

基于相关分析的粗差理论*

施 闯 刘经南

(武汉测绘科技大学地学测量工程学院,武汉市珞喻路 129 号, 430079)

摘 要 提出了基于相关分析的可靠性研究和粗差分析理论,并应用于相关观测量的多粗差分析。

关键词 相关分析法;相关观测量;粗差分析;误差处理

分类号 P207

现代测量手段趋于向数据采集的自动化和快速发展。然而由于各种因素的影响,由仪器和计算机自动采集的数据中,常带有一定数量的粗差观测量。这些粗差观测量如果直接参与数据处理,将大大降低成果应有的精度,甚至导致错误的估计。现有的粗差分析方法在实用上多限于独立观测量或等权独立观测量的单个或多个粗差的条件,而对于相关观测量的多粗差问题,被认为是一个难题。对此,本文提出了用相关分析的方法,解决相关观测量的多粗差问题。

1 原 理

在相关观测量的最小二乘平差中,一个观测量的误差,通过观测量之间随机特性的相关性和图形几何条件的关联性,反映于其它观测量的平差改正数之中。当一个观测量含有粗差时,它必将或多或少地影响到其它观测量的改正数。那么粗差观测量是否与改正数向量 V 之间存在着某种必然的联系?

1.1 观测量误差对改正数向量的影响

由可靠性矩阵^[1]:

$$R = Q_{VV}P_{II} = E - A(A^T P_{II} A)^{-1} A^T P_{II} \quad (1)$$

且有: $RX = -V$ (2)

其中, A 是平差系统的图形设计矩阵; P_{II} 是观测量的权矩阵; X 是观测值误差; V 是观测值改正数向量; Q_{VV} 是观测值改正数的协因数阵。将 (2) 式展开整理得:

$$RX = \begin{bmatrix} r_{11} \\ r_{21} \\ \vdots \\ r_{n1} \end{bmatrix} X_1 + \begin{bmatrix} r_{12} \\ r_{22} \\ \vdots \\ r_{n2} \end{bmatrix} X_2 + \cdots + \begin{bmatrix} r_{1n} \\ r_{2n} \\ \vdots \\ r_{nn} \end{bmatrix} X_n$$

$$= -[v_1, v_2, \cdots, v_n]^T \quad (3)$$

式中, X_i 为数值变量 ($i = 1, 2, \cdots, n$)。令

$$F_i = [r_{1i}, r_{2i}, \cdots, r_{ni}]^T, (i = 1, 2, \cdots, n)$$

则有:

$$X_{F1} + X_{F2} + \cdots + X_{Fn} = -V \quad (4)$$

称 F_i 为观测量 i 的误差 X 对改正数向量 V 的影响向量。

F_i 由平差系统的图形设计矩阵 A 和观测量的权矩阵 P_{II} 所决定,它反映了观测量 i 的误差对改正数向量 V 的内在的影响关系和作用程度。

由 (4) 式可以看出: 观测量的改正数向量 V 可以表示为 X_{Fi} ($i = 1, 2, \cdots, n$) 的向量和; 各观测量的误差 X 在其影响向量 F_i 的作用下反映于改正数向量 V 之中; X 起到了对向量 F_i 的缩放作用。

1.2 粗差与改正数向量的相关关系

当某观测量出现粗差时,无论是将其归于函数模型还是随机模型的误差^[1],都将表现为该观测量的 $\|X_{Fi}\|_2$ 值,即向量 X_{Fi} 的模较大,且在 (4) 式左半部分各项中占有优势。这时,改正数向量 V 将突出地表现为与该观测量的影响向量 F_i 有较强的相关性。

图 1 显示出观测量中含有一个粗差时,粗差观测量 ($i = 18$) 的影响向量 F_{18} 的各分量 $r_{11}, r_{21}, \cdots, r_{n1}$ ($i = 18, n = 25$) 与 V 对应分量 v_1, v_2, \cdots, v_n ($n = 25$) 的关系。

