

The use of log-ratio methodology in cell-wise outlier diagnostics

Jan Walach¹, P.Filzmoser¹, K.Hron²

¹TU Wien, ²Palacky University Olomouc



TECHNISCHE
UNIVERSITÄT
WIEN
Vienna University of Technology

June 11, 2018

Outline

The screenshot displays the TeamViewer 13 interface with an R script editor and a console window. The script in the editor is as follows:

```
1 # Outline of DSSV 2018 presentation -----
2 setwd("/home/jan/Dropbox/Git/Batch2/Prezentace")
3 set.seed(1314233448)
4 for (i in 1:7)
5 {
6   p <- sample(c('Talking', 'Talking', 'Describing a plot', 'Getting
7   print(p)
8   for (j in 1:length(dir(paste0(getwd(), '/', p))))
9   {
10    temp = list.files(path = paste0(getwd(), '/', p), pattern="*")
11    mySlides = lapply(temp, read.delim)
12  }
13 }
```

The console window shows the output of the script's execution:

```
[1] "Talking"
[1] "Getting lost in a complicated formulas"
[1] "Talking"
[1] "Talking"
[1] "Getting lost in a complicated formulas"
[1] "Describing a plot"
[1] "Talking"
>
```

The right-hand side of the interface shows the Environment and History panels. The Environment panel lists variables: p2, p3, p4, p5, pp, and s. The History panel shows a list of operations. Below these panels is a file explorer showing the current directory structure, including folders like 'Prezentace' and 'res_OPLSDA.Rdata', and files like 'WeightMatrixCubeW.pdf' and 'SampleSizeResults.pdf'.

Data

- Focus on metabolomics data
- Two groups (e.g. control/disease patients)
- Usually low number of observations
- Flat data: many variables

Goal

- Clustering/Classification
- Exploratory analysis
- Identification of non-standard samples/values
- Find informative variables

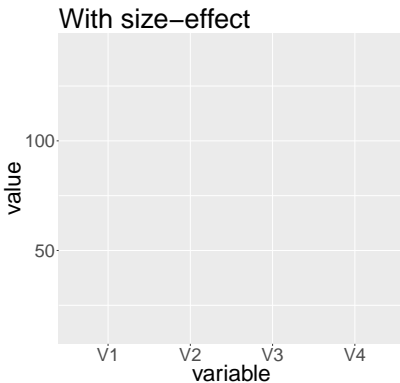
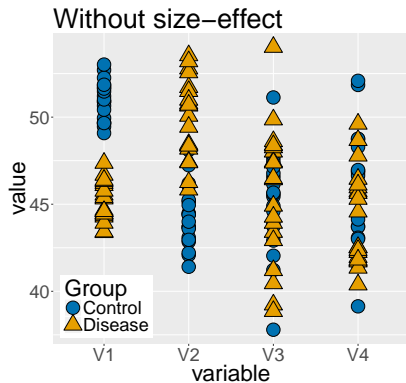
The use of log-ratio methodology in cell-wise outlier diagnostics

The use of log-ratio methodology in cell-wise outlier diagnostics

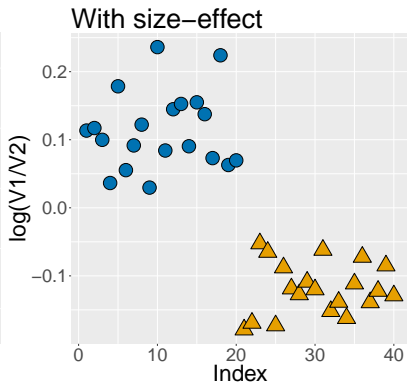
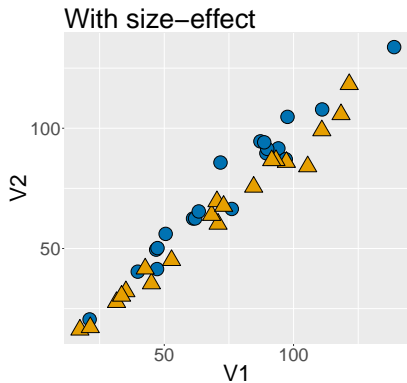
Pawlowsky-Glahn, V., Egozcue, R. Tolosana-Delgado, J.J.:
Modeling and Analysis of Compositional Data. Chichester: Wiley,
2015.

