

Outlier Detection with Application to Geochemistry

Peter Filzmoser

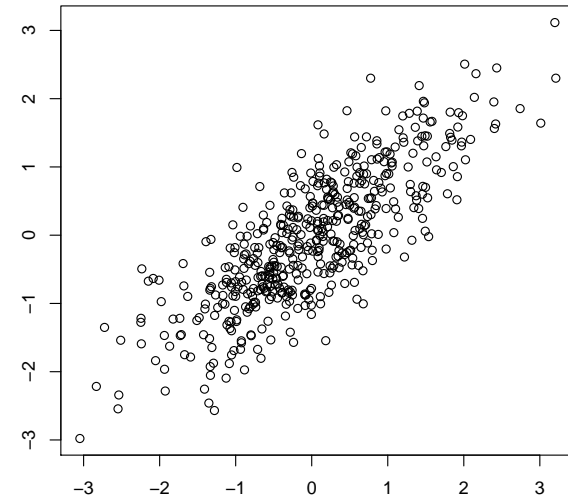
Department of Statistics and Probability Theory
Vienna University of Technology, Austria

Vienna, Austria

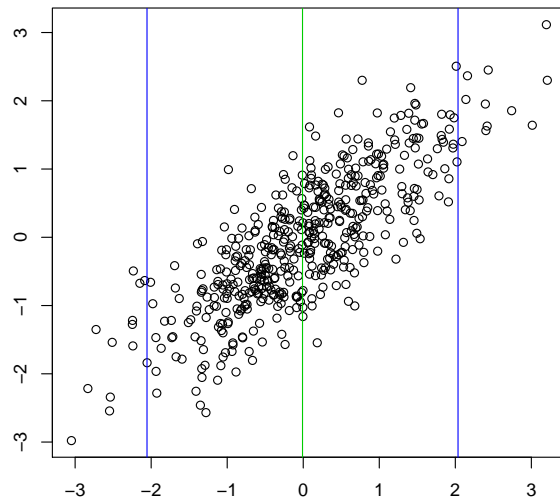
June 16, 2006



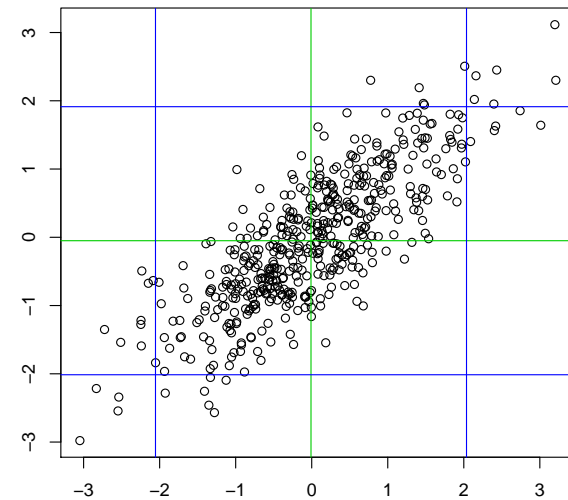
Vienna University of Technology

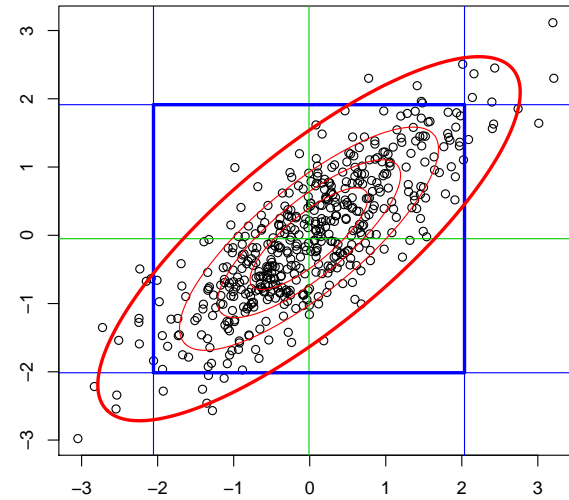
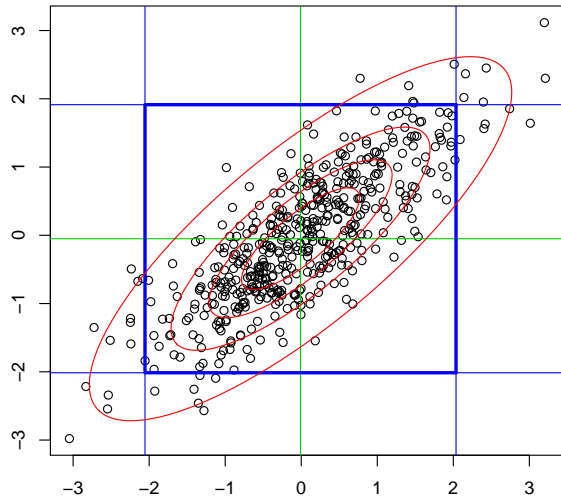
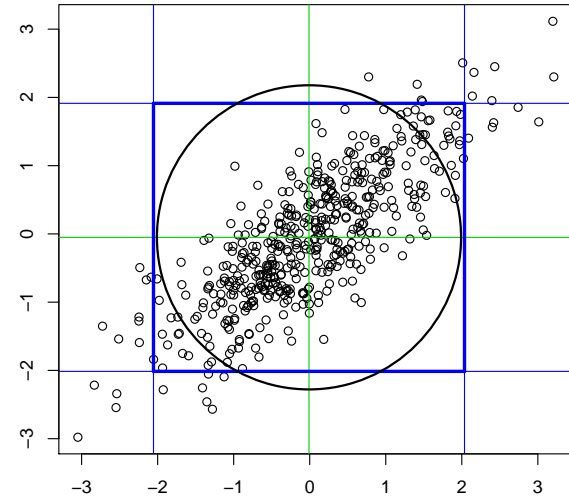
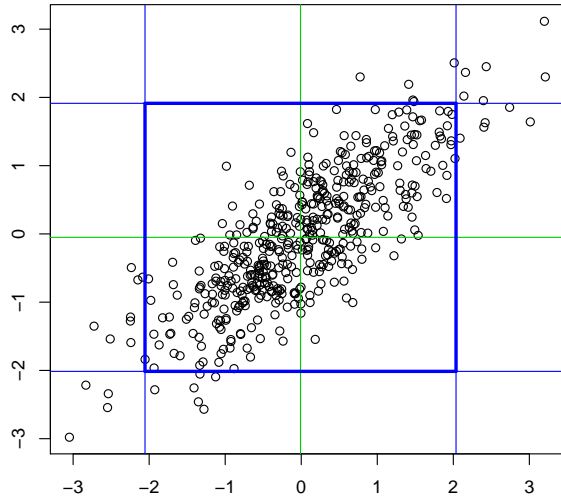


Univariate versus Multivariate Outliers



Univariate versus Multivariate Outliers