当观测量中含有多个粗差时,就表现为这种相关关系的叠加。图 2 中的数据是实测数据 (相关观测量) 中,加入了 3 倍观测量中误差的粗差所得到的结果。

下面对向量 F_i 与 V 的相关性进行定量的分析。

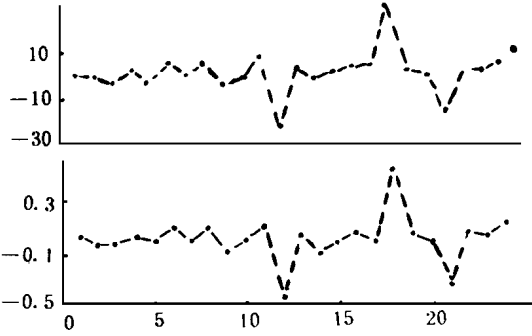


图 1 含单个粗差时,粗差的影响向量
 F_{18} 与 V 的关系

Fig. 1 Relation Between the Outlier
Influence Vector F_{18} and the Residual
Vector V (Single Outlier)

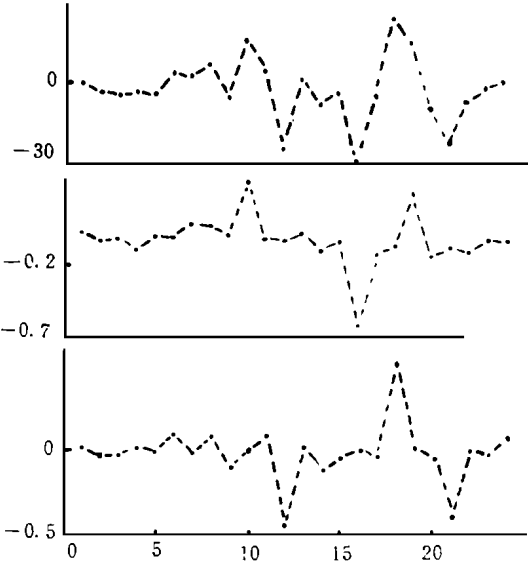


图 2 含两粗差时,粗差的影响向量
 F_{18} F_{16} 与 V 的关系

Fig. 2 Relation Between the Outlier Influence
Vector F_{18} , F_{16} and the Residual
Vector V (Two Outliers)

如图 1 图 2 所示,若要定量地描述两者的相关程度,可以考虑用它们对应分量 F_{ji} ($F_{ji} = r_{ji}$) 与 v_j ($j = 1, 2, \dots, n$) 差的平方和的最小值 Q_0 来表示:

$$Q_0 = \min_{a,b} \frac{1}{n} \sum_{j=1}^n (v_j - a - bF_{ji})^2 \quad (5)$$

如果存在某个 a 和 b , 使得 $Q_0 = 0$, 则可以说 F_i 与 V 完全相关, 否则就用 Q_0 的大小来描述两者的相关程度, 为了求 Q_0 , 对函数

$$Q(a,b) = \frac{1}{n} \sum_{j=1}^n (v_j - a - bF_{ji})^2 \quad (6)$$

求关于 a b 的偏导数, 并令其等于 0, 即

$$\frac{\partial Q}{\partial a} = - \frac{2}{n} \sum_{j=1}^n (v_j - a - bF_{ji}) = 0$$

$$\frac{\partial Q}{\partial b} = - \frac{2}{n} \sum_{j=1}^n ((v_j - a)F_{ji} - bF_{ji}^2) = 0$$

解得:

$$b = \frac{\sum_{j=1}^n (F_{ji} - \bar{F}_i)(v_j - \bar{v})}{\sum_{j=1}^n (F_{ji} - \bar{F}_i)^2} \quad (7)$$

$$a = \bar{v} - b\bar{F}_i$$

将式 (7) 代入式 (6) 得:

$$Q_0 = \frac{1}{n} \sum_{j=1}^n (v_j - \bar{v})^2 (1 - d_{F_i,V}^2) \quad (8)$$

其中记:

$$d_{F_i,V} = \frac{\sum_{j=1}^n (F_{ji} - \bar{F}_i)(v_j - \bar{v})}{\left(\sum_{j=1}^n (F_{ji} - \bar{F}_i)^2 \sum_{j=1}^n (v_j - \bar{v})^2 \right)^{1/2}} \quad (9)$$

称 $d_{F_i,V}$ 为 F_i 与 V 的相关系数。 $d_{F_i,V}$ 是一个无量纲的量, 它定量地反映了 F_i 与 V 之间的相关程度

相关系数 $d_{F_i,V}$ 具有下列性质:

$|d_{F_i,V}| \leq 1$; $|d_{F_i,V}|$ 越大, 则 F_i 与 V 越相关;
 $|d_{F_i,V}| = 0$ 时, F_i 与 V 不相关

