

ICS 03.120.30
A 41



中华人民共和国国家标准

GB/T 4883—2008
代替 GB/T 4883—1985

数据的统计处理和解释 正态样本离群值的判断和处理

Statistical interpretation of data—
Detection and treatment of outliers in the normal sample

2008-07-16 发布

2009-01-01 实施

中华人民共和国国家质量监督检验检疫总局
中国国家标准化管理委员会 发布

507



前 言

本标准代替 GB/T 4883—1985。本标准与 GB/T 4883—1985 相比较,技术内容的变化主要包括:

- 增加了术语、定义和符号一章;
- 将“正态样本异常值的判断和处理”改为“正态样本离群值的判断和处理”;
- 将术语“检出异常值”和“高度异常值”分别改为“歧离值”和“统计离群值”,并进一步明确了二者的含义及相互差异;
- 增加了检出水平和剔除水平的定义;
- 检出水平由原标准中“检出水平 α 一般取为1%,5%或10%”改为“除非根据本标准达成协议的各方另有约定外, α 值应为0.05”;
- 明确规定剔除水平 α' 为“除非根据本标准达成协议的各方另有约定外, α' 值应为0.01”;
- 增加了各种情形“统计离群值”的检验步骤;
- 将“没有异常值”和“没有高度异常值的离群值”分别改为“未发现离群值”和“未发现统计离群值”;
- 增加了奈尔(Nair)统计量、格拉布斯(Grubbs)统计量、狄克逊(Dixon)统计量、偏度统计量、峰度统计量的符号;
- 作狄克逊(Dixon)检验时,将样本量由30扩充到100,此内容作为附录C。

本标准的附录A为规范性附录,附录B和附录C为资料性附录。

本标准由中国标准化研究院提出。

本标准由全国统计方法应用标准化技术委员会归口。

本标准起草单位:中国标准化研究院、中国科学院数学与系统科学研究院、宁波工程学院、北京大学、无锡市产品质量监督检验所、北京师范大学。

本标准主要起草人:于振凡、丁文兴、陈敏、荆广珠、房祥忠、吴建国、崔恒建、陈玉忠。

本标准所代替标准的历次版本的发布情况为:

- GB/T 4883—1985。

引 言

从事科学研究、工农业制造以及管理工作都离不开数据,而对这些数据的整理、分析和解释都离不开统计方法。统计学是研究数字资料的整理、分析和正确解释的一门学科。人们各自从不同的来源取得各种数字资料,这些数字资料通常都是杂乱无章的,必须经过整理和简缩才能利用,使用完善的统计方法就可使数据整理、排列的有条有理,用图形或少量的几个重要参数,就可把一大堆数据的特征表达出来,这样既可避免不正确的解释,又可将获得满意数据的成本降到最低限度,提高了经济效益。《数据的统计处理和解释》含有多项国家标准,它们是:

- 统计容忍区间的确定(GB/T 3359)
- 均值的估计和置信区间(GB/T 3360)
- 在成对观测值情形下两个均值的比较(GB/T 3361)
- 二项分布参数的估计与检验(GB/T 4088)
- 泊松分布参数的估计与检验(GB/T 4089)
- 正态性检验(GB/T 4882)
- 正态样本离群值的判断和处理(GB/T 4883)
- 正态分布均值和方差的估计与检验(GB/T 4889)
- 正态分布均值和方差检验的功效(GB/T 4890)
- I型极值分布样本离群值的判断和处理(GB/T 6380)
- 伽玛分布(皮尔逊Ⅲ型分布)的参数估计(GB/T 8055)
- 指数分布样本离群值的判断和处理(GB/T 8056)

对于《数据的统计处理和解释 正态样本离群值的判断和处理》尚无相应的国际标准,但在一些关于测量的国际标准和技术文件中(例如:ISO 5725《测量方法与结果的准确度》、ISO 导则 98《用蒙特卡罗方法评定不确定度》)都采用了本标准中规定的一些正态样本离群值的判断和处理的方法。

数据的统计处理和解释

正态样本离群值的判断和处理

1 范围

本标准适用于来自正态分布的样本中离群值的判断和处理。

2 规范性引用文件

下列文件中的条款通过本标准的引用而成为本标准的条款。凡是注日期的引用文件,其随后所有的修改单(不包括勘误的内容)或修订版均不适用于本标准,然而,鼓励根据本标准达成协议的各方研究是否可使用这些文件的最新版本。凡是不注日期的引用文件,其最新版本适用于本标准。

GB/T 4882—2001 数据的统计处理和解释 正态性检验

GB/T 19000—2000 质量管理体系 基础和术语

ISO 3534-1:2006 统计学词汇及符号 第1部分:一般统计术语与用于概率的术语

ISO 3534-2:2006 统计学词汇及符号 第2部分:应用统计

3 术语、定义和符号

ISO 3534-1:2006、ISO 3534-2:2006 和 GB/T 19000—2000 确定的术语和定义以及下列术语、定义和符号适用于本标准。为便于参考,某些术语直接引自上述标准。

3.1 术语和定义

3.1.1

离群值 outlier

样本中的一个或几个观测值,它们离开其他观测值较远,暗示它们可能来自不同的总体。

注:离群值按显著性的程度分为歧离值和统计离群值。

3.1.2

统计离群值 statistical outlier

在剔除水平下统计检验为显著的离群值。

3.1.3

歧离值 straggler

在检出水平(3.1.4)下显著,但在剔除水平(3.1.5)下不显著的离群值。

3.1.4

检出水平 detection level

为检出离群值而指定的统计检验的显著性水平。

注:除非根据本标准达成协议的各方另有约定, α 值应为0.05。

3.1.5

剔除水平 deletion level

为检出离群值是否高度离群而指定的统计检验的显著性水平。

注:剔除水平 α' 的值应不超过检出水平 α 的值。除非根据本标准达成协议的各方另有约定, α' 值应为0.01。

3.2 符号和缩略语

n 样本量(观测值个数)

\bar{x} 样本均值

α 检验离群值所使用的显著性水平,简称检出水平

- α' 检验统计离群值所使用的显著性水平,简称剔除水平($\alpha' < \alpha$)
- $x_{(i)}$ 观测值自小到大排序后的第 i 个值
- σ 总体标准差
- s 样本标准差
- R_n 奈尔(Nair)上统计量
- R'_n 奈尔(Nair)下统计量
- G_n 格拉布斯(Grubbs)上统计量
- G'_n 格拉布斯(Grubbs)下统计量
- D_n 狄克逊(Dixon)上统计量
- D'_n 狄克逊(Dixon)下统计量
- b_1 偏度统计量
- b_k 峰度统计量

4 离群值判断

4.1 来源与判断

离群值按产生原因分为两类:

- 第一类离群值是总体固有变异性的极端表现,这类离群值与样本中其余观测值属于同一总体;
- 第二类离群值是由于试验条件和试验方法的偶然偏离所产生的结果,或产生于观测、记录、计算中的失误,这类离群值与样本中其余观测值不属于同一总体。

对离群值的判定通常可根据技术上或物理上的理由直接进行,例如当试验者已经知道试验偏离了规定的试验方法,或测试仪器发生问题等。当上述理由不明确时,可用本标准规定的方法。

4.2 离群值的三种情形

本标准在下述不同情形下判断样本中的离群值:

- 上侧情形:根据实际情况或以往经验,离群值都为高端值;
- 下侧情形:根据实际情况或以往经验,离群值都为低端值;
- 双侧情形:根据实际情况或以往经验,离群值可为高端值,也可为低端值。

注:1) 上侧情形和下侧情形统称单侧情形;
2) 若无法认定单侧情形,按双侧情形处理。

4.3 检出离群值个数的上限

应规定在样本中检出离群值个数的上限(与样本量相比应较小),当检出离群值个数超过了这个上限时,对此样本应作慎重的研究和处理。

4.4 单个离群值情形

- 依实际情况或以往经验选定,选定适宜的离群值检验规则(见第6章、第7章、第8章);
- 确定适当的显著性水平;
- 根据显著性水平及样本量,确定检验的临界值;
- 由观测值计算相应统计量的值,根据所得值与临界值的比较结果作出判断。

4.5 判定多个离群值的检验规则

在允许检出离群值的个数大于1的情况下,重复使用4.4规定的检验规则进行检验。若没有检出离群值,则整个检验停止;若检出离群值,当检出的离群值总数超过上限(4.3)时,检验停止,对此样本应慎重处理,否则,采用相同的检出水平和相同的规则,对除去已检出的离群值后余下的观测值继续检验。

5 离群值处理

5.1 处理方式

处理离群值的方式有:

- 保留离群值并用于后续数据处理;