Standard methods are based on the **Mahalanobis distances (MD)**:

$$MD_i := d(\mathbf{x}_i, \mathbf{t}, \mathbf{C}) = \{(\mathbf{x}_i - \mathbf{t})^\top \mathbf{C}^{-1} (\mathbf{x}_i - \mathbf{t})\}^{1/2}$$

for a sample $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ and estimators of location \mathbf{t} and covariance \mathbf{C} .

⇒ **Robust estimates of location and covariance are needed!**

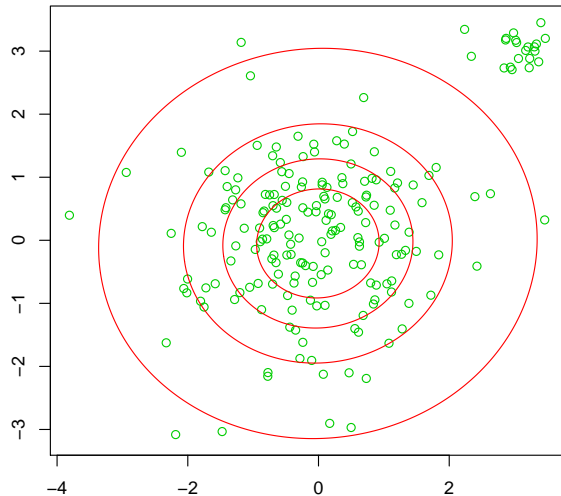
Outlier detection:

Outliers will typically have large distance. If multivariate normal distribution is assumed, MD_i^2 is approx. χ_p^2 distributed.