一个粗差的出现, 会对整个平差系统产生较大的影响, 这种影响可以通过 F_i 的作用, 反映于改正数向量 V 之中。 要想从 V 中发现粗差的踪迹, 必须从分析 F_i 与 V 的关系入手, 而 $d_{F_i,V}$ 正是定量地反映了两者之间的相关程度, 这也就是 $X F_i$ 与 V 的相关程度。 当 $|d_{F_i,V}|$ 越接近于 1 时, 相关性越强, 说明改正数向量 V 的变化受来自观测量 i 的误差的影响越显著。 如果存在粗差, 则该观测量为粗差的可能性也就越大。 反之, $|d_{F_i,V}|$ 越接近于 0 时, $X F_i$ 对 V 的影响越不显著, 则它为粗差的可能性就越小。

当观测量中含有多个粗差时, 改正数向量 V 的变化规律, 将表现为来自这些粗差对其影响的叠加 (如图 2 所示)。 各粗差观测量的 F_i 都将显示出与 V 有着相对显著的相关性, 相关性的 大小取决于粗差个数、 分布和粗差值的大小。 所以当相关观测量中含有多个粗差时, 只要粗差是可测的, 通过 $d_{F_i,V}$ 仍能分析出粗差所在的位置。

2 相关系数 $d_{F_i,V}$ 的检验

如果 F_i V 的分量由随机变量 f_i v 组成, 则其相关系数可由下式给出:

$$\tilde{d}_{F_i,V} = \frac{\text{cov}(f_i, v)}{\sqrt{D(f_i) D(v)}} \quad (10)$$

其中, $\text{cov}(f_i, v) = E\{[f_i - E(f_i)][v - E(v)]\}$

$$D(f_i) = E(f_i^2) - [E(f_i)]^2$$

$$D(v) = E(v^2) - [E(v)]^2$$

在实际问题中,如果有 n 个观测量,计算 F_k V 的相关系数是用式 (9) 求 $d_{F_i, V}$ 来估计 $\tilde{d}_{F_i, V}$

如果 f_k V 服从二维正态分布,则当 $\tilde{d}_{F_i, V} = 0$ 时, $d_{F_i, V}$ 的概率密度函数为:

$$\tilde{d}_{F_i, V}(d_{F_i, V}) = \frac{1}{\pi} \frac{\Gamma\left(\frac{n-1}{2}\right)}{\Gamma\left(\frac{n-2}{2}\right)} (1 - d_{F_i, V}^2)^{\frac{n-4}{2}} \tag{11}$$

其中, $\Gamma\left(\frac{n-1}{2}\right)$ 、 $\Gamma\left(\frac{n-2}{2}\right)$ 是 Gamma 函数

检验时,取零假设 $H_0: \tilde{d}_{F_i, V} = 0$,即 F_k V 不相关,则有统计量:

$$t = d_{F_i, V}(n - 2)^{1/2} / (1 - d_{F_i, V}^2)^{1/2} \tag{12}$$

t 是遵从 $n - 2$ 自由度的 t 分布

给定置信度 $\alpha = 0.001$,则由 (12) 式可得:

若
$$d_{F_i, V} > \frac{t_\alpha}{t_\alpha^2 + (n - 2)} \tag{13}$$

成立,则拒绝零假设,认为 F_i 与 V 相关是显著的;否则接受零假设

3 相关分析法的几何解释

由式 (9) 定义的相关系数可以表达为下列形式:

$$d_{F_i, V} = \frac{F_i^T H V}{\|H F_i\| \|H V\|} = \cos\theta \tag{14}$$

式中, H 为中心化矩阵, θ 为向量 $H F_i$ 和 $H V$ 的夹角。这就是说, F_i 与 V 的相关系数的几何意义是其经过中心化变换后的向量 $H F_i$ 与 $H V$ 的夹角的余弦。 F_i 与 V 的相关性越强, $H F_i$ 与 $H V$ 的夹角越接近于 0° 或 180° 。

根据式 (4),若 i 出现粗差,即 $\|X F_i\|_2$ 显著大于其它值,这时 $X F_j$ ($j = 1, 2, \dots, n$) 的向量之和必然使 $H V$ 的方向接近于 $H F_i$,即 F_i 与 V 具有较强的相关性。