- b) 在找到实际原因时修正离群值,否则予以保留;
- c) 剔除离群值,不追加观测值;
- d) 剔除离群值,并追加新的观测值或用适宜的插补值代替。

5.2 处理规则

对检出的离群值,应尽可能寻找其技术上和物理上的原因,作为处理离群值的依据。应根据实际问题的性质,权衡寻找和判定产生离群值的原因所需代价、正确判定离群值的得益及错误剔除正常观测值的风险,以确定实施下述三个规则之一:

- a) 若在技术上或物理上找到了产生离群值的原因,则应剔除或修正;若未找到产生它的物理上和技术上的原因,则不得剔除或进行修正。
- b) 若在技术上或物理上找到产生离群值的原因,则应剔除或修正;否则,保留歧离值,剔除或修正统计离群值;在重复使用同一检验规则检验多个离群值的情形,每次检出离群值后,都要再检验它是否为统计离群值。若某次检出的离群值为统计离群值,则此离群值及在它前面检出的离群值(含歧离值)都应被剔除或修正。
- c) 检出的离群值(含歧离值)都应被剔除或进行修正。

5.3 备案

被剔除或修正的观测值及其理由应予记录,以备查询。

6 已知标准差情形离群值的判断规则

6.1 一般原则

当已知标准差时,使用奈尔(Nair)检验法,奈尔检验法的样本量 $3 \leq n \leq 100$ 。

6.2 离群值的判断规则

6.2.1 上侧情形

- a) 计算出统计量 R_n 的值:

$$R_n = (x_{(n)} - \bar{x}) / \sigma$$

其中 σ 是已知的总体标准差, \bar{x} 是样本均值, $\bar{x} = (x_1 + \dots + x_n) / n$;

- b) 确定检出水平 α , 在表 A.1 中查出临界值 $R_{1-\alpha}(n)$;
- c) 当 $R_n > R_{1-\alpha}(n)$ 时,判定 $x_{(n)}$ 为离群值,否则判未发现 $x_{(n)}$ 是离群值;
- d) 对于检出的离群值 $x_{(n)}$,确定剔除水平 α^* ,在表 A.1 中查出临界值 $R_{1-\alpha^*}(n)$ 。当 $R_n > R_{1-\alpha^*}(n)$ 时,判定 $x_{(n)}$ 为统计离群值,否则判未发现 $x_{(n)}$ 是统计离群值(即 $x_{(n)}$ 为歧离值)。

6.2.2 下侧情形

- a) 计算出统计量 R'_n 的值:

$$R'_n = (\bar{x} - x_{(1)}) / \sigma$$

其中 σ 是已知的总体标准差, \bar{x} 是样本均值;

- b) 确定检出水平 α , 在表 A.1 中查出临界值 $R_{1-\alpha}(n)$;
- c) 当 $R'_n > R_{1-\alpha}(n)$ 时,判定 $x_{(1)}$ 为离群值,否则判未发现 $x_{(1)}$ 是离群值;
- d) 对于检出的离群值 $x_{(1)}$,确定剔除水平 α^* ,在表 A.1 中查出临界值 $R_{1-\alpha^*}(n)$ 。当 $R'_n > R_{1-\alpha^*}(n)$ 时,判定 $x_{(1)}$ 为统计离群值,否则判未发现 $x_{(1)}$ 是统计离群值(即 $x_{(1)}$ 为歧离值)。

6.2.3 双侧情形

- a) 计算出统计量 R_n 与 R'_n 的值;
- b) 确定检出水平 α , 在表 A.1 中查出临界值 $R_{1-\alpha/2}(n)$;
- c) 当 $R_n > R'_n$, 且 $R_n > R_{1-\alpha/2}(n)$ 时,判定最大值 $x_{(n)}$ 为离群值;当 $R'_n > R_n$, 且 $R'_n > R_{1-\alpha/2}(n)$ 时,判定最小值 $x_{(1)}$ 为离群值;否则判未发现离群值;当 $R_n = R'_n$ 时,同时对最大值和最小值进行检验;

- d) 对于检出的离群值 $x_{(1)}$ 或 $x_{(n)}$, 确定剔除水平 α^* , 在表 A.1 中查出临界值 $R_{1-\alpha^*/2}(n)$, 当 $R'_n > R_{1-\alpha^*/2}(n)$ 时, 判定 $x_{(1)}$ 为统计离群值, 否则未发现 $x_{(1)}$ 是统计离群值 (即 $x_{(1)}$ 为歧离值); 当 $R'_n > R_{1-\alpha^*/2}(n)$ 时, 判定 $x_{(n)}$ 为统计离群值, 否则未发现 $x_{(n)}$ 是统计离群值 (即 $x_{(n)}$ 为歧离值)。

6.3 使用奈尔(Nair)检验法的示例

对某种化纤的纤维干收缩率测试 25 个样品, 其数据经排列后为(单位 %):

3.13	3.49	4.01	4.48	4.61	4.76	4.98	5.25	5.32	5.39	5.42	5.57	5.59
5.59	5.63	5.63	5.65	5.66	5.67	5.69	5.71	6.00	6.03	6.12	6.76	

经验表明这种化纤的纤维干收缩率服从正态分布, 已知 $\sigma = 0.65$, 检查这些数据中是否存在下侧离群值。

规定至多检出三个离群值, 采用 5.2 中 b) 的处理方式。

1) 确定检出水平 $\alpha = 0.05$, 对 25 个样品, 经计算得 $\bar{x} = 5.2856$, $R'_{25} = (\bar{x} - x_{(1)})/\sigma = (5.2856 - 3.13)/0.65 = 3.316$ 。在表 A.1 中查出临界值 $R_{0.95}(25) = 2.815$, 因 $R'_n > R_{0.95}(25)$, 故判定 $x_{(1)} = 3.13$ 是离群值。

对于检出的离群值 $x_{(1)} = 3.13$, 确定剔除水平 $\alpha^* = 0.01$, 在表 A.1 中查出临界值 $R_{0.99}(25) = 3.284$, 因 $R'_n > R_{0.99}(25)$, 故判定 $x_{(1)} = 3.13$ 是统计离群值。

2) 取出观测值为 3.13 的数据后, 在余下的 24 个观测值中计算均值 $\bar{x} = 5.375$, 这时最小值为 $x_{(2)} = 3.49$, 计算得 $R'_{24} = (5.375 - 3.49)/0.65 = 2.90$ 。在表 A.1 中查出临界值 $R_{0.95}(24) = 2.8$, 因 $R'_{24} > R_{0.95}(24)$, 故判定 $x_{(2)} = 3.49$ 是离群值。

对于检出的离群值 $x_{(2)} = 3.49$, 确定剔除水平 $\alpha^* = 0.01$, 在表 A.1 中查出临界值 $R_{0.99}(24) = 3.269$, 因 $R'_{24} < R_{0.99}(24)$, 故判定未发现 $x_{(2)} = 3.49$ 是统计离群值 (即 $x_{(2)} = 3.49$ 为歧离值)。

3) 取出观测值为 3.13、3.49 的数据后, 余下 23 个观测值的样本均值为 5.457, 这时最小值为 $x_{(3)} = 4.01$ 。计算得 $R'_{23} = (5.457 - 4.01)/0.65 = 2.227$, 在表 A.1 中查出临界值 $R_{0.95}(23) = 2.784$, 因 $R'_{23} < R_{0.95}(23)$, 故判定“未发现 $x_{(3)} = 4.01$ 是离群值”。

本例检出 $x_{(1)} = 3.13$ 和 $x_{(2)} = 3.49$ 是离群值, 其中 $x_{(1)} = 3.13$ 是统计离群值, $x_{(2)} = 3.49$ 是歧离值。应参照 5.2 中规定的规则考虑是否剔除。

7 未知标准差情形离群值的判断规则(限定检出离群值的个数不超过 1 时)

7.1 一般原则

在未知标准差的情形下可使用格拉布斯(Grubbs)检验法和狄克逊(Dixon)检验法。可根据实际要求选定其中一种检验法(见附录 B)。

7.2 格拉布斯(Grubbs)检验法

7.2.1 上侧情形

- a) 计算出统计量 G_n 的值:

$$G_n = (x_{(n)} - \bar{x})/s \dots\dots\dots(1)$$

$$s = \left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2} \dots\dots\dots(2)$$

其中 \bar{x} 和 s 是样本均值和样本标准差;

- b) 确定检出水平 α , 在表 A.2 中查出临界值 $G_{1-\alpha}(n)$;
 c) 当 $G_n > G_{1-\alpha}(n)$ 时, 判定 $x_{(n)}$ 为离群值, 否则未发现 $x_{(n)}$ 是离群值;
 d) 对于检出的离群值 $x_{(n)}$, 确定剔除水平 α^* , 在表 A.2 中查出临界值 $G_{1-\alpha^*}(n)$ 。当 $G_n > G_{1-\alpha^*}(n)$ 时, 判定 $x_{(n)}$ 为统计离群值, 否则未发现 $x_{(n)}$ 是统计离群值 (即 $x_{(n)}$ 为歧离值)。