⇒ suspect observations: $MD_i^2 > \chi_{p,0.975}^2$

- does not account for different sample size
- χ_p^2 -approximation is poor

Example: Simulated data with outliers



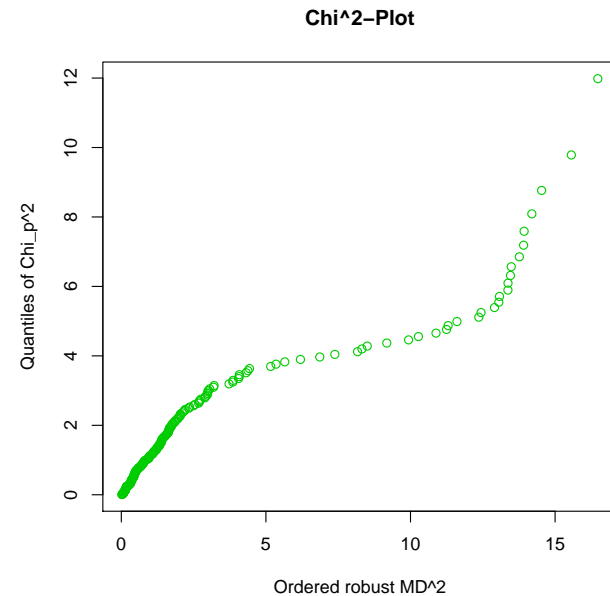
Chi-square plot:

Plot robust MD_i^2 against quantiles of χ_p^2 .

⇒ iterative deletion of points with large distance until a straight line appears.

Drawback: no automatic procedure, needs user interaction.

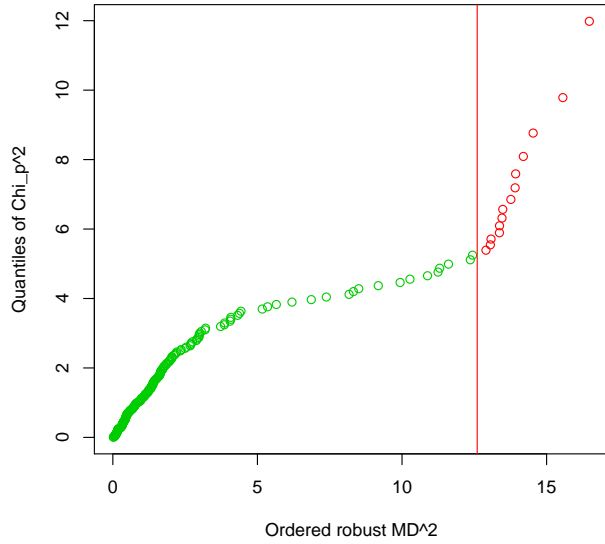
Iterative deletion of outliers:



Iterative deletion of outliers:



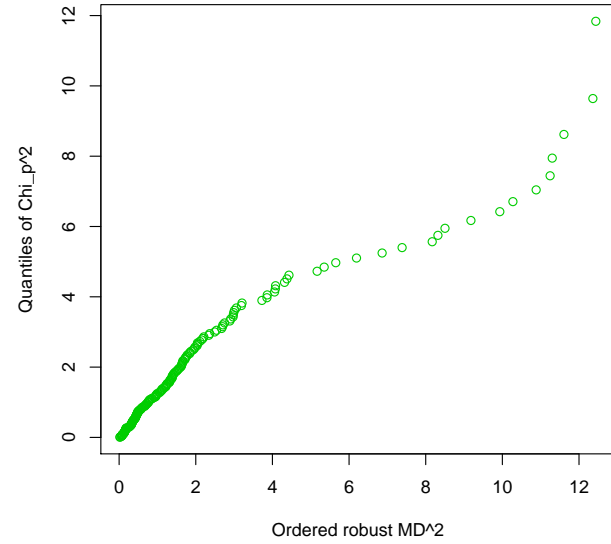
Chi²-Plot



Iterative deletion of outliers:



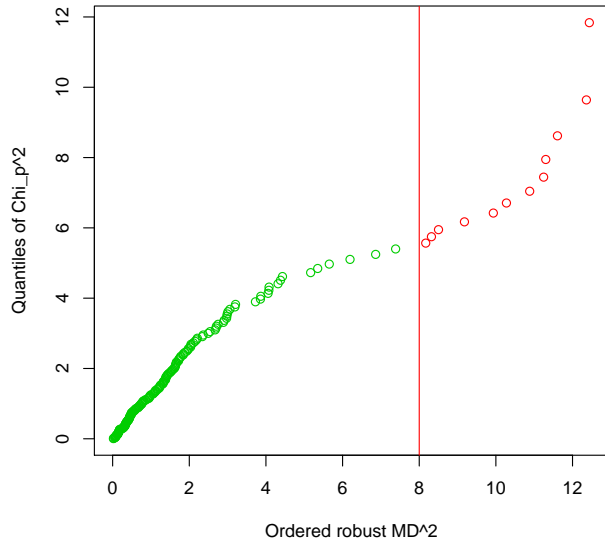
Chi²-Plot



Iterative deletion of outliers:



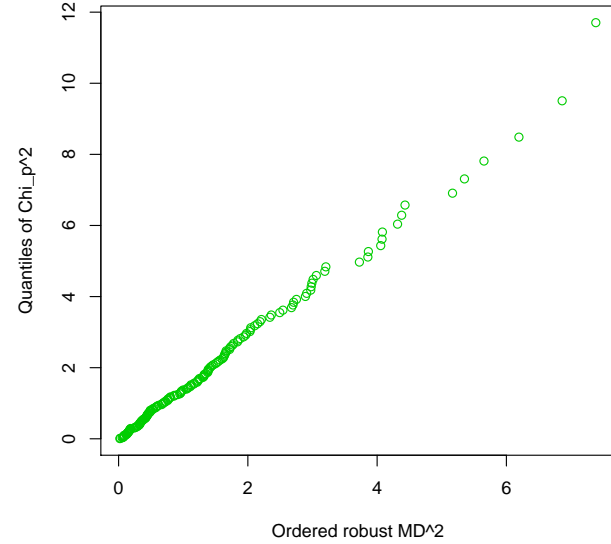
Chi²-Plot

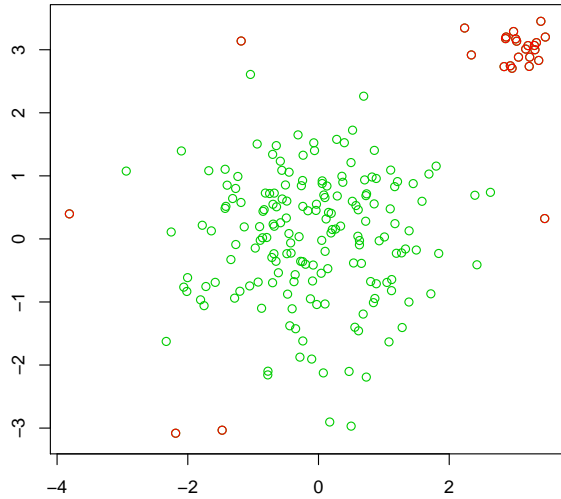


Iterative deletion of outliers:



Chi²-Plot





$G(u)$... theoretical distribution function of χ_p^2 ,
 $G_n(u)$... empirical distribution function of MD_i^2 .

