

可剔除多个粗差的F—T法

徐培亮

摘 要

本文首先导出了严密的F和T统计量，F统计量可用于粗差的整体检验，而T统计量可用于F检验后的各单个粗差的检验。然后给出一个算例，且与[1, 7]中的方法进行比较。可以看到，文中提出的F-T法在理论上是严密的，在实用上是可行的，因此，它优于现有的几种方法。

【关键词】 粗差剔除，严密统计量

一、引 言

在根据观测值求解待估参数以前，应该剔除粗差。当观测值中不存在粗差，或观测值中的粗差已经被剔除时，最小二乘平差的结果具有优良的统计特性。但当观测值中存在粗差时，最小二乘法所具有的优良性质是虚假的。因此，最小二乘估计的最大缺点是对粗差的抵御力差，受粗差的影响较大。近几年来，学者们提出了许多剔除粗差的方法，主要有：使用单位权方差先验值的Baarda的w检验^[2]，Stefanovic的 χ^2 检验^[7]，使用单位权方差估值的 τ 检验^[9]，t和F检验^[11]。此外，还有各种各样的近似统计量检验法^[3, 6]，稳健估计法^[4, 5]，向前向后选择法^[8]等。当观测值中存在多个粗差时，有些方法不严密，有些只能自动剔除大粗差，有些则只能进行整体检验，不便作各单个粗差的检验。本文导出了严密的F和T统计量，它们可分别用于整体检验和单个粗差的检验。特别是T统计量具有很多优点。

二、统计量的推导

设有数学模型为：

$$l_1 = AX + \varepsilon_1 \quad (1)$$

$$l_2 = A_2X + \Delta L_2 + \varepsilon_2 \quad (2)$$

本文1985年9月收到。

其中, A_1 是列满秩矩阵, 秩为 t , 就是说, 由 l_1 可以构成基本网形, 而且是没有粗差的观测值向量; l_2 是经过初步判别后认为存在粗差的观测值向量; ε_1 和 ε_2 分别为 l_1 和 l_2 的随机误差向量; X 是 t 维未知参数向量。

式 (1) 的最小二乘解为:

$$\hat{X} = (A_1^T P_1 A_1)^{-1} A_1^T P_1 l_1 \quad (3)$$

$$D(\hat{X}) = (A_1^T P_1 A_1)^{-1} \sigma^2 \quad (4)$$

$$V_1 = A_1 \hat{X} - l_1 \quad (5)$$

$$\hat{\sigma}^2 = V_1^T P_1 V_1 / (n - m - t) \quad (6)$$

其中, P_1 是对应于观测向量 l_1 的权矩阵, n 是观测值总数, m 是 l_2 的维数。

利用 (3) 式可求得 l_2 的预测向量,

$$\hat{l}_2 = A_2 (A_1^T P_1 A_1)^{-1} A_1^T P_1 l_1 \quad (7)$$

对于 l_2 中的每个元素 l_i , $i > n - m$, 预测值为:

$$\hat{l}_i = A_i (A_1^T P_1 A_1)^{-1} A_1^T P_1 l_1 \quad (8)$$

其中, A_i 是 A_2 的一行向量。

由实测的 l_2 和 (7) 式可求得预测残差为:

$$\hat{V}_2 = l_2 - \hat{l}_2 = l_2 - A_2 \hat{X} \quad (9)$$

对于 \hat{V}_2 中的每个元素 \hat{v}_i 为:

$$\hat{v}_i = l_i - A_i \hat{X} \quad (10)$$

当 l_1 与 l_2 相互独立, 则有:

$$D(\hat{V}_2) = D(l_2) + D(\hat{l}_2) = \{P_2^{-1} + A_2 (A_1^T P_1 A_1)^{-1} A_2^T\} \sigma^2 \quad (11)$$

$$D(\hat{v}_i) = \{p_i^{-1} + A_i (A_1^T P_1 A_1)^{-1} A_i^T\} \sigma^2 \quad (12)$$