7.2.2 下侧情形

- a) 计算出统计量 G'_n 的值:

$$G'_n = (\bar{x} - x_{(1)})/s \dots\dots\dots(3)$$

$$s = \left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2} \dots\dots\dots(4)$$

其中 \bar{x} 和 s 是样本均值和样本标准差；

- b) 确定检出水平 α , 在表 A. 2 中查出临界值 $G_{1-\alpha}(n)$;
- c) 当 $G'_n > G_{1-\alpha}(n)$ 时, 判定 $x_{(1)}$ 为离群值, 否则判未发现 $x_{(1)}$ 是离群值;
- d) 对于检出的离群值 $x_{(1)}$, 确定剔除水平 α^* , 在表 A. 2 中查出临界值 $G_{1-\alpha^*}(n)$ 。当 $G'_n > G_{1-\alpha^*}(n)$ 时, 判定 $x_{(1)}$ 为统计离群值, 否则判未发现 $x_{(1)}$ 是统计离群值(即 $x_{(1)}$ 为歧离值)。

7.2.3 双侧情形

- a) 计算出统计量 G_n 和 G'_n 的值。
- b) 确定检出水平 α , 在表 A. 2 中查出临界值 $G_{1-\alpha/2}(n)$ 。
- c) 当 $G_n > G'_n$ 且 $G_n > G_{1-\alpha/2}(n)$, 判定 $x_{(n)}$ 为离群值; 当 $G'_n > G_n$ 且 $G'_n > G_{1-\alpha/2}(n)$, 判定 $x_{(1)}$ 为离群值; 否则判未发现离群值。当 $G'_n = G_n$ 时, 应重新考虑限定检出离群值的个数。
- d) 对于检出的离群值 $x_{(1)}$ 或 $x_{(n)}$, 确定剔除水平 α^* , 在表 A. 2 中查出临界值 $G_{1-\alpha^*/2}(n)$, 当 $G'_n > G_{1-\alpha^*/2}(n)$ 时, 判定 $x_{(1)}$ 为统计离群值, 否则判未发现 $x_{(1)}$ 是统计离群值(即 $x_{(1)}$ 为歧离值); 当 $G_n > G_{1-\alpha^*/2}(n)$ 时, 判定 $x_{(n)}$ 为统计离群值, 否则判未发现 $x_{(n)}$ 是统计离群值(即 $x_{(n)}$ 为歧离值)。

7.2.4 使用格拉布斯(Grubbs)检验法的示例

对某种砖的抗压强度测试 10 个样品, 其数据经排列后为(单位: MPa):

4. 7, 5. 4, 6. 0, 6. 5, 7. 3, 7. 7, 8. 2, 9. 0, 10. 1, 14. 0

经验表明这种砖的抗压强度服从正态分布, 检查这些数据中是否存在上侧离群值。

本例中, 样本量 $n=10$, $\bar{x}=7.89$, $s^2=7.312$, $s=2.704$ 。计算得:

$$G_{10} = (x_{(10)} - \bar{x})/s = (14 - 7.89)/2.704 = 2.260$$

确定检出水平 $\alpha=0.05$, 在表 A. 2 中查出临界值 $G_{0.95}(10)=2.176$, 因 $G_{10} > G_{0.95}(10)$, 判定 $x_{(10)} = 14.0$ 为离群值。

对于检出的离群值 $x_{(10)} = 14.0$, 确定剔除水平 $\alpha^* = 0.01$, 在表 A. 2 中查出临界值 $G_{0.99}(10) = 2.410$, 因 $G_{10} < G_{0.99}(10)$, 故判为未发现 $x_{(10)} = 14.0$ 是统计离群值(即 $x_{(10)}$ 为歧离值)。

7.3 狄克逊(Dixon)检验法

当使用狄克逊检验法时, 若样本量 $3 \leq n \leq 30$, 其临界值见表 A. 3; 若样本量 $30 < n \leq 100$, 其检验方法见附录 C。

7.3.1 单侧情形

- a) 计算出下述统计量的值:

样 本 量	检验高端离群值	检验低端离群值
$n: 3 \sim 7$	$D_n = r_{10} = \frac{x_{(n)} - x_{(n-1)}}{x_{(n)} - x_{(1)}}$	$D'_n = r'_{10} = \frac{x_{(2)} - x_{(1)}}{x_{(n)} - x_{(1)}}$
$n: 8 \sim 10$	$D_n = r_{11} = \frac{x_{(n)} - x_{(n-1)}}{x_{(n)} - x_{(2)}}$	$D'_n = r'_{11} = \frac{x_{(2)} - x_{(1)}}{x_{(n-1)} - x_{(1)}}$
$n: 11 \sim 13$	$D_n = r_{21} = \frac{x_{(n)} - x_{(n-2)}}{x_{(n)} - x_{(2)}}$	$D'_n = r'_{21} = \frac{x_{(3)} - x_{(1)}}{x_{(n-1)} - x_{(1)}}$
$n: 14 \sim 30$	$D_n = r_{22} = \frac{x_{(n)} - x_{(n-2)}}{x_{(n)} - x_{(3)}}$	$D'_n = r'_{22} = \frac{x_{(3)} - x_{(1)}}{x_{(n-2)} - x_{(1)}}$

- b) 确定检出水平 α , 在表 A. 3 中查出临界值 $D_{1-\alpha}(n)$ 。
- c) 检验高端值, 当 $D_n > D_{1-\alpha}(n)$ 时, 判定 $x_{(n)}$ 为离群值; 检验低端值, 当 $D'_n > D_{1-\alpha}(n)$ 时, 判定 $x_{(1)}$ 为离群值; 否则判未发现离群值。

- d) 对于检出的离群值 $x_{(1)}$ 或 $x_{(n)}$, 确定剔除水平 α^* , 在表 A. 3 中查出临界值 $D_{1-\alpha^*}(n)$ 。检验高端值, 当 $D_n > D_{1-\alpha^*}(n)$ 时, 判定 $x_{(n)}$ 为统计离群值, 否则判未发现 $x_{(n)}$ 是统计离群值 (即 $x_{(n)}$ 为歧离值); 检验低端值, 当 $D'_n > D_{1-\alpha^*}(n)$ 时, 判定 $x_{(1)}$ 为统计离群值, 否则判未发现 $x_{(1)}$ 是统计离群值 (即 $x_{(1)}$ 为歧离值)。

7.3.2 双侧情形

- a) 计算出统计量 D_n 与 D'_n 的值, 这里 D_n 与 D'_n 由 7.3.1 的 a) 给出;
- b) 确定检出水平 α , 在表 A. 3' 中查出临界值 $\tilde{D}_{1-\alpha}(n)$;
- c) 当 $D_n > D'_n, D_n > \tilde{D}_{1-\alpha}(n)$ 时, 判定 $x_{(n)}$ 为离群值; 当 $D'_n > D_n, D'_n > \tilde{D}_{1-\alpha}(n)$ 时, 判定 $x_{(1)}$ 为离群值; 否则判未发现离群值;
- d) 对于检出的离群值 $x_{(1)}$ 或 $x_{(n)}$, 确定剔除水平 α^* , 在表 A. 3' 中查出临界值 $\tilde{D}_{1-\alpha^*}(n)$ 。当 $D_n > D'_n$ 且 $D_n > \tilde{D}_{1-\alpha^*}(n)$ 时, 判定 $x_{(n)}$ 为统计离群值, 否则判未发现 $x_{(n)}$ 是统计离群值 (即 $x_{(n)}$ 为歧离值); 当 $D'_n > D_n$ 且 $D'_n > \tilde{D}_{1-\alpha^*}(n)$ 时, 判定 $x_{(1)}$ 为统计离群值, 否则判未发现 $x_{(1)}$ 是统计离群值 (即 $x_{(1)}$ 为歧离值)。

7.3.3 使用狄克逊(Dixon)检验法的示例

射击 16 发子弹, 射程数据经排列后为(单位: m):

1 125	1 248	1 250	1 259	1 273	1 279	1 285	1 285
1 293	1 300	1 305	1 312	1 315	1 324	1 325	1 350

经验表明子弹射程服从正态分布, 根据实际中的关注不同, 分别对低端值和高端值进行检验。

- a) 检验低端值 $x_{(1)} = 1 125$ 是否为离群值

本例中, 样本量 $n = 16$, 计算

$$D'_{16} = r'_{22} = \frac{x_{(3)} - x_{(1)}}{x_{(14)} - x_{(1)}} = \frac{1 250 - 1 125}{1 324 - 1 125} = \frac{125}{189} = 0.661 4$$

确定检出水平 $\alpha = 0.05$, 在表 A. 3 中查出临界值 $D_{0.95}(16) = 0.505$, 因 $D'_{16} > D_{0.95}(16)$, 故判定最小值 $x_{(1)} = 1 125$ 为离群值。