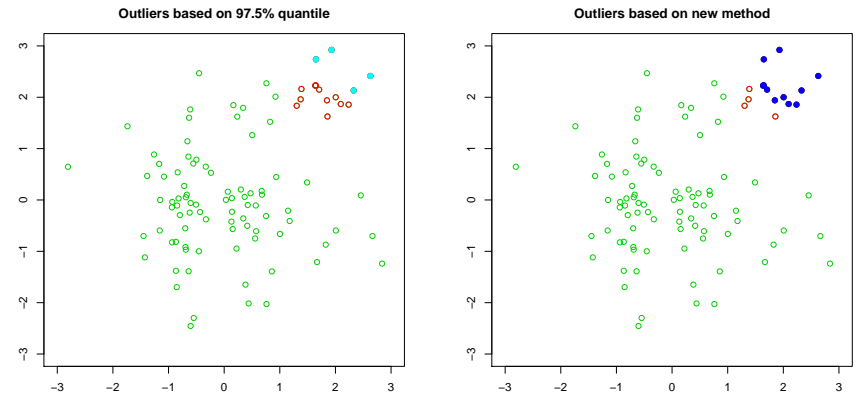
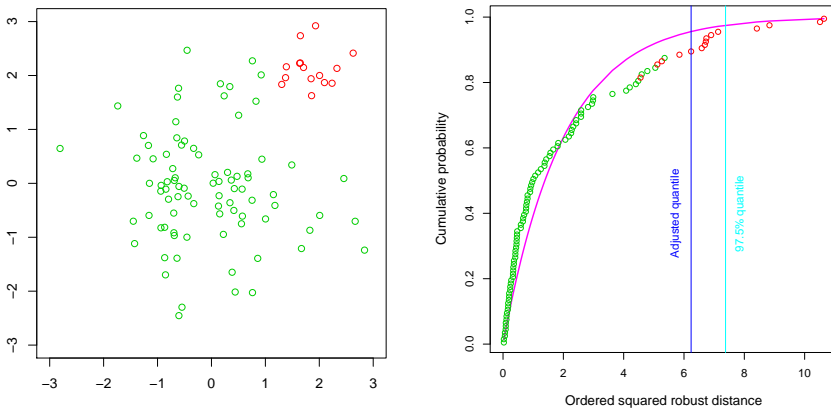
For $\eta = \chi_{p,1-\alpha}^2$ define

$$p_n(\eta) = \sup_{u \geq \eta} \{G(u) - G_n(u)\}^+.$$

Then a measure of outliers in the sample is

$$\alpha_n(\eta) = \begin{cases} 0 & \text{if } p_n(\eta) \leq p_{crit}(\eta, n, p) \\ p_n(\eta) & \text{if } p_n(\eta) > p_{crit}(\eta, n, p). \end{cases}$$

$p_{crit}(\eta, n, p)$ can be obtained by simulations.

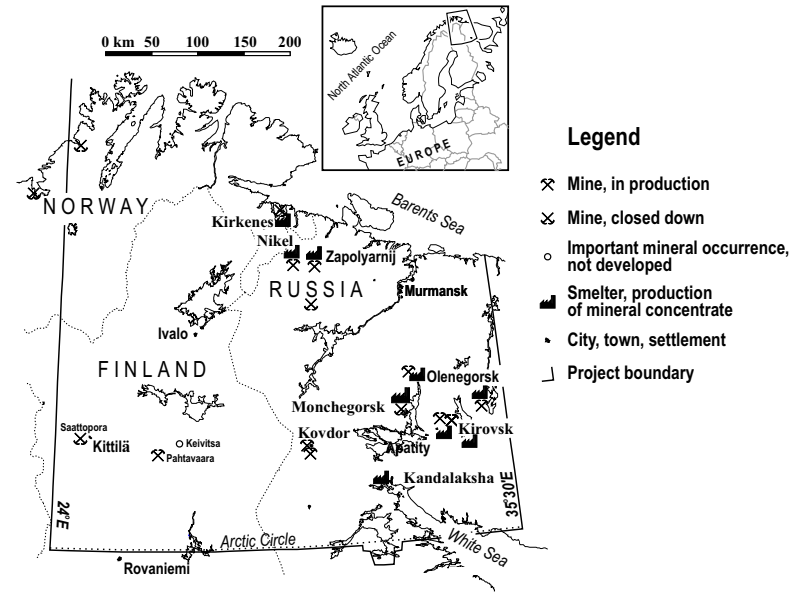


Consider the O-horizon (organic surface soil) of the Kola data set.