Size-effect

- Problem: “Size effect”
- Different sample volume and/or sample concentration
- Special treatment required



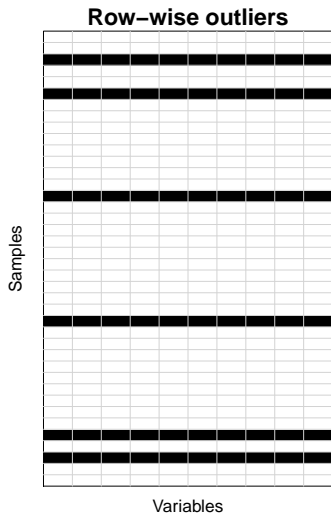
Size-effect: Use of log-ratios



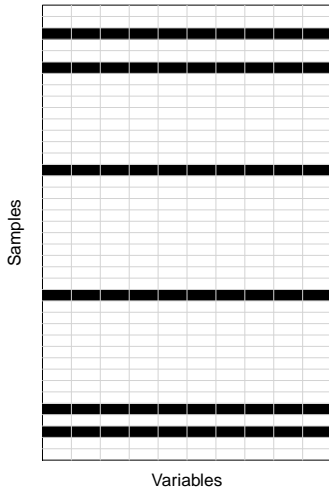
The use of log-ratio methodology in **cell-wise outlier diagnostics**

The use of log-ratio methodology in **cell-wise outlier diagnostics**

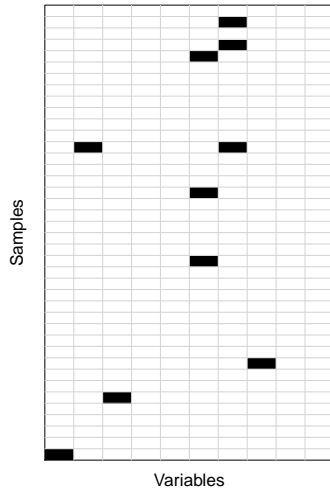
Peter J. Rousseeuw, Wannes Van Den Bossche. "Detecting deviating data cells." Technometrics 60.2 (2018): 135-145.



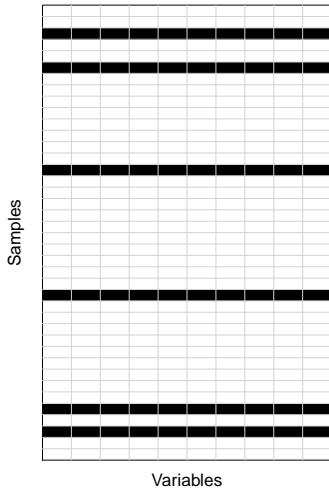
Row-wise outliers



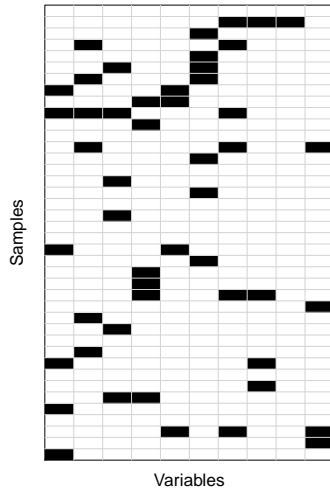
Cell-wise outliers



Row-wise outliers



Cell-wise outliers



Method:

Cell-wise rPLR

Steps:

- 1 Computing log-ratios
- 2 Robust centering + scaling
- 3 Applying weighting function
- 4 Projection to original space

Data

- Data matrix X with dimensions $n \times D$

Based on idea of robust variation matrix

- Matrix T , elements t_{jk} , for $j, k = 1, \dots, D$

$$t_{jk} = \text{var} \left[\ln \left(\frac{x_{1j}}{x_{1k}} \right), \ln \left(\frac{x_{2j}}{x_{2k}} \right), \dots, \ln \left(\frac{x_{nj}}{x_{nk}} \right) \right], \quad (1)$$

- where var denotes a variance
- Elements inside $\text{var}()$ will be robustly centred + scaled and each will get a weight
- There are in total $\frac{D \times (D-1)}{2}$ possible log-ratios

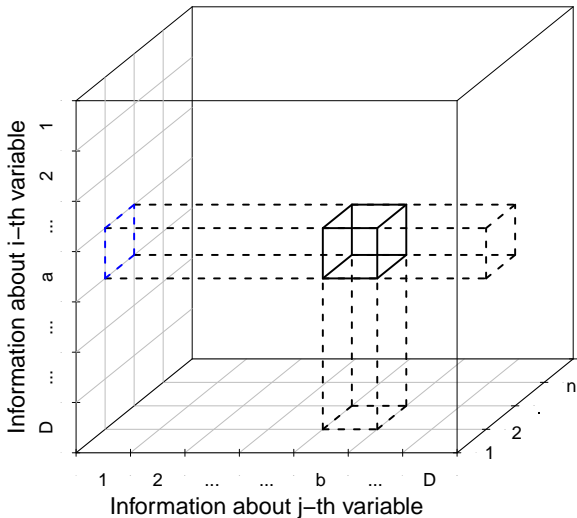
$$t_{jk} = \text{var} \left[\ln \left(\frac{x_{1j}}{x_{1k}} \right), \ln \left(\frac{x_{2j}}{x_{2k}} \right), \dots, \ln \left(\frac{x_{nj}}{x_{nk}} \right) \right],$$