对于检出的离群值 $x_{(1)} = 1 125$, 确定剔除水平 $\alpha^* = 0.01$, 在表 A. 3 中查出临界值 $D_{0.99}(16) = 0.597$, 因 $D'_{16} > D_{0.99}(16)$, 故判定最小值 $x_{(1)} = 1 125$ 为统计离群值。

- b) 双侧情形

计算 $D'_{16} = 0.661 4$ 和

$$D_{16} = r_{22} = \frac{x_{(16)} - x_{(14)}}{x_{(16)} - x_{(3)}} = \frac{1 350 - 1 324}{1 350 - 1 250} = \frac{26}{100} = 0.26$$

确定检出水平 $\alpha = 0.05$, 在表 A. 3' 查出临界值 $\tilde{D}_{0.95}(16) = 0.547$ 。因 $D'_{16} > D_{16}$ 且 $D'_{16} > \tilde{D}_{0.95}(16)$, 故判定最小值 $x_{(1)} = 1 125$ 为离群值。

对于检出的离群值 $x_{(1)} = 1 125$, 确定剔除水平 $\alpha^* = 0.01$, 在表 A. 3' 查出临界值 $\tilde{D}_{0.99}(16) = 0.627$ 。因 $D'_{16} > D_{16}$ 且 $D'_{16} > \tilde{D}_{0.99}(16)$, 故判定最小值 $x_{(1)} = 1 125$ 为统计离群值。

8 未知标准差情形离群值的判断规则(限定检出离群值的个数大于 1 时)

8.1 一般原则

当限定检出离群值的个数大于 1 时, 可使用偏度—峰度检验法或狄克逊(Dixon)检验法的重复使用方法, 可根据实际要求选定其中一种检验法(见附录 B)。

8.2 偏度—峰度检验法

8.2.1 使用条件

考查样本诸观测值, 确认它们的样本主体来自正态总体, 而极端值应较明显的偏离样本主体。

8.2.2 单侧情形——偏度检验法

a) 计算偏度统计量 b_1 的值

$$b_1 = \frac{\sqrt{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^{3/2}} = \frac{\sqrt{n} \left[\sum_{i=1}^n x_i^3 - 3\bar{x} \sum_{i=1}^n x_i^2 + 2n(\bar{x})^3 \right]}{\left[\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right]^{3/2}} \dots\dots\dots (5)$$

- b) 确定检出水平 α , 在表 A.4 中查出临界值 $b_{1-\alpha}(n)$ 。
- c) 对上侧情形, 当 $b_1 > b_{1-\alpha}(n)$ 时, 判定最大值 $x_{(n)}$ 为离群值; 否则判未发现 $x_{(n)}$ 是离群值; 对下侧情形, 当 $-b_1 > b_{1-\alpha}(n)$ 时, 判定最小值 $x_{(1)}$ 为离群值; 否则判未发现 $x_{(1)}$ 是离群值。
- d) 对于检出的离群值 $x_{(1)}$ 或 $x_{(n)}$, 确定剔除水平 α^* , 在表 A.4 中查出临界值 $b_{1-\alpha^*}(n)$ 。对上侧情形, 当 $b_1 > b_{1-\alpha^*}(n)$ 时, 判定 $x_{(n)}$ 为统计离群值, 否则判未发现 $x_{(n)}$ 是统计离群值 (即 $x_{(n)}$ 为歧离值); 对下侧情形, 当 $-b_1 > b_{1-\alpha^*}(n)$ 时, 判定 $x_{(1)}$ 为统计离群值, 否则判未发现 $x_{(1)}$ 是统计离群值 (即 $x_{(1)}$ 为歧离值)。

8.2.3 双侧情形——峰度检验法

a) 计算峰度统计量 b_k 的值

$$b_k = \frac{n \sum_{i=1}^n (x_i - \bar{x})^4}{\left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^2} = \frac{n \left[\sum_{i=1}^n x_i^4 - 4\bar{x} \sum_{i=1}^n x_i^3 + 6\bar{x}^2 \sum_{i=1}^n x_i^2 - 3n\bar{x}^4 \right]}{\left[\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right]^2} \dots\dots\dots (6)$$

- b) 确定检出水平 α , 在表 A.5 中查出临界值 $b'_{1-\alpha}(n)$ 。
- c) 当 $b_k > b'_{1-\alpha}(n)$ 时, 判定离均值 \bar{x} 最远的观测值为离群值; 否则判未发现离群值。
- d) 对于检出的离群值, 确定剔除水平 α^* , 在表 A.5 中查出临界值 $b'_{1-\alpha^*}(n)$ 。当 $b_k > b'_{1-\alpha^*}(n)$ 时, 判定离均值 \bar{x} 最远的观测值为统计离群值, 否则判未发现该离群值是统计离群值 (即该离群值为歧离值)。

8.2.4 重复使用峰度检验法的示例

本例为离群值问题早期研究中的著名实例(1883年)。观测金星垂直半径的15个观测数据的离差经排列后为(单位: s)。

-1.40	-0.44	-0.30	-0.24	-0.22	-0.13	-0.05	0.06
0.10	0.18	0.20	0.39	0.48	0.63	1.01	

由问题的背景需要判断 $x_{(1)} = -1.40$ 和 $x_{(15)} = 1.01$ 是否离群。

根据 GB/T 4882—2001, 使用正态概率纸进行正态性检验。

将上述数据点在正态概率纸上 (见图 1), 此时, 样本的诸点近似在一条直线近旁两侧, 当画出适宜的直线后, 样本的低端向上而高端向下偏离, 故可用偏度—峰度检验法。

计算得:

$$\begin{array}{cccc} \sum_{i=1}^{15} x_i & \sum_{i=1}^{15} x_i^2 & \sum_{i=1}^{15} x_i^3 & \sum_{i=1}^{15} x_i^4 \\ 0.27 & 4.254\ 5 & -1.417\ 671 & 5.170\ 248\ 05 \\ \bar{x} = 0.27/15 = 0.018, b_k = 4.386 \end{array}$$

确定检出水平 $\alpha = 0.05$, 在表 A.5 中查出临界值 $b'_{0.95}(15) = 4.13$, 因 $b_k = 4.386 > b'_{0.95}(15) = 4.13$, 判定距离均值 0.018 最远的 $x_{(1)} = -1.40$ 为离群值。

对于检出的离群值 $x_{(1)} = -1.40$, 确定剔除水平 $\alpha^* = 0.01$, 在表 A.5 中查出临界值 $b'_{0.99}(15) = 5.30$, 因 $b_k = 4.386 < b'_{0.99}(15) = 5.30$, 故判未发现该离群值 $x_{(1)} = -1.40$ 是统计离群值 (即 $x_{(1)} = -1.40$ 为歧离值)。

取出 $x_{(1)} = -1.40$ 之后,对余下 14 个值进行计算如下:

$\sum_{i=1}^{14} x_i$	$\sum_{i=1}^{14} x_i^2$	$\sum_{i=1}^{14} x_i^3$	$\sum_{i=1}^{14} x_i^4$
0.27	4.254 5	-1.417 671	5.170 248 05
+1.40	-1.960 0	+2.744 000	-3.841 600 00
1.67	2.294 5	1.326 329	1.328 648 05

$\bar{x} = 1.67/14 = 0.119 3$,再计算 $b_k = 2.816 4$ 。确定检出水平 $\alpha = 0.05$,在表 A. 5 中查出临界值 $b'_{0.95}(14) = 4.11$,而 $b_k = 2.816 4 < b'_{0.95}(14) = 4.11$,故不能再检出离群值。

所以,本例只检出一个歧离值 $x_{(1)} = -1.40$ 。

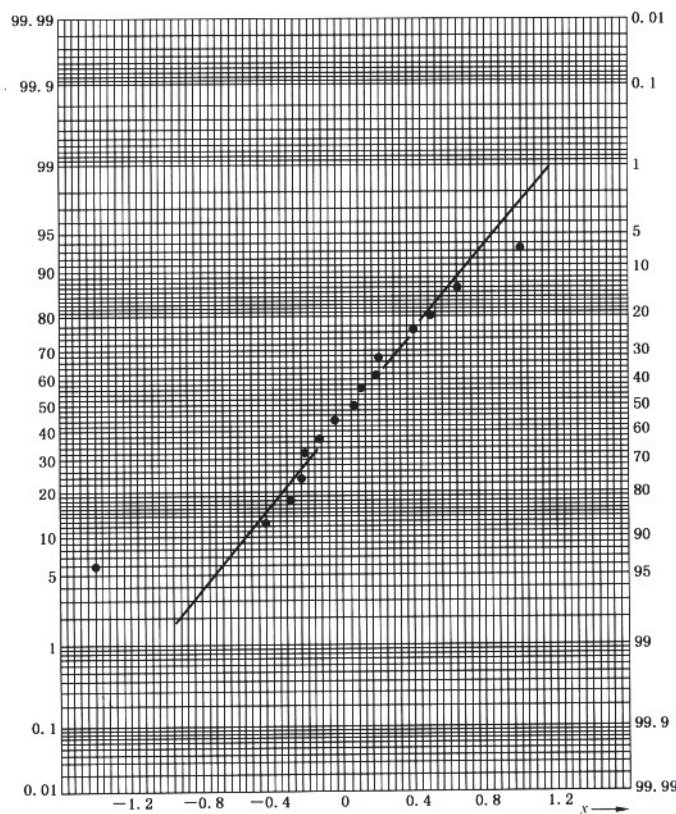


图 1 正态概率纸

8.3 狄克逊(Dixon)检验法

8.3.1 狄克逊(Dixon)检验法的规则

见 7.3。

8.3.2 重复使用狄克逊(Dixon)检验法的示例

数据同 8.2.4。计算

$$D_{15} = r_{22} = \frac{x_{(15)} - x_{(13)}}{x_{(15)} - x_{(3)}} = \frac{1.01 - 0.48}{1.01 + 0.30} = \frac{0.53}{1.31} = 0.406$$

$$D'_{15} = r'_{22} = \frac{x_{(3)} - x_{(1)}}{x_{(13)} - x_{(1)}} = \frac{-0.30 + 1.40}{0.48 + 1.40} = \frac{1.10}{1.88} = 0.585$$