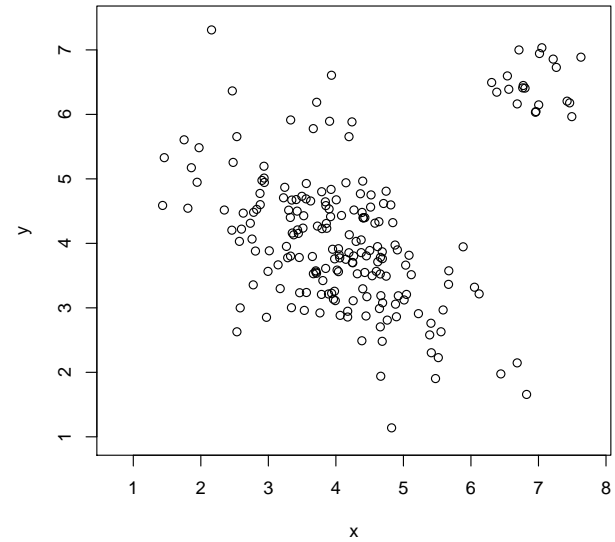
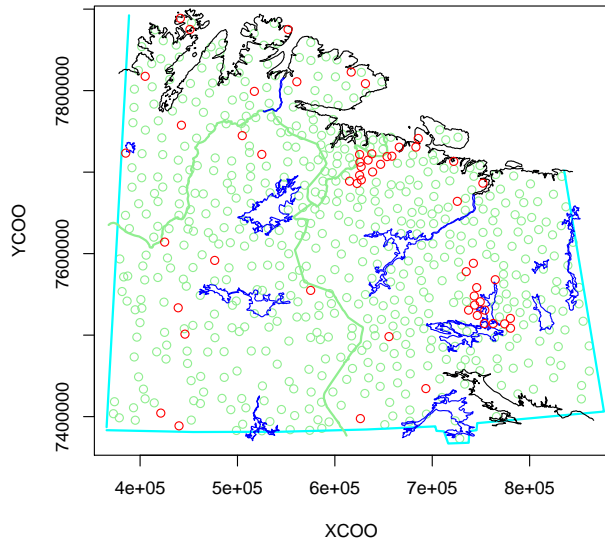
Take (more or less) typical elements for "pollution":