同样,观测量中存在多个粗差时,由于观测量的个数远远大于粗差的个数, F 向量空间的维数也远远大于粗差个数。这时, $X F_i$ ($i = 1, 2, \dots, n$) 的向量和,同样会使 $H V$ 在 n 维空间中的方向接近于各粗差测量的 $H F_i$ 的方向。

当然有可能会出 现多个粗差共同作用的结果,使 $H V$ 的方向更接近于某个非粗差观测量的 $H F_i$,但这种情况仅限于某些特殊的图形结构和

粗差组合。在粗差分析中,可以通过反向搜索的方法,去伪存真。

4 平差系统中是否含有粗差

以 GPS 网为例,对验后单位权中误差 $\hat{\sigma}_0^2$ 进行 i^2 检验。 $\hat{\sigma}_0^2$ 的理论值为 1,则有 i^2 检验量:

$$i^2 = V^T P V \sim i^2(f)$$

取零假设: $H_0 = \hat{\sigma}_0^2 = 1$; 选置信水平 $1 - \alpha$, $\alpha = 0.005$,如果:

$$i_{p_1}^2 \leq V^T P V \leq i_{p_2}^2 \text{ 或 } \frac{V^T P V}{i_{p_1}^2} \leq \hat{\sigma}_0^2 \leq \frac{V^T P V}{i_{p_2}^2}$$

其中, $p_1 = (1 - \alpha)/2$, $p_2 = (1 + \alpha)/2$ 成立,则接受零假设。否则,拒绝零假设,认为平差系统中可能存在粗差

5 粗差的可测性和可区分性

在一些情况下,某些观测量的粗差是不可能被平差系统自身所发现的;而有些粗差即使被发现,也不可能区分它存在于哪个观测量之中。这就是粗差的可测性和可区分性问题。

基于本文中的理论,不难得出下列结论:

1) 观测量对改正数向量的影响向量 F_i 的模 $\|F_i\|_2$ 越接近于 0,则平差系统对该观测量的误差越不敏感,该观测量的粗差越不易被发现。当 $F_i = 0$ 时,该观测量的粗差是不可测的,称为粗差不可测观测量。

2) 对于两个观测量,其影响向量 F_k F_j 的相关系数为:

$$d_{F_i, F_j} = \frac{\sum_{k=1}^n (F_{ki} - \bar{F}_i)(F_{kj} - \bar{F}_j)}{\left(\sum_{k=1}^n (F_{ki} - \bar{F}_i)^2\right)^{1/2} \left(\sum_{k=1}^n (F_{kj} - \bar{F}_j)^2\right)^{1/2}} \tag{15}$$

当 $|d_{F_i, F_j}|$ 越接近于 1 时,这两个观测量的粗差的可区分性越差;若 $|d_{F_i, F_j}| = 1$,则这两个观测量的粗差是不可区分的,称为粗差不可区分观测量。

6 粗差分析实现的方法和步骤 (以 GPS 网为例)

6.1 同时探测多个粗差

1) 对观测量进行最小二乘平差,求得观测量改正数向量 V 和验后单位权中误差 $\hat{\sigma}_0$

- 2)对 $\hat{\epsilon}_0^2$ 进行 i^2 检验,若不通过,则系统中可能含有粗差
- 3)计算可靠性矩阵 R 及各观测量对改正数的影响向量 F_i
- 4)计算各观测量的 F_i 与 V 的相关系数 $d_{F_i,V}$
- 5)对各 $d_{F_i,V}$ 的显著性进行 $t_{(n-2)}$ 检验,若显著,则对该观测量做粗差标记

6)对所有粗差标记的观测量进行剔除或降权处理后,重复 1)~ 6)。若 2)中 $\hat{\epsilon}_0^2$ 检验通过,则进行 7)。

7)将有粗差标记的观测量,逐一恢复到平差系统中,重复 1)、2)。若 $\hat{\epsilon}_0^2$ 检验仍能通过,则该观测量为误判,应确认恢复;否则确认为粗差。在这一步中,应注意不可区分粗差的特殊情况。

8)最后确认所有被标记的观测量中含有粗差,进行适当的处理后,完成粗差分析,重新进行平差。

6.2 逐个探测多粗差

- 1)、2)、3)、4)步同 § 6.1;
- 5)对 $|d_{F_i,V}|$ 最大的一个观测量作粗差标记;
- 6)对该观测量进行剔除或降权处理后,重复 1)~ 6),若其中 2) $\hat{\epsilon}_0^2$ 检验通过,则进行 7);
- 7)、8)同 § 6.1