对双侧问题,确定检出水平 $\alpha = 0.05$,在表 A. 3' 中查出临界值 $\tilde{D}_{0.95}(15) = 0.565$,由于 $D'_{15} > D_{15}$ 且 $D'_{15} > \tilde{D}_{0.95}(15)$,故判定最小值 $x_{(1)} = -1.40$ 为离群值。

对于检出的离群值 $x_{(1)} = -1.40$,确定剔除水平 $\alpha^* = 0.01$,在表 A. 3' 中查出临界值 $\tilde{D}_{0.99}(15) = 0.646$ 。因为 $D'_{15} < \tilde{D}_{0.99}(15)$,故判未发现 $x_{(1)} = -1.40$ 是统计离群值(即 $x_{(1)} = -1.40$ 为歧离值)。

取出这个观测值后还剩余 14 个值($n=14$),使用

$$D_{14} = r_{22} = \frac{x_{(14)} - x_{(12)}}{x_{(14)} - x_{(3)}} = \frac{1.01 - 0.48}{1.01 + 0.24} = \frac{0.53}{1.25} = 0.424$$

$$D'_{14} = r'_{22} = \frac{x_{(3)} - x_{(1)}}{x_{(12)} - x_{(1)}} = \frac{-0.24 + 0.44}{0.48 + 0.44} = \frac{0.20}{0.92} = 0.217$$

对于上述确定的检出水平 $\alpha = 0.05$,在表 A. 3' 中查出临界值 $\tilde{D}_{0.95}(14) = 0.586$,因为 $D'_{14} < \tilde{D}_{0.95}(14)$,故不能继续检出离群值。

所以,本例只检出一个歧离值 $x_{(1)} = -0.140$ 。

附 录 A
(规范性附录)
统 计 数 值 表

奈尔(Nair)检验的临界值表见表 A.1, 格拉布斯(Grubbs)检验的临界值表见表 A.2, 狄克逊(Dixon)检验的临界值表见表 A.3, 偏度检验的临界值表见表 A.4, 峰度检验的临界值表见表 A.5。

表 A.1 奈尔(Nair)检验的临界值表

<i>n</i>	0.90	0.95	0.975	0.99	0.995	<i>n</i>	0.90	0.95	0.975	0.99	0.995
3	1.497	1.738	1.955	2.215	2.396	36	2.722	2.944	3.150	3.403	3.584
4	1.696	1.941	2.163	2.431	2.618	37	2.732	2.953	3.159	3.412	3.592
5	1.835	2.080	2.304	2.574	2.764	38	2.741	2.962	3.167	3.420	3.600
6	1.939	2.184	2.408	2.679	2.870	39	2.750	2.971	3.176	3.428	3.608
7	2.022	2.267	2.490	2.761	2.952	40	2.759	2.980	3.184	3.436	3.616
8	2.091	2.334	2.557	2.828	3.019	41	2.768	2.988	3.192	3.444	3.623
9	2.150	2.392	2.613	2.884	3.074	42	2.776	2.996	3.200	3.451	3.630
10	2.200	2.441	2.662	2.931	3.122	43	2.784	3.004	3.207	3.458	3.637
11	2.245	2.484	2.704	2.973	3.163	44	2.792	3.011	3.215	3.465	3.644
12	2.284	2.523	2.742	3.010	3.199	45	2.800	3.019	3.222	3.472	3.651
13	2.320	2.557	2.776	3.043	3.232	46	2.808	3.026	3.229	3.479	3.657
14	2.352	2.589	2.806	3.072	3.261	47	2.815	3.033	3.235	3.485	3.663
15	2.382	2.617	2.834	3.099	3.287	48	2.822	3.040	3.242	3.491	3.669
16	2.409	2.644	2.860	3.124	3.312	49	2.829	3.047	3.249	3.498	3.675
17	2.434	2.668	2.883	3.147	3.334	50	2.836	3.053	3.255	3.504	3.681
18	2.458	2.691	2.905	3.168	3.355	51	2.843	3.060	3.261	3.509	3.687
19	2.480	2.712	2.926	3.188	3.374	52	2.849	3.066	3.267	3.515	3.692
20	2.500	2.732	2.945	3.207	3.392	53	2.856	3.072	3.273	3.521	3.698
21	2.519	2.750	2.963	3.224	3.409	54	2.862	3.078	3.279	3.526	3.703
22	2.538	2.768	2.980	3.240	3.425	55	2.868	3.084	3.284	3.532	3.708
23	2.555	2.784	2.996	3.256	3.440	56	2.874	3.090	3.290	3.537	3.713
24	2.571	2.800	3.011	3.270	3.455	57	2.880	3.095	3.295	3.542	3.718
25	2.587	2.815	3.026	3.284	3.468	58	2.886	3.101	3.300	3.547	3.723
26	2.602	2.829	3.039	3.298	3.481	59	2.892	3.106	3.306	3.552	3.728
27	2.616	2.843	3.053	3.310	3.493	60	2.897	3.112	3.311	3.557	3.733
28	2.630	2.856	3.065	3.322	3.505	61	2.903	3.117	3.316	3.562	3.737
29	2.643	2.869	3.077	3.334	3.516	62	2.908	3.122	3.321	3.566	3.742
30	2.656	2.881	3.089	3.345	3.527	63	2.913	3.127	3.326	3.571	3.746
31	2.668	2.892	3.100	3.356	3.538	64	2.919	3.132	3.330	3.575	3.751
32	2.679	2.903	3.111	3.366	3.548	65	2.924	3.137	3.335	3.580	3.755
33	2.690	2.914	3.121	3.376	3.557	66	2.929	3.142	3.339	3.584	3.759
34	2.701	2.924	3.131	3.385	3.566	67	2.934	3.146	3.344	3.588	3.763
35	2.712	2.934	3.140	3.394	3.575	68	2.938	3.151	3.348	3.593	3.767
						69	2.943	3.155	3.353	3.597	3.771
						70	2.948	3.160	3.357	3.601	3.775

表 A.1 (续)

<i>n</i>	0.90	0.95	0.975	0.99	0.995	<i>n</i>	0.90	0.95	0.975	0.99	0.995
71	2.952	3.164	3.361	3.605	3.779	86	3.014	3.223	3.417	3.658	3.831
72	2.957	3.169	3.365	3.609	3.783	87	3.017	3.226	3.421	3.661	3.834
73	2.961	3.173	3.369	3.613	3.787	88	3.021	3.230	3.424	3.665	3.837
74	2.966	3.177	3.373	3.617	3.791	89	3.024	3.233	3.427	3.668	3.840
75	2.970	3.181	3.377	3.620	3.794	90	3.028	3.236	3.430	3.671	3.843
76	2.974	3.185	3.381	3.624	3.798	91	3.031	3.240	3.433	3.674	3.846
77	2.978	3.189	3.385	3.628	3.801	92	3.035	3.243	3.437	3.677	3.849
78	2.983	3.193	3.389	3.631	3.805	93	3.038	3.246	3.440	3.680	3.852
79	2.987	3.197	3.393	3.635	3.808	94	3.042	3.249	3.443	3.683	3.854
80	2.991	3.201	3.396	3.638	3.812	95	3.045	3.253	3.446	3.685	3.857
81	2.995	3.205	3.400	3.642	3.815	96	3.048	3.256	3.449	3.688	3.860
82	2.999	3.208	3.403	3.645	3.818	97	3.052	3.259	3.452	3.691	3.863
83	3.002	3.212	3.407	3.648	3.821	98	3.055	3.262	3.455	3.694	3.865
84	3.006	3.216	3.410	3.652	3.825	99	3.058	3.265	3.458	3.697	3.868
85	3.010	3.219	3.414	3.655	3.828	100	3.061	3.268	3.460	3.699	3.871