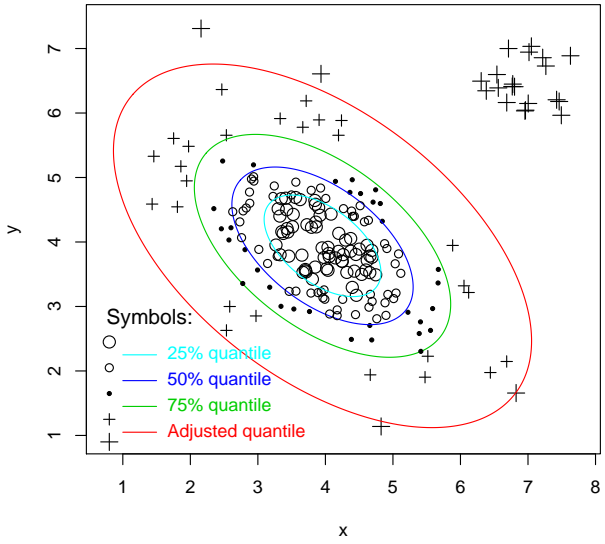
As, Cd, Co, Cu, Mg, Pb, Zn

Question: *Where are the multivariate outliers?*

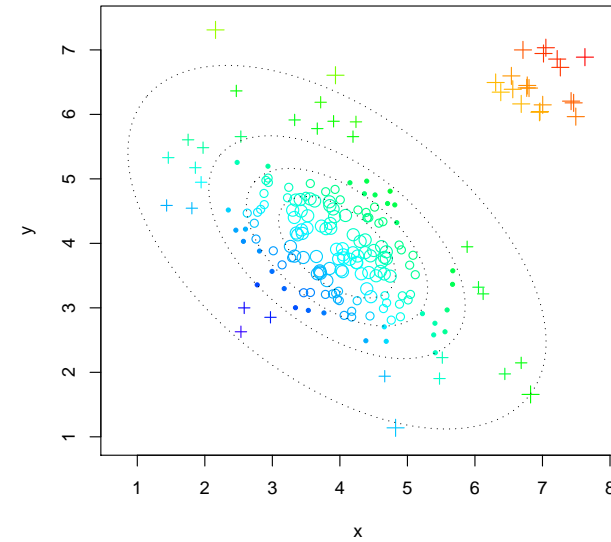
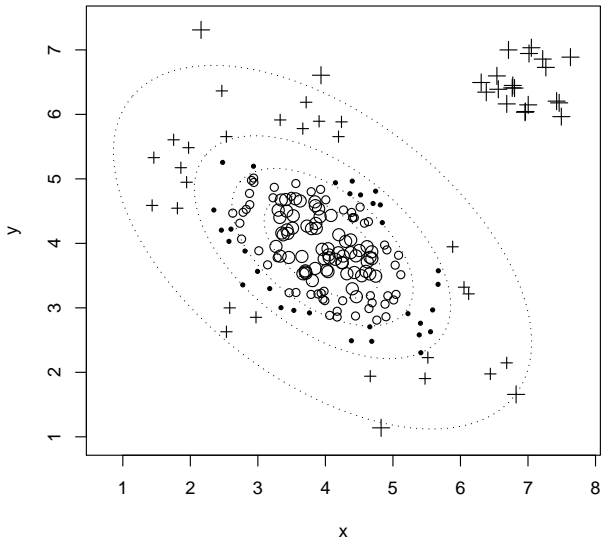
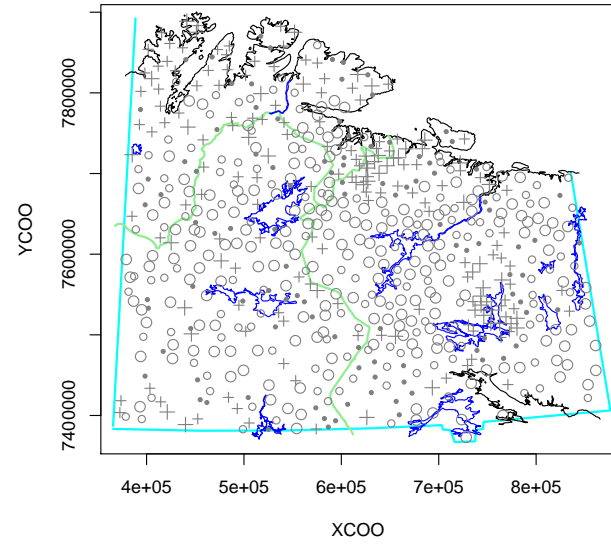


RED points are outliers

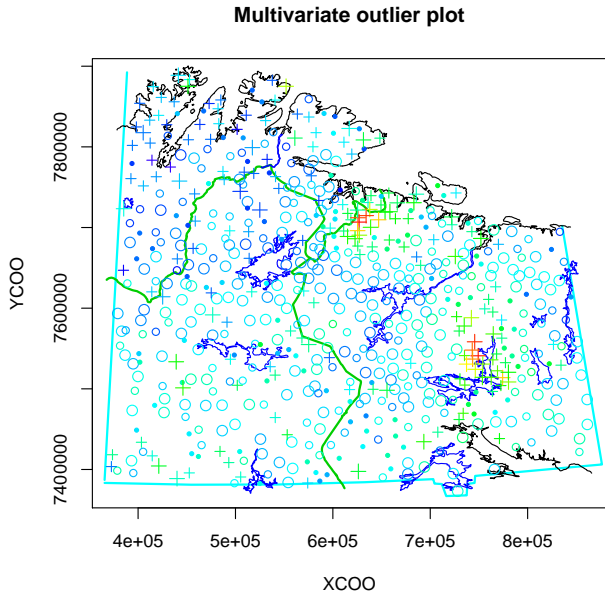