检验 l_2 中是否存在粗差, 实际上相当于检验原假设

$$H_0: \Delta L_2 = 0 \quad (13)$$

当原假设成立时, \hat{V}_2 是数学期望为零向量, 方差为 (11) 式的正态分布的随机向量; 当原假设不成立时, \hat{V}_2 是期望值为 ΔL_2 、方差为 (11) 式的正态分布随机向量。易知, 二次型

$$U = \hat{V}_2^T D^{-1}(\hat{V}_2) \hat{V}_2 \quad (14)$$

在原假设成立时是一个 χ^2 分布的随机变量, 其非中心参数为零, 自由度为 m ; 否则其非中心参数为 $\delta = \Delta L_2^T D^{-1}(\hat{V}_2) \Delta L_2$ 。事实上, 它就是 Stefanovic 的 χ^2 检验统计量。

因为二次型 U 中含有未知的 σ^2 参数, 如果其先验值未知, 则不便于检验。但是, 由

于 V_1 与 \hat{X} 相互独立, $\hat{\sigma}^2$ 与 \hat{X} 相互独立, 因此, $\hat{\sigma}^2$ 与 \hat{V}_2 相互独立, 从而 U 与 $\hat{\sigma}^2$ 也相互独立。二次型 U 反映的是粗差的影响, $\hat{\sigma}^2$ 反映的是随机误差的影响, 原假设 (13) 实际上也相当于检验 $E(U)$ 与 $E(\hat{\sigma}^2) = \sigma^2$ 是否存在显著性差异。利用 U 与 $\hat{\sigma}^2$ 构成检验统计量, 记为 F

$$F = \frac{U/m}{\hat{\sigma}^2/\sigma^2} = \frac{\hat{V}_2^T \{P_2^{-1} + A_2(A_1^T P_1 A_1)^{-1} A_2^T\}^{-1} \hat{V}_2}{m \hat{\sigma}^2} \quad (15)$$

其中, $\hat{\sigma}^2(n-m-t)/\sigma^2 \sim \chi^2(n-m-t)$ 。因此, 当原假设成立时, F 统计量严格地服从中心 F 分布, 即 $F \sim F(m, n-m-t)$; 当原假设不成立时, F 统计量服从于非中心参数为 $\delta = \Delta L_2^T D^{-1}(\hat{V}_2) \Delta L_2$ 的 F 分布。

利用 (15) 式可以方便地对 l_2 整体上是否存在粗差作统计检验。取显著性水平 α_F , 当有 $F < F_{\alpha_F}(m, n-m-t)$ 时, 接受原假设, 即在检验置信水平为 $(1-\alpha_F)$ 下认为 l_2 整体上不存在粗差; 反之, 否定原假设, 这时认为 l_2 整体上存在粗差。为了进一步判别是哪些 l_i 要对否定原假设负责, 即哪些 l_i 含有粗差, 还需要进一步依次对每个 l_i 进行检验。

检验单个 l_i 是否存在粗差, 实际上相当于检验原假设

$$H_0: \quad \Delta L_i = 0 \quad (16)$$

由式 (10) 和 (12) 知, 当原假设成立时, \hat{v}_i 服从期望值为零, 方差为 (12) 式的正态分布; 当原假设不成立时, 其期望值为 ΔL_i 。因为 \hat{V}_2 与 $\hat{\sigma}^2$ 相互独立, 所以 \hat{v}_i 与 $\hat{\sigma}^2$ 相互独立。利用 \hat{v}_i 和 $\hat{\sigma}^2$ 构成检验统计量, 记为 T

$$T = \frac{\hat{v}_i \sqrt{D(\hat{v}_i)}}{\sqrt{\hat{\sigma}^2/\sigma^2}} = \frac{\hat{v}_i}{\{P_i^{-1} + A_i(A_1^T P_1 A_1)^{-1} A_i^T\}^{\frac{1}{2}} \hat{\sigma}} \quad (17)$$

显然, 在原假设成立条件下, T 统计量服从中心 t 分布, 且其自由度为 $(n-m-t)$; 否则, 其非中心参数为 $\Delta L_i / \{P_i^{-1} + A_i(A_1^T P_1 A_1)^{-1} A_i^T\}^{\frac{1}{2}} \sigma$ 。因此, 作为检验统计量 T , 可以用于单个 l_i 是否存在粗差的检验。

