

数据分析之 – 离群值（**Outliers**）

一：什么是Outliers

Outliers是统计学专业术语，是指相比一组数据中的其它数据的极限值

二：极限值意味什么

1. 决定哪些值是**Outliers**是一个主观行为，有一些基准数据来决定是否一个值是一个**Outliers**，这些基准是任意选择的，比如 $P \leq 0.5$ 就是一个任意选择的基准

2. 一个基准是用**BoxPlot**来决定适度离群值（**mild Outliers**）和极限离群值（**extreme Outliers**），适度离群值是任何值1.5倍大于基于剩下所有的值的**IQR**，极限离群值是任何值3倍大于剩下所有的值的**IQR**，**IQR**（**Interquartile Range**）代表四分位数间距，是这些值中的50%中间值，分别是Q1-25%, Median-50%, Q3-75%, $IQR = Q3 - Q1$

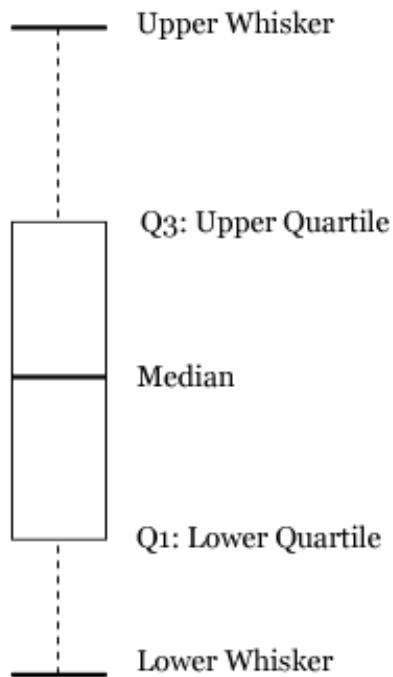
三：使用Box Plot来发现Outliers

一个典型的**Box Plot**是基于以下五个值计算而来的

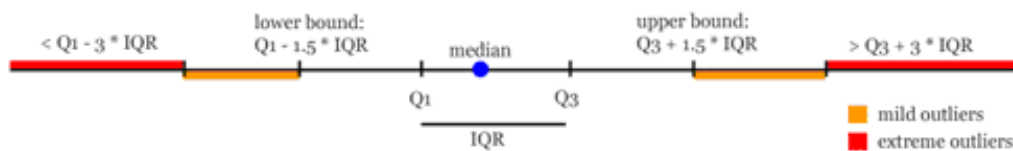
- a. 一组样本的最小值
- b. 一组样本的最大值
- c. 一组样本的中值
- d. 下四分位数（**Lower Quartile / Q1**）
- e. 上四分位数（**Upper Quartile / Q3**）

