由半参数模型解算的残差

Fig. 4 Residual Estimated from Semiparametric Model

3 结语

观测值中存在系统误差的现象较普遍,如果不加区别地建立 G-M 模型将会导致错误的结论。目前,已有多种对含有系统误差的观测值平差的方法,其解决问题的关键在于建立正确的模型和在此模型下的算法。半参数估计的自然样条函数法是有效地解决半参数估计的方法。自然样条函数具有最光滑等优良特性,能够描述特性复杂的系统误差,用这种方法可以放宽观测的外界条件和设备要求,而且利用数据处理结果可以研究系统误差的特性和规律。

半参数模型下的算法有多种,半参数估计的自然样条函数法只是其中的一种,其中光滑参数 α 的选取存在着客观依据。也就是说,不同的观

测数据对应着客观存在的光滑参数 α ,在这个对应的 α 取值下,可以求得未知量的最佳估值。但光滑参数 α 的具体取值、估计量的统计分析等许多问题,还值得在理论与实践上进行更深入的研究。

参考文献

- 1 Jia M. Mitigation of Systematic Errors of GPS Positioning Using Vector Semiparametric Models. The 13th Int. Tech. Meeting of the Satellite Division of the U. S. Inst. of Navigation, Salt Lake City Utah, 1938~1947
- 2 孙海燕,吴云.半参数回归与模型精化.武汉大学学报·信息科学版,2002,27(2):172~174
- 3 王仁宏.数值逼近理论.北京:高等教育出版社,2000. 220~229

4 Green P J, Silverman B W. Nonparametric Regression and Generalized Linear Models. London: Chapman & Hall, 1994

5 Fischer B, Hegland M. Collocation, Filtering and Nonparametric Regression. ZfV, 1999(1): 17~24

6 崔希璋, 於宗侑, 陶本藻, 等. 广义测量平差(新版). 武

汉: 武汉大学出版社, 2001

7 武汉测绘科技大学测量平差教研室. 测量平差基础(第三版). 北京: 测绘出版社, 1996. 83~95

第一作者简介: 吴云, 讲师, 硕士. 现从事现代测量数据处理理论及应用的研究.

E-mail: ywu@sgg.wtusm.edu.cn

Semiparametric Regression with Cubic Spline

WU Yun¹ SUN Haiyan¹ MA Xuezhong¹

(1 School of Geodesy and Geomatics, Wuhan University, 129 Luoyu Road, Wuhan 430079, China)

Abstract: Systematic errors contained in observations are always complicated smooth function varying with some variables. This paper describes this systematic errors using natural cubic spline, which is nonparametric component in semiparametric regression model. Penalised least squares technique implemented in the procedure reduces to unique solution. According to simulating tests, the semiparametric regression model and the penalised least squares technique can better separate systematic errors from observations compared with the parametric model and the least squares technique.

Key words: semiparametric regression; systematic error; natural cubic spline; smoothing parameter; penalised least squares technique

About the first author: WU Yun lecturer, master. Her major research is on the theory and application of surveying data processing.
E-mail: ywu@sgg.wtusm.edu.cn

(责任编辑: 洪远)

俄罗斯《文摘杂志》2003 年收录 《武汉大学学报·信息科学版》情况

截至 2003 年 12 月, 俄罗斯《文摘杂志》在 2003 年度中累计收录《武汉大学学报·信息科学版》发表的论文共 49 篇, 是上一年的 2.13 倍; 占同期发表论文数的 30.2%, 比上年提高 11 个百分点。具体如下:

- 1 黄声享 00.25(6): 485~490, 孙海燕 00.25(6): 496~499
- 1 李德仁 01.26(1): 1~7, 朱庆 01.26(1): 8~11, 马洪超 01.26(1): 18~23, 童小华 01.26(1): 64~69, 刘耀林 01.26(1): 75~81, 王新生 01.26(1): 82~85
- 1 王仁享 01.26(2): 95~100, 陶本藻 01.26(2): 101~104, 童小华 01.26(2): 105~111, 李桂苓 01.26(2): 118~121, 吴凡 01.26(2): 170~176
- 1 刘经南 01.26(3): 189~195, 王密 01.26(3): 205~208, 黄加纳 01.26(3): 213~216, 孙海燕 01.26(3): 222~225, 边馥苓 01.26(3): 232~238, 黄培之 01.26(3): 247~252
- 1 陈俊勇 01.26(4): 283~289, 卜方玲 01.26(4): 315~319, 闫军 01.26(4): 320~324, 刘慧敏 01.26(4): 325~330, 李英冰 01.26(4): 343~348, 张正禄 01.26(4): 354~360, 毛建华 01.26(4): 364~368
- 1 熊汉江 01.26(5): 393~398, 王峰 01.26(5): 425~429, 贾永红 01.26(5): 430~434, 张世强 01.26(5): 435~440, 张小红 01.26(5): 451~454, 钟业勋 01.26(5): 465~468
- 1 王任享 01.26(6): 487~490, 陶本藻 01.26(6): 504~508, 杨元喜 01.26(6): 509~513, 李建成 01.26(6): 514~517, 管铮 01.26(6): 529~532, 刘雁春 01.26(6): 533~538, 石磐 01.26(6): 549~554, 许才军 01.26(6): 555~561
- 1 陈南 02.27(2): 143~147, 安如 02.27(2): 188~193, 柴登峰 02.27(2): 199~202, 万晓霞 02.27(2): 203~207
- 1 李德仁 02.27(3): 221~233
- 1 熊兴华 02.27(5): 516~521, 王树根 02.27(5): 543~550
- 1 张剑清 02.27(6): 555~559, 张永军 02.27(6): 566~571