比较这两种方法,前者适用于观测量间相关性较弱的情况下的多粗差探测,而后者适用于观测量之间相关性较强时的粗差探测。

7 实践和算例

高精度 GPS 网的基线处理,大都采用网解方式,一个同步网的基线解的协方差阵是满阵,即在一个同步网中所有基线向量之间及基线向量各分量之间都是相关的,而整个 GPS 网可由若干个同步网组成。

7.1 模拟粗差

图 3 是一实测高精度 GPS 网,由 4 个同步网组成,每个同步网含两条基线向量。在同步网 B9502601.LCX 中基线 TN08-TN09 的 ΔX 和 ΔY 分量上,加入 3 倍观测量平差值中误差的粗差,则粗差分析中得出表 1。

若简单地从观测量改正数 V 及其精度入手,则难以发现粗差所在位置。根据 (13) 式,对 $d_{F_i,V}$ 进行 t_{n-2} 检验,这时 $n=22$,取 $\mathbb{T}=0.001$ 时,检验量为 0.628。不难发现: $d_{4,V}=0.849>0.628$, $d_{5,V}=0.691>0.628$,正是粗差所在的观测量。

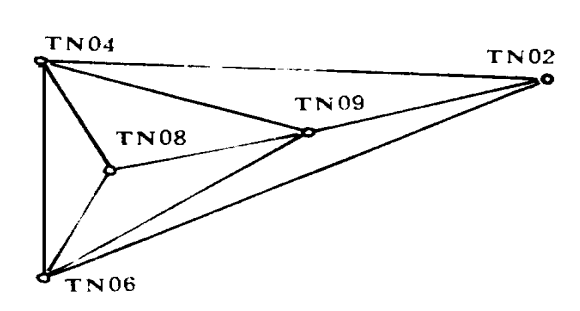


图 3 GPS网图

Fig. 3 GPS Network Plot

表 1 粗差分析结果

Tab. 1 Results of Outlier Analysis

序号	测站	观测量	V/m	精度 $/m$	d	同步网
1	TN08	DX	$-0.0044 \pm .0107$.378	B9502601.LCX	
		DY	$0.0031 \pm .0154$.273		
		DZ	$0.0016 \pm .0091$.208		
2	TN08#	DX	$0.0173 \pm .0121$.849	B9502601.LCX	
		DY	$0.0190 \pm .0300$.691		
		DZ	$-0.0004 \pm .0098$.190		
3	TN02	DX	$0.0121 \pm .0106$.378	B9502701.LCX	
		DY	$-0.0017 \pm .0218$.273		
		DZ	$-0.0052 \pm .0120$.208		
4	TN02	DX	$-0.0026 \pm .0149$.104	B9502701.LCX	
		DY	$0.0018 \pm .0287$.080		
		DZ	$-0.0004 \pm .0152$.199		
5	TN09	DX	$0.0191 \pm .0138$.266	B9502602.LCX	
		DY	$0.0361 \pm .0343$.155		
		DZ	$0.0086 \pm .0176$.004		
6	TN09	DX	$0.0190 \pm .0152$.332	B9502702.LCX	
		DY	$0.0368 \pm .0366$.365		
		DZ	$0.0096 \pm .0185$.043		
7	TN08	DX	$-0.0165 \pm .0218$.254	B9502703.LCX	
		DY	$0.0038 \pm .0357$.363		
		DZ	$0.0101 \pm .0323$.230		
8	TN08	DX	$-0.0350 \pm .0237$.529	B9502703.LCX	
		DY	$-0.0288 \pm .0580$.145		
		DZ	$-0.0027 \pm .0469$.257		