表 A.2 格拉布斯(Grubbs)检验的临界值表

<i>n</i>	0.90	0.95	0.975	0.99	0.995	<i>n</i>	0.90	0.95	0.975	0.99	0.995
3	1.148	1.153	1.155	1.155	1.155	26	2.502	2.681	2.841	3.029	3.157
4	1.425	1.463	1.481	1.492	1.496	27	2.519	2.698	2.859	3.049	3.178
5	1.602	1.672	1.715	1.749	1.764	28	2.534	2.714	2.876	3.068	3.199
6	1.729	1.822	1.887	1.944	1.973	29	2.549	2.730	2.893	3.085	3.218
7	1.828	1.938	2.020	2.097	2.139	30	2.563	2.745	2.908	3.103	3.236
8	1.909	2.032	2.126	2.221	2.274	31	2.577	2.759	2.924	3.119	3.253
9	1.977	2.110	2.215	2.323	2.387	32	2.591	2.773	2.938	3.135	3.270
10	2.036	2.176	2.290	2.410	2.482	33	2.604	2.786	2.952	3.150	3.286
11	2.088	2.234	2.355	2.485	2.564	34	2.616	2.799	2.965	3.164	3.301
12	2.134	2.285	2.412	2.550	2.636	35	2.628	2.811	2.979	3.178	3.316
13	2.175	2.331	2.462	2.607	2.699	36	2.639	2.823	2.991	3.191	3.330
14	2.213	2.371	2.507	2.659	2.755	37	2.650	2.835	3.003	3.204	3.343
15	2.247	2.409	2.549	2.705	2.806	38	2.661	2.846	3.014	3.216	3.356
16	2.279	2.443	2.585	2.747	2.852	39	2.671	2.857	3.025	3.228	3.369
17	2.309	2.475	2.620	2.785	2.894	40	2.682	2.866	3.036	3.240	3.381
18	2.335	2.504	2.651	2.821	2.932	41	2.692	2.877	3.046	3.251	3.393
19	2.361	2.532	2.681	2.854	2.968	42	2.700	2.887	3.057	3.261	3.404
20	2.385	2.557	2.709	2.884	3.001	43	2.710	2.896	3.067	3.271	3.415
21	2.408	2.580	2.733	2.912	3.031	44	2.719	2.905	3.075	3.282	3.425
22	2.429	2.603	2.758	2.939	3.060	45	2.727	2.914	3.085	3.292	3.435
23	2.448	2.624	2.781	2.963	3.087	46	2.736	2.923	3.094	3.302	3.445
24	2.467	2.644	2.802	2.987	3.112	47	2.744	2.931	3.103	3.310	3.455
25	2.486	2.663	2.822	3.009	3.135	48	2.753	2.940	3.111	3.319	3.464
						49	2.760	2.948	3.120	3.329	3.474
						50	2.768	2.956	3.128	3.336	3.483

表 A.2 (续)

<i>n</i>	0.90	0.95	0.975	0.99	0.995	<i>n</i>	0.90	0.95	0.975	0.99	0.995
51	2.775	2.964	3.136	3.345	3.491	76	2.922	3.111	3.287	3.502	3.654
52	2.783	2.971	3.143	3.353	3.500	77	2.927	3.117	3.291	3.507	3.658
53	2.790	2.978	3.151	3.361	3.507	78	2.931	3.121	3.297	3.511	3.663
54	2.798	2.986	3.158	3.368	3.516	79	2.935	3.125	3.301	3.516	3.669
55	2.804	2.992	3.166	3.376	3.524	80	2.940	3.130	3.305	3.521	3.673
56	2.811	3.000	3.172	3.383	3.531	81	2.945	3.134	3.309	3.525	3.677
57	2.818	3.006	3.180	3.391	3.539	82	2.949	3.139	3.315	3.529	3.682
58	2.824	3.013	3.186	3.397	3.546	83	2.953	3.143	3.319	3.534	3.687
59	2.831	3.019	3.193	3.405	3.553	84	2.957	3.147	3.323	3.539	3.691
60	2.837	3.025	3.199	3.411	3.560	85	2.961	3.151	3.327	3.543	3.695
61	2.842	3.032	3.205	3.418	3.566	86	2.966	3.155	3.331	3.547	3.699
62	2.849	3.037	3.212	3.424	3.573	87	2.970	3.160	3.335	3.551	3.704
63	2.854	3.044	3.218	3.430	3.579	88	2.973	3.163	3.339	3.555	3.708
64	2.860	3.049	3.224	3.437	3.586	89	2.977	3.167	3.343	3.559	3.712
65	2.866	3.055	3.230	3.442	3.592	90	2.981	3.171	3.347	3.563	3.716
66	2.871	3.061	3.235	3.449	3.598	91	2.984	3.174	3.350	3.567	3.720
67	2.877	3.066	3.241	3.454	3.605	92	2.989	3.179	3.355	3.570	3.725
68	2.883	3.071	3.246	3.460	3.610	93	2.993	3.182	3.358	3.575	3.728
69	2.888	3.076	3.252	3.466	3.617	94	2.996	3.186	3.362	3.579	3.732
70	2.893	3.082	3.257	3.471	3.622	95	3.000	3.189	3.365	3.582	3.736
71	2.897	3.087	3.262	3.476	3.627	96	3.003	3.193	3.369	3.586	3.739
72	2.903	3.092	3.267	3.482	3.633	97	3.006	3.196	3.372	3.589	3.744
73	2.908	3.098	3.272	3.487	3.638	98	3.011	3.201	3.377	3.593	3.747
74	2.912	3.102	3.278	3.492	3.643	99	3.014	3.204	3.380	3.597	3.750
75	2.917	3.107	3.282	3.496	3.648	100	3.017	3.207	3.383	3.600	3.754

表 A.3 单侧狄克逊(Dixon)检验的临界值表

<i>n</i>	统 计 量	0.90	0.95	0.99	0.995
3		0.885	0.941	0.988	0.994
4		0.679	0.765	0.889	0.920
5	$r_{10} = \frac{x_{(n)} - x_{(n-1)}}{x_{(n)} - x_{(1)}} \text{ 或 } r'_{10} = \frac{x_{(2)} - x_{(1)}}{x_{(n)} - x_{(1)}}$	0.557	0.642	0.782	0.823
6		0.484	0.562	0.698	0.744
7		0.434	0.507	0.637	0.680
8		0.479	0.554	0.681	0.723
9	$r_{11} = \frac{x_{(n)} - x_{(n-1)}}{x_{(n)} - x_{(2)}} \text{ 或 } r'_{11} = \frac{x_{(2)} - x_{(1)}}{x_{(n-1)} - x_{(1)}}$	0.441	0.512	0.635	0.676
10		0.410	0.477	0.597	0.638

表 A.3 (续)

n	统计量	0.90	0.95	0.99	0.995
11	$r_{21} = \frac{x_{(n)} - x_{(n-2)}}{x_{(n)} - x_{(2)}} \text{ 或 } r'_{21} = \frac{x_{(3)} - x_{(1)}}{x_{(n-1)} - x_{(1)}}$	0.517	0.575	0.674	0.707
12		0.490	0.546	0.642	0.675
13		0.467	0.521	0.617	0.649
14	$r_{22} = \frac{x_{(n)} - x_{(n-2)}}{x_{(n)} - x_{(3)}} \text{ 或 } r'_{22} = \frac{x_{(3)} - x_{(1)}}{x_{(n-2)} - x_{(1)}}$	0.491	0.546	0.640	0.672
15		0.470	0.524	0.618	0.649
16		0.453	0.505	0.597	0.629
17		0.437	0.489	0.580	0.611
18		0.424	0.475	0.564	0.595
19		0.412	0.462	0.550	0.580
20		0.401	0.450	0.538	0.568
21		0.391	0.440	0.526	0.556
22		0.382	0.431	0.516	0.545
23		0.374	0.422	0.507	0.536
24		0.367	0.413	0.497	0.526
25		0.360	0.406	0.489	0.519
26		0.353	0.399	0.482	0.510
27		0.347	0.393	0.474	0.503
28		0.341	0.387	0.468	0.496
29		0.337	0.381	0.462	0.489
30	0.332	0.376	0.456	0.484	