我们将利用上面导出的 F 和 T 统计量进行检验来剔除粗差的方法称为 $F-T$ 法。

三、计算步骤及算例分析

按 $F-T$ 法剔除粗差的计算步骤是:

1. 初步判别粗差。初步判别粗差的方法可以用最小二乘法和稳健估计法等。最小二乘判别法的依据是样本中位数, 而稳健估计判别法的依据是残差的大小。稳健估计法具有数值计算的优点, 但不能为各种检验建立严密的统计量, 因此, 难以对粗差作出正确的判断。此外, 如果某一观测值含有粗差, 则其改正数与该值具有的权没有很大关系, 且估计与初值有关, 据文献 [4], 当初值不理想时, 稳健估计法可能收敛于局部最优, 而不是整体最优。所以, 利用稳健估计法还不能对粗差做最后的诊断。但它是一种有效的初步判别法。当然还可

以使用文献[8]中的向后选择法等。

2. 根据初步判别结果划分平差模型(1)和预测模型(2), 求解平差值 \hat{X} 和预测残差 \hat{V}_2 及其方差——协方差矩阵。

3. 利用式(15)做 l_2 的整体性检验。当否定原假设时, 还须进一步应用 T 检验判别哪些 l_i 存在粗差。

4. 根据检验结果求平差值, 作为最后结果。

下面用一个例子来说明上述方法的应用

本例取自[5]。共21组数据, 待估参数4个, 观测权阵 $P=I$ 。

稳健估计的残差结果取自[5], 21个残差:

6.11 1.04 6.31 8.24 -1.24 -0.71 -0.33
 0.67 -0.97 0.14 0.79 0.24 -2.71 -1.44
 1.33 0.11 -0.42 0.08 0.63 1.87 -8.91

由此可见, 观测值1、3、4、21和13对应的残差较大, 按要求, 归为预测模型, 其它观测值构成平差模型(1)。

(1)式的最小二乘平差残差结果及预测残差共21组数据如下*:

-**5.91** -0.82 -**6.13** -**8.30** 0.81 1.26 0.15
 -0.84 0.87 0.30 -0.54 0.11 **3.11** 1.54
 -1.31 -0.03 0.71 0.06 -0.58 -1.69 **9.32**

其中黑体字表示由(7)式算得的预测残差。

单位权方差估计 $\hat{\sigma}^2 = 1.0504$

l_2 的整体性检验, 由(15)式算得:

$$F = 31.65$$

给定 $\alpha_F = 0.05$, 查表得 $F_{0.05}(5, 12) = 3.11$ 。因 $F > F_{0.05}(5, 12)$, 所以否定原假设, 即认为 l_2 整体上存在粗差。须进一步地做单个观测值是否存在粗差的检验。取 $\alpha_T = 0.01$, 把检验结果列于表1中。

表 1

点号	T 计算值	$T_{0.01}(12)$	$T_{0.05}(12)$	检验结果
1	4.4436			$T_{0.01}(12) < T $
3	5.0138			
4	7.4574			
21	7.2391	3.0545	2.1788	
13	2.7243			$T_{0.01}(12) > T $

* 最小二乘平差残差和稳健估计残差的符号相反。

从表 1 的检验结果知, 在 $\alpha_T = 0.01$ 时, 点 1、3、4 和 21 存在粗差, 而点 13 不存在粗差。但是, 当 $\alpha_T = 0.05$ 时, 点 13 也能被判为粗差。对于象测值 13 这样的临界粗差, 用稳健估计难于作出正确的判断。并且, 稳健估计法做出的判断没有概率意义。由于稳健估计法与初值有关, 可能收敛于局部最优, 所以稳健估计法可能误判粗差, 特别是那些临界值粗差。可是, 它与 F—T 法结合能够准确地在概率意义上进行粗差的判别。