7.2 实践算例

图 4 是国家高精度 GPS B 级网某测区中,任意提取的一个子块。

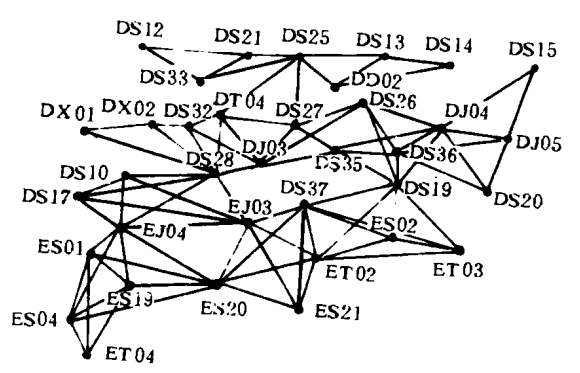
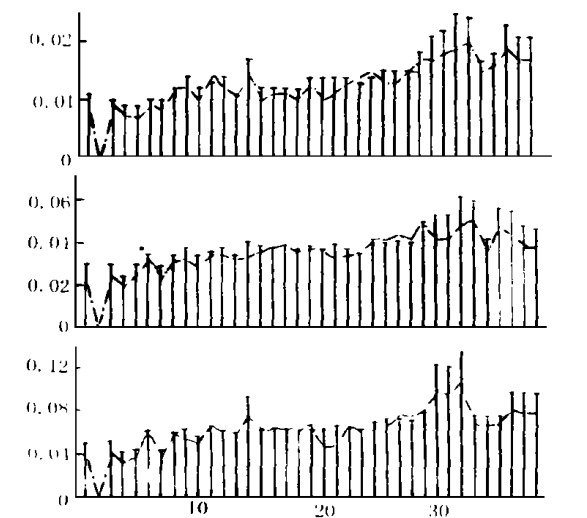


图 4 GPS网图

Fig. 4 GPS Network Plot

对原始基线观测量进行最小二乘无约束平差,其验后单位权中误差 $\hat{\sigma}_0^2=1.374$, χ^2 检验不通过,说明平差系统中含有小的粗差。按 § 6.2 节逐个探测的方法进行粗差分析,最后确认 3 个同步网中含有粗差基线。对其进行简单的剔除后,重新平差,则 $\hat{\sigma}_0^2=1.076$, χ^2 检验通过。

图 5 是粗差处理前后平差结果的坐标分量精度比较。



T和曲线分别表示粗差处理前后的坐标分量精度

图 5 粗差处理前后坐标分量 rms 比较

Fig. 5 Comparisons of Component rms

8 结 论

1) 观测量误差对改正数向量的影响向量 F_i

由平差系统的图形设计矩阵和观测量权矩阵共同决定,它反映了观测量的误差对改正数向量的内在的影响关系和作用程度。

2) 无论将粗差归于函数模型的误差还是随机模型的误差,都将在 $\|X F_i\|_2$ 上有显著的反映。

3) $\|F_i\|_2$ 的大小反映了该观测量粗差的可测性程度。

4) 观测量之间的相关系数 d_{F_i, F_j} 反映了两个观测量粗差的可区分性程度。

5) 若观测量中含有粗差,则其对应的 F_i 与 V 将表现为较强的相关性,这时 $d_{F_i, V}$ 的 $t_{(n-2)}$ 检验显著。

本文所涉及的理论、算法和实践,由武汉测绘科技大学研制的 GPSN/AS 软件实现。

参 考 文 献

1 李德仁.误差处理和可靠性理论.北京:测绘出版社,1988.

2 崔希璋,刘经南.国家高精度 GPS 网数据处理个别问题探讨.武汉测绘科技大学学报,1994,19(增刊).

3 胡定国,张润楚.多元数据分析方法.天津:南开大学出版社,1990.

4 方开泰,张尧庭.广义多元分析.北京:科学出版社,1993.

5 熊西文等译.数值分析的理论及其应用.上海:上海科技出版社,1980.

6 孙继广.矩阵扰动分析.北京:科学出版社,1987.

7 施 闯.国家高精度 GPS 网数据处理方案、模型和软件研究. [学位论文]. 武汉: 武汉测绘科技大学, 1995.

8 程云鹏.矩阵论.西安:西北工业大学出版社,1989.

Correspondence Based Outlier Analysis

Shi Chuang Liu Jingnan

(School of Geoscience and Surveying Engineering, WUTUSM, 129 Luoyu Road, Wuhan, China: 430079)

Abstract The specialty of reliability for correlative observables is studied. The residual vector can be expressed as a linearization equation of observable's errors and column components of the reliability matrix (the influence vectors). The correspondence coefficient of the influence vector and the residual vector is used to study on the condition that the observable's error is outlier. In this paper, the theory of correspondence analysis is proposed and discussed. This theory is used to solve the problem of mult-outlier analysis in correlative observables.

Key words correspondence analysis; correlative observable; outlier analysis; error discussing