表 A.3' 双侧狄克逊(Dixon)检验的临界值表

n	统计量	0.95	0.99	n	统计量	0.95	0.99
3	r_{10} 和 r'_{10} 中较大者	0.970	0.994	17	r_{22} 和 r'_{22} 中较大者	0.527	0.614
4		0.829	0.926	18		0.513	0.602
5		0.710	0.821	19		0.500	0.582
6		0.628	0.740	20		0.488	0.570
7	0.569	0.680	21	0.479		0.560	
8	r_{11} 和 r'_{11} 中较大者	0.608	0.717	22		0.469	0.548
9		0.564	0.672	23		0.460	0.537
10		0.530	0.635	24		0.449	0.522
11	r_{21} 和 r'_{21} 中较大者	0.619	0.709	25		0.441	0.518
12		0.583	0.660	26		0.436	0.509
13		0.557	0.638	27		0.427	0.504
14	r_{22} 和 r'_{22} 中较大者	0.587	0.669	28		0.420	0.497
15		0.565	0.646	29		0.415	0.489
16		0.547	0.629	30		0.409	0.480

表 A.4 偏度检验的临界值表

n	0.95	0.99	n	0.95	0.99
8	0.99	1.42	40	0.59	0.87
9	0.97	1.41	45	0.56	0.82
10	0.95	1.39	50	0.53	0.79
12	0.91	1.34	60	0.49	0.72
15	0.85	1.26	70	0.46	0.67
20	0.77	1.15	80	0.43	0.63
25	0.71	1.06	90	0.41	0.60
30	0.66	0.98	100	0.39	0.57
35	0.62	0.92			

表 A.5 峰度检验的临界值表

n	0.95	0.99	n	0.95	0.99
8	3.70	4.53	40	4.05	5.02
9	3.86	4.82	45	4.02	4.94
10	3.95	5.00	50	3.99	4.87
12	4.05	5.20	60	3.93	4.73
15	4.13	5.30	70	3.88	4.62
20	4.17	5.38	80	3.84	4.52
25	4.14	5.29	90	3.80	4.45
30	4.11	5.20	100	3.77	4.37
35	4.08	5.11			

附录 B

(资料性附录)

选择离群值判断方法和处理规则的指南

B.1 判定和处理离群值的目的

B.1.1 三种不同的目的

B.1.1.1 识别与诊断

主要目的是找出离群值,从而进行质量控制、新规律探索、技术考察等工作。

B.1.1.2 估计参数

主要目的在于估计总体的某个参数,寻找离群值的目的在于确定这些值是否计入样本,以便准确估计其参数。

B.1.1.3 检验假设

主要目的在于判定总体是否符合所考察的要求,寻找离群值的目的主要在于确定这些值是否计入样本,以使判定结果计量准确。

B.1.2 判定离群值的不同目的引起的不同的选择

B.1.2.1 以识别为目的

选择判断离群值的主要标准在于判定准确性,要根据所判定错误带来的风险不同,选择适宜的规则。

B.1.2.2 以估计和检验为目的

要判定离群值,就应把判定和处理离群值的方法和进一步作估计或检验的准确性统一起来考虑。如使用格拉布斯(Grubbs)检验法作估计,实际是一种新估计量

$$\hat{\mu} = \begin{cases} (x_1 + \cdots + x_n) / n, & \text{当 } G_n \leq G_{1-\alpha} \cdot (n) \\ (x_{(1)} + \cdots + x_{(n-1)}) / (n-1), & \text{当 } G_n > G_{1-\alpha} \cdot (n) \end{cases}$$

有时也可以不经过判定离群值的步骤,而采用稳健的方法。

例如:在塑料材料中,有时使用截割均值,把12个观测值的最大值与最小值舍去,以余下的10个观测值作算术平均以估计 μ 。(体操比赛评分时,也把诸裁判报出的最高分和最低分舍去,以余下的几个评分的平均值报出),并不需要追查舍去的一定是离群值,而这种估计也很好地预防了离群值的不利影响。

B.2 对各种检验法的选择

本标准第7章、第8章给出了三种检验法,在选用检验方法时应主要考虑下述几点。

B.2.1 限定检出离群值的个数不超过1时

B.2.1.1 当 n 较小时,格拉布斯(Grubbs)检验法具有判定离群值的功效最优性,而狄克逊(Dixon)检验法正确判定离群值的功效与格拉布斯(Grubbs)检验法相差甚微;建议使用格拉布斯(Grubbs)检验法。

B.2.1.2 当 n 较大时,同时在正态概率纸上,若样本主体是基本在一条直线的近旁;建议使用偏度—峰度检验法。

B.2.1.3 当 n 较大时,同时在正态概率纸上,若样本主体不是基本在一条直线的近旁,使用格拉布斯(Grubbs)检验法。

B.2.2 限定检出离群值的个数大于1时

重复使用同一检验法可能犯判多为少(只检出一部分离群值)的错误,而不易犯判少为多(错将一部分非离群的观测值判为离群值)的错误。这两类错误的概率以重复使用偏度—峰度检验法为少(可以证

明,它也具有正确判定离群值的功效优良性)。但计算相对复杂得多;重复使用狄克逊(Dixon)检验法的效果次之,而重复使用格拉布斯(Grubbs)的功效则较差。

偏度—峰度检验法又是正态性检验的优良检验法,不来自正态分布的样本都可能被它拒绝。但这不只是正态样本主体加离群值的模型,所以使用偏度—峰度检验法时,要满足规定的使用条件。比如,在正态概率纸上,若样本主体不是基本在一条直线的近邻两侧;或是样本主体基本在一条直线近邻两侧,而高端值相对于这条直线而言向上而不是向下偏离(如图 B.1),低端值相对于这条直线而言向下而不是向上偏离,则采用偏度—峰度检验法就可能把一部分非离群观测值误判为离群值。

B.2.2.1 当 n 较小时,重复使用狄克逊检验法。

B.2.2.2 当 n 较大时,且在正态概率纸上,若样本主体是基本在一条直线的近旁;可重复偏度—峰度检验法。

B.2.2.3 当 n 较大时,且在正态概率纸上,若样本主体不是基本在一条直线的近旁,建议重复使用格拉布斯检验法。

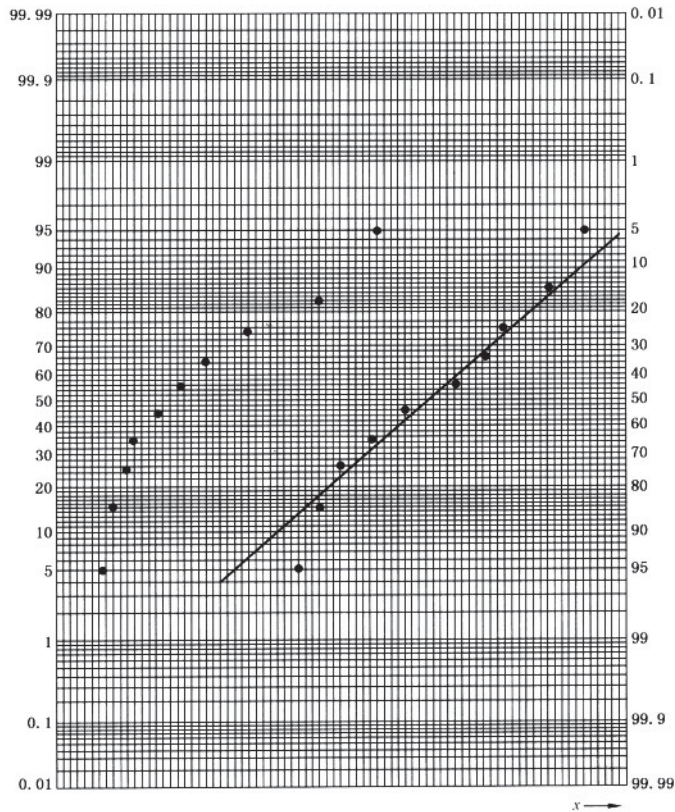


图 B.1 正态概率纸

B.3 重视检出的离群值给出的信息

B.3.1 在一段时间后,考察检出离群值的全体,往往能明显地发现其物理原因和系统倾向,如离群值出自某个测试者为多,说明此人的操作有系统偏离。

B.3.2 若各个样本中出现离群值较为经常,又常不能明确其物理原因,则应怀疑分布的正态性假设。此时可更细微的确定统计分布及选择适宜的统计量形式。

因此,标准使用者应完善判定和处理离群值的记录,并作定期分析。

附 录 C
(资料性附录)
当 $n > 30$ 时的狄克逊(Dixon)检验

C.1 单侧情形

a) 计算出下述统计量的值

样 本 量	检验高端离群值	检验低端离群值
$n: 31 \sim 100$	$D_n = r_{22} = \frac{x_{(n)} - x_{(n-2)}}{x_{(n)} - x_{(3)}}$	$D'_n = r'_{22} = \frac{x_{(3)} - x_{(1)}}{x_{(n-2)} - x_{(1)}}$