与 [1] 中的 F 检验法或 [7] 中的 χ^2 检验法相比, 在完成了整体性检验后, 当对各单个粗差进行检验时, 则这两种方法都须重新平差才能检验, 而 F—T 法的 T 统计量则便于直接在整体检验后使用, 具有使单个粗差的检测变得容易的优点。

与 Heck 的 t 检验和 Pope 的 τ 检验相比, 当观测值中存在多个粗差时, 它们不能进行整体检验, 因而也不可能讨论多个粗差的检测与各单一粗差的检测之间的显著性水平间的关系。并且它们都不是严格统计量, 单位权方差的估计是有偏的。

现把 $t_i = V_i / \sqrt{q_{ii}} \hat{\sigma}_i$ 其中 $(n-t-1) \hat{\sigma}_i^2 = V^T P V - V_i^2 / q_{ii}$, t 为 X 的维数, V_i 是全部观测值的最小二乘平差的残差 i 分量) 的计算结果列于下:

-1.2095	0.7051	-1.6179	-2.0518	0.5305	0.9632	0.8259
0.4737	1.0486	-0.4262	-0.8783	-0.9667	0.4687	0.0169
-0.8006	-0.2912	0.5996	0.1487	0.1972	-0.4431	3.3305

为了进行比较, 取 $\alpha = 0.01$, 查表得 $T_{0.01}(16) = 2.9208$, 因此, 首次检验结果只检测出观测值 21 存在粗差。其根本原因在于, 当观测值中存在多个粗差时, 残差向量 V 受到所有粗差的影响, 从而导致单位权方差估计 $\hat{\sigma}^2$ 有偏, 有时甚至与其真值相差甚远。反之, 本文的 F—T 法能提供基本上无偏的估值。

四、结 论

本文提出的 F—T 法的 F 统计量和 T 统计量分别严格地服从 F 分布和 t 分布, F 统计量可用于整体检验, T 统计量可用于单个粗差的检验。

当单位权方差的先验值已知时, 本文的 U 统计量就是 Stefanovic 的 χ^2 检验统计量^[7]。由于本文另导出了严密的 T 统计量, 因此, 与 [1] 的 F 检验或 [7] 中 χ^2 检验相比, 易于进行单个粗差的检验。当观测值中存在多个粗差时, 与 t 检验和 τ 检验相比, 本文的方法可进行整体检验。且使用的单位权方差基本是无偏的, 用它可构成 χ^2 随机变量, 而 t 检验或 τ 检验使用的单位权方差是有偏的, 有时甚至与其真值相差很大。因此, 严格地说, 当观测值中存在多个粗差时, t 或 τ 检验统计量不严密。可见, F—T 法优于以上几种方法。

参 考 文 献

- [1] 陈永奇, 测量工程中的个别问题, 武汉测绘学院, 1984.
- [2] W. Baarda, A compass for the land surveyor, 1968.
- [3] 李德仁, 利用验后方差估计发现观测值中粗差的一种方法, 测绘译丛, No.1, 1984.

- [4] 陈希孺、王松桂, 近代实用回归分析, 广西人民出版社, 1984.
- [5] D.F.Andrews, A robust method for multiple linear regression, *Technometrics* Vol.16 No. 4, 1974.
- [6] K.Kubik, An error theory for the Danish method, Symposium of ISP, Comm. 3, Helsinki, 1982.
- [7] P.Stefanovic, Blunders and least squares, *ITC Journal*, 1, 1978.
- [8] B.Benciolini, Luigi, Mussio, F.Sansó, An approach to gross detection more conservative than Baarda snooping, Symposium of ISP, Comm3, Helsinki, 1982.
- [9] J.Pope, The statistics of residuals and the detection of outliers, NOAA technical report Nos 66. NGS, 1976.

A F—T Method for Outliers

Xu Peiliang

Abstract

This paper derives rigorous F and T statistics which are used for the global test of outliers and for the test of each outlier immediately and directly after F-test respectively. Then an example is given for the purpose of comparison with the methods in [1,7]. It can be seen that the method presented in this paper is rigorous in theory, feasible in practice and better than the methods [1, 7].

[Key words] detection of outliers, rigorous statistics