根据这五个值构建出来基本的**Box Plot**，某些图形软件还会显示平均值， $IQR = Q3 - Q1$

显然超出上下四分位数的值可以看做为**Outliers**。我们通过眼睛就可以很好的观察到这些**Outliers**值的点。



一个显示适度 and 极限Outliers值的Box plot显示如下：



四：示例说明及JfreeChart的实现

假设一组数据为：2,4,6,8,12,14,16,18,20,25,45

中值 Median = 14

Q1-下四分位数 ($11 * 0.25 = 3$) = 7

Q3-上四分位数 ($11 * 0.75 = 9$) = 19

IQR ($Q3 - Q1$) = 12

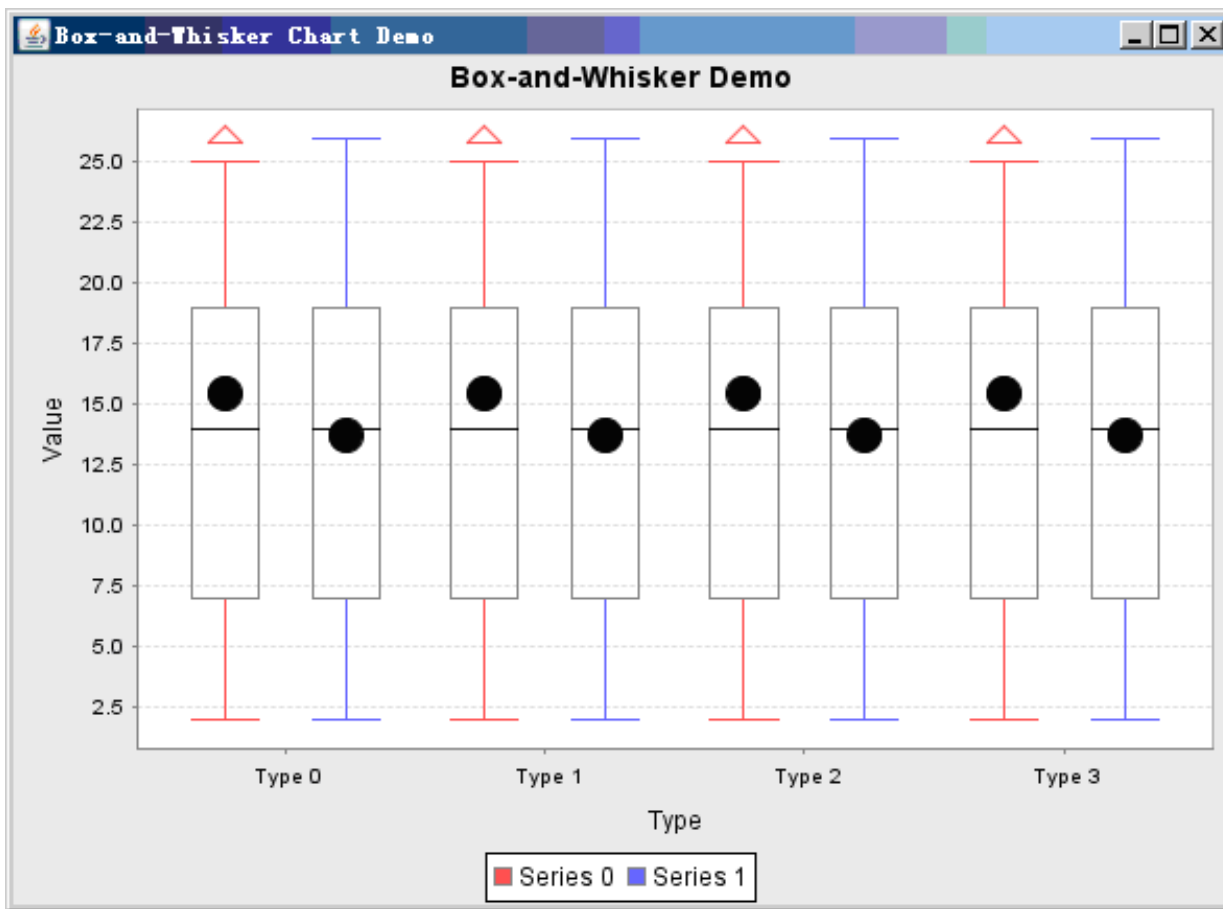
$1.5 * IQR = 18$

最小值 ($6 - 1.5 * IQR$) = 2

最大值 ($20 + 1.5 * IQR$) = 25

很显然值45是一个适度Outliers

对比的一组数据为：2,4,6,8,12,14,16,18,20,25,26



从图上可以看出Series0的数据存在Outliers，一个红色三角形已经表明

同样Series1的数据是一组非常好的数据，没有Outliers.

下面是Java源代码：

```
[java]
01. package com.dataanalysis.plots;
02. import java.awt.Font;
03. import java.util.ArrayList;
04. import java.util.List;
05. import org.jfree.chart.ChartPanel;
06. import org.jfree.chart.JFreeChart;
07. import org.jfree.chart.axis.CategoryAxis;
08. import org.jfree.chart.axis.NumberAxis;
09. import org.jfree.chart.labels.BoxAndWhiskerToolTipGenerator;
```

```

10. import org.jfree.chart.plot.CategoryPlot;
11. import org.jfree.chart.renderer.category.BoxAndWhiskerRenderer;
12. import org.jfree.data.statistics.BoxAndWhiskerCategoryDataset;
13. import org.jfree.data.statistics.DefaultBoxAndWhiskerCategoryDataset;
14. import org.jfree.ui.ApplicationFrame;
15. import org.jfree.ui.RefineryUtilities;
16. public class BoxAndWhiskerDemo extends ApplicationFrame {
17.     /**
18.      *
19.      */
20.     private static final long serialVersionUID = -3205574763811416266L;
21.     /**
22.      * Creates a new demo.
23.      *
24.      * @param title the frame title.
25.      */
26.     public BoxAndWhiskerDemo(final String title) {
27.         super(title);
28.
29.         final BoxAndWhiskerCategoryDataset dataset = createSampleDataset();
30.         final CategoryAxis xAxis = new CategoryAxis("Type");
31.         final NumberAxis yAxis = new NumberAxis("Value");
32.         yAxis.setAutoRangeIncludesZero(false);
33.         final BoxAndWhiskerRenderer renderer = new BoxAndWhiskerRenderer();
34.         renderer.setFillBox(false);
35.         renderer.setToolTipGenerator(new BoxAndWhiskerToolTipGenerator());
36.         final CategoryPlot plot = new CategoryPlot(dataset, xAxis, yAxis, renderer);
37.         final JFreeChart chart = new JFreeChart(
38.             "Box-and-Whisker Demo",
39.             new Font("SansSerif", Font.BOLD, 14),
40.             plot,
41.             true
42.         );
43.         final ChartPanel chartPanel = new ChartPanel(chart);
44.         chartPanel.setPreferredSize(new java.awt.Dimension(450, 270));
45.         setContentPane(chartPanel);
46.     }
47.     /**
48.      * Creates a sample dataset.
49.      *
50.      * @return A sample dataset.
51.      */
52.     private BoxAndWhiskerCategoryDataset createSampleDataset() {
53.
54.         final int seriesCount = 2;
55.         final int categoryCount = 4;
56.         double[] data = null;
57.         final DefaultBoxAndWhiskerCategoryDataset dataset
58.             = new DefaultBoxAndWhiskerCategoryDataset();
59.         for (int i = 0; i < seriesCount; i++) {
60.             if(i == 0) {
61.                 data = new double[]{2,4,6,8,12,14,16,18,20,25,45};

```

```
62.         } else {
63.             data = new double[]{2,4,6,8,12,14,16,18,20,25,26};
64.         }
65.
66.         for (int j = 0; j < categoryCount; j++) {
67.             final List list = new ArrayList();
68.             for (int k = 0; k < data.length; k++) {
69.                 list.add(new Double(data[k]));
70.             }
71.             dataset.add(list, "Series " + i, " Type " + j);
72.         }
73.
74.     }
75.     return dataset;
76. }
77. /**
78.  * For testing from the command line.
79.  *
80.  * @param args ignored.
81.  */
82. public static void main(final String[] args) {
83.     final BoxAndWhiskerDemo demo = new BoxAndWhiskerDemo("Box-and-Whisker Chart Demo");
84.     demo.pack();
85.     RefineryUtilities.centerFrameOnScreen(demo);
86.     demo.setVisible(true);
87. }
88. }
```