Three-way weight matrix

- All weights can be stored in the three-way matrices W , with D rows, D columns and n slices
- Information about observations is in certain slices
- Information about cells in rows of slices

Projection to original dimensions

Weight matrix



Diagnostics

- Average of all weights for each observation and each involved variable,

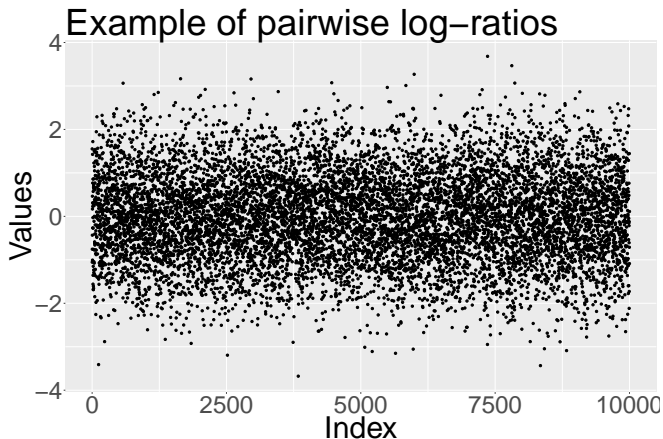
$$m_{ij} = \frac{1}{D} \sum_{k=1}^D w_{jki} \quad (2)$$

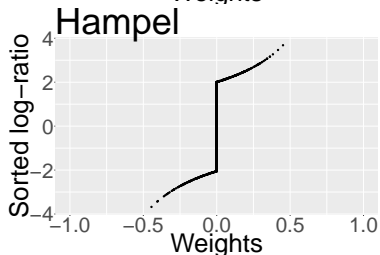
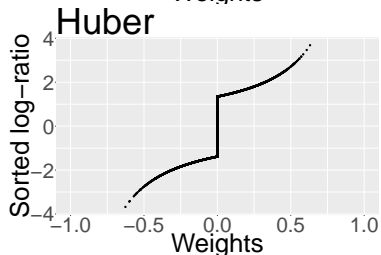
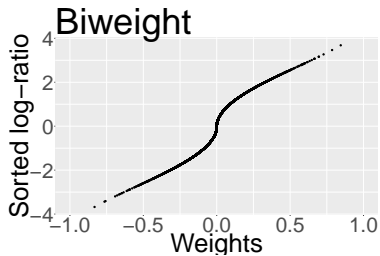
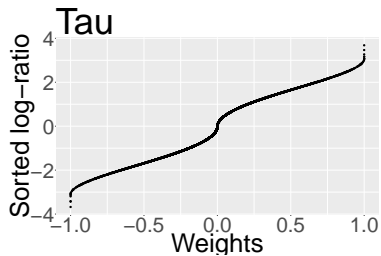
- Or Median

$$m_{ij} = \operatorname{median}_k (w_{jki}) \quad (3)$$

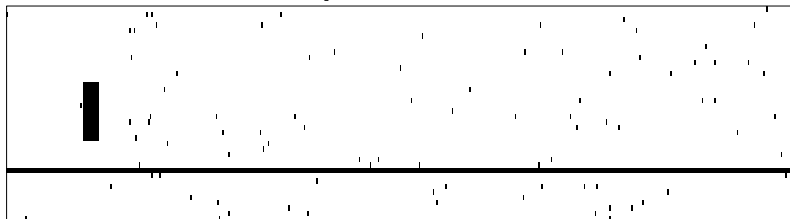
Weighting functions:

- τ function
- Tukey's biweight function
- Huber function
- Hampel function

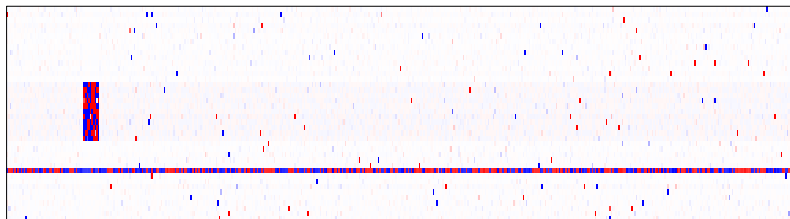




Imputed outliers

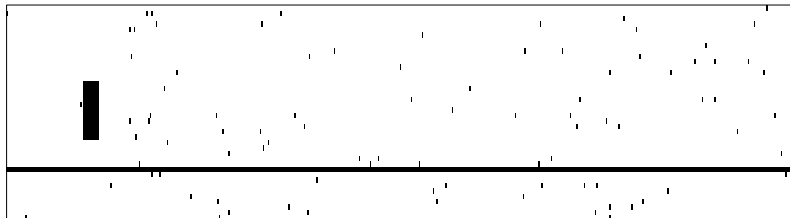


Identified outliers

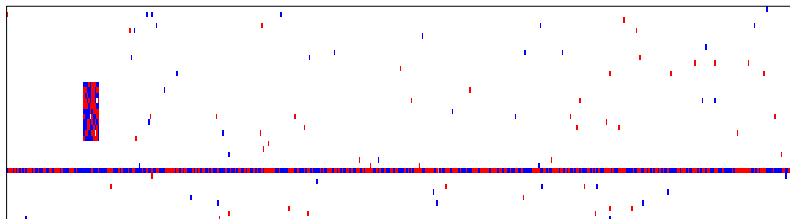


Diagnostics example: Hampel

Imputed outliers



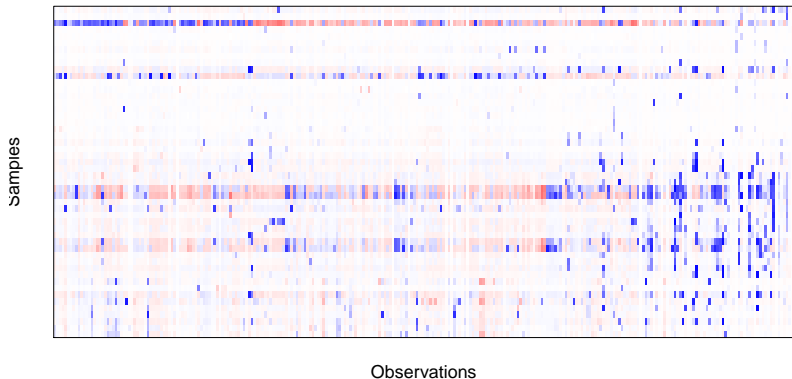
Identified outliers



Metabolomic dataset MCADD

- MCADD (Medium chain acyl-CoA dehydrogenase deficiency)
- $n_1 = 25$, $n_2 = 25$, $D = 273$

MCAD diagnostics

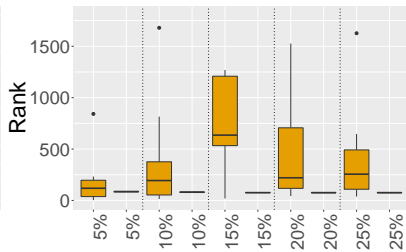
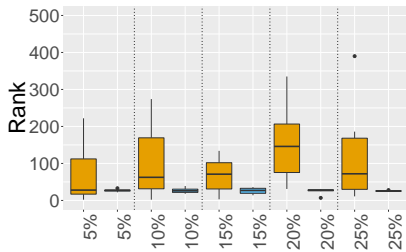
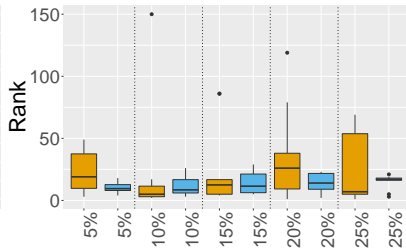
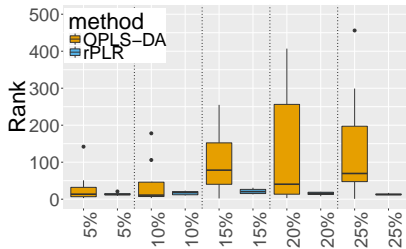


- Let's consider two group problem with n_1 , n_2 samples
- For each variable j :






$$V_j = |\text{median}_{i=1,\dots,n_1} m_{ij} - \text{median}_{i=1,\dots,n_2} m_{ij}| \quad (4)$$

- Permutation test can be used in order to set a cut-off values

Metabolomic dataset



References

-  Pawlowsky-Glahn, V., Egozcue, R. Tolosana-Delgado, J.J.: *Modeling and Analysis of Compositional Data*. Chichester: Wiley, 2015.
-  Peter J. Rousseeuw, Wannes Van Den Bossche. "Detecting deviating data cells." *Technometrics* 60.2 (2018): 135-145.
-  Walach, J., Filzmoser, P., Hron, K., Walczak, B., Najdekr, L.: Robust biomarker identification in a two-class problem based on pairwise log-ratios. *Chemometrics and Intelligent Laboratory Systems*, 171 (2017), 277-285
-  P. Filzmoser, B. Walczak, What can go wrong at the data normalization step for identification of biomarkers?, *Journal of Chromatography A* 1362 (2014) 194–205.
-  https://github.com/walachja/Cell-wise_rPLR