- b) 确定检出水平 α , 在表 C.1 中查出临界值 $D_{1-\alpha}(n)$ 。
- c) 检验高端值, 当 $D_n > D_{1-\alpha}(n)$ 时, 判定 $x_{(n)}$ 为离群值; 检验低端值, 当 $D'_n > D_{1-\alpha}(n)$ 时, 判定 $x_{(1)}$ 为离群值; 否则判未发现离群值。
- d) 对于检出的离群值 $x_{(1)}$ 或 $x_{(n)}$, 确定剔除水平 α^* , 在表 C.1 中查出临界值 $D_{1-\alpha^*}(n)$ 。检验高端值, 当 $D_n > D_{1-\alpha^*}(n)$ 时, 判定 $x_{(n)}$ 为统计离群值, 否则判未发现 $x_{(n)}$ 是统计离群值(即 $x_{(n)}$ 为歧离值); 检验低端值, 当 $D'_n > D_{1-\alpha^*}(n)$ 时, 判定 $x_{(1)}$ 为统计离群值, 否则判未发现 $x_{(1)}$ 是统计离群值(即 $x_{(1)}$ 为歧离值)。

C.2 双侧情形

- a) 计算出统计量 D_n 与 D'_n 的值, 这里 D_n 与 D'_n 由 C.1 的 a) 给出。
- b) 确定检出水平 α , 在表 C.2 查出临界值 $\tilde{D}_{1-\alpha}(n)$ 。
- c) 当 $D_n > D'_n$, $D_n > \tilde{D}_{1-\alpha}(n)$ 时, 判定 $x_{(n)}$ 为离群值; 当 $D'_n > D_n$, $D'_n > \tilde{D}_{1-\alpha}(n)$ 时, 判定 $x_{(1)}$ 为离群值; 否则判未发现离群值。
- d) 对于检出的离群值 $x_{(1)}$ 或 $x_{(n)}$, 确定剔除水平 α^* , 在表 C.2 中查出临界值 $\tilde{D}_{1-\alpha^*}(n)$ 。当 $D_n > D'_n$ 且 $D_n > \tilde{D}_{1-\alpha^*}(n)$ 时, 判定 $x_{(n)}$ 为统计离群值, 否则判未发现 $x_{(n)}$ 是统计离群值(即 $x_{(n)}$ 为歧离值); 当 $D'_n > D_n$ 且 $D'_n > \tilde{D}_{1-\alpha^*}(n)$ 时, 判定 $x_{(1)}$ 为统计离群值, 否则判未发现 $x_{(1)}$ 是统计离群值(即 $x_{(1)}$ 为歧离值)。

表 C.1 单侧狄克逊(Dixon)检验的临界值表

31		0.327	0.371	0.450	0.478
32		0.323	0.367	0.445	0.473
33		0.319	0.362	0.441	0.468
34		0.315	0.358	0.436	0.463
35		0.311	0.354	0.432	0.458
36		0.308	0.350	0.427	0.454
37		0.305	0.347	0.423	0.450
38		0.301	0.343	0.419	0.446
39		0.298	0.340	0.416	0.442
40		0.296	0.337	0.413	0.439
41		0.293	0.334	0.409	0.435
42		0.290	0.331	0.406	0.432
43		0.288	0.328	0.403	0.429
44		0.285	0.326	0.400	0.425
45		0.283	0.323	0.397	0.423
46		0.281	0.321	0.394	0.420
47		0.279	0.318	0.391	0.417
48		0.277	0.316	0.389	0.414
49	$r_{22} = \frac{x_{(n)} - x_{(n-2)}}{x_{(n)} - x_{(3)}} \text{ 或 } r'_{22} = \frac{x_{(3)} - x_{(1)}}{x_{(n-2)} - x_{(1)}}$	0.275	0.314	0.386	0.412
50		0.273	0.312	0.384	0.409
51		0.271	0.310	0.382	0.407
52		0.269	0.308	0.379	0.405
53		0.267	0.306	0.377	0.402
54		0.265	0.304	0.375	0.400
55		0.264	0.302	0.373	0.398
56		0.262	0.300	0.371	0.396
57		0.261	0.298	0.369	0.394
58		0.259	0.297	0.367	0.392
59		0.258	0.295	0.366	0.391
60		0.256	0.294	0.363	0.388
61		0.255	0.292	0.362	0.387
62		0.253	0.291	0.361	0.385
63		0.252	0.289	0.359	0.383
64		0.251	0.288	0.357	0.382
65		0.250	0.287	0.355	0.380
66		0.249	0.285	0.354	0.379

表 C.1 (续)

67		0.247	0.284	0.353	0.377
68		0.246	0.283	0.351	0.376
69		0.245	0.282	0.350	0.374
70		0.244	0.280	0.348	0.372
71		0.243	0.279	0.347	0.371
72		0.242	0.278	0.346	0.370
73		0.241	0.277	0.344	0.368
74		0.240	0.276	0.343	0.368
75		0.239	0.275	0.342	0.366
76		0.238	0.274	0.341	0.365
77		0.237	0.273	0.340	0.364
78		0.236	0.272	0.338	0.363
79		0.235	0.271	0.337	0.361
80		0.234	0.270	0.336	0.360
81		0.233	0.269	0.335	0.359
82		0.232	0.268	0.334	0.358
83		0.232	0.267	0.333	0.356
84	$r_{22} = \frac{x_{(n)} - x_{(n-2)}}{x_{(n)} - x_{(n-1)}} \text{ 或 } r_{22} = \frac{x_{(3)} - x_{(1)}}{x_{(n-2)} - x_{(1)}}$	0.231	0.266	0.332	0.356
85		0.230	0.265	0.331	0.355
86		0.229	0.264	0.330	0.353
87		0.228	0.263	0.329	0.352
88		0.228	0.262	0.328	0.352
89		0.227	0.262	0.327	0.351
90		0.226	0.261	0.326	0.350
91		0.225	0.260	0.325	0.349
92		0.225	0.259	0.324	0.348
93		0.224	0.259	0.323	0.347
94		0.223	0.258	0.323	0.346
95		0.223	0.257	0.322	0.345
96		0.222	0.256	0.321	0.344
97		0.221	0.255	0.320	0.344
98		0.221	0.255	0.320	0.343
99		0.220	0.254	0.319	0.341
100		0.219	0.254	0.318	0.341

表 C.2 双侧狄克逊(Dixon)检验的临界值表

n	统计量	0.95	0.99	n	统计量	0.95	0.99
31		0.403	0.473	66		0.316	0.377
32		0.399	0.468	67		0.315	0.375
33		0.395	0.463	68		0.313	0.376
34		0.39	0.46	69		0.313	0.375
35		0.388	0.458	70		0.312	0.375
36		0.438	0.442	71		0.31	0.373
37		0.38	0.45	72		0.309	0.373
38		0.377	0.447	73		0.308	0.371
39		0.375	0.442	74		0.306	0.37
40		0.37	0.438	75		0.305	0.368
41		0.367	0.433	76		0.304	0.363
42		0.364	0.432	77		0.304	0.363
43		0.362	0.428	78		0.303	0.362
44		0.359	0.425	79		0.303	0.361
45		0.357	0.422	80		0.302	0.358
46		0.353	0.419	81		0.301	0.358
47		0.352	0.416	82		0.301	0.355
48	r_{22} 和 r'_{22} 中较大者	0.35	0.413	83	r_{22} 和 r'_{22} 中较大者	0.301	0.355
49		0.346	0.412	84		0.298	0.353
50		0.343	0.409	85		0.297	0.351
51		0.342	0.407	86		0.297	0.351
52		0.34	0.405	87		0.296	0.349
53		0.338	0.402	88		0.295	0.349
54		0.337	0.4	89		0.294	0.347
55		0.335	0.399	90		0.293	0.347
56		0.334	0.399	91		0.291	0.344
57		0.33	0.396	92		0.29	0.344
58		0.329	0.393	93		0.289	0.343
59		0.327	0.39	94		0.289	0.343
60		0.325	0.389	95		0.288	0.343
61		0.323	0.387	96		0.288	0.342
62		0.321	0.385	97		0.286	0.34
63		0.32	0.383	98		0.285	0.34
64		0.319	0.382	99		0.285	0.339
65		0.318	0.379	100		0.284	0.339

参 考 文 献

- [1] W. J. Dixon. Processing data for outliers. *Biometrics*, 1953, 9(1), 74~89. University of Oregon
- [2] W. J. Dixon. Analysis of extreme values; *Annals of Mathematical Statistics*. 1950, 21(4) . 488~506
- [3] W. J. Dixon. Ratios involving extreme values; *Annals of Mathematical Statistics*, 1951, 22(1). 68~78
- [4] Surendra P. Verma and Alfredo Quiroz-Ruiz. Critical values for six Dixon tests for outliers in normal samples up to sizes 100, and applications in science and engineering. *Revista Mexicana de Ciencias Geológicas*. 2006, 23(2). 133~161
- [5] C. E. Efstathiou. Estimation of Type I Error Probability from Experimental Dixon's "Q" Parameter on Testing for Outliers within small Size Data Sets. *Talanta*, 2006. 69(5)
- [6] V. Bartlett and T. Lewis. *Outliers in Statistical data*. Chichester. John Wiley. Third edition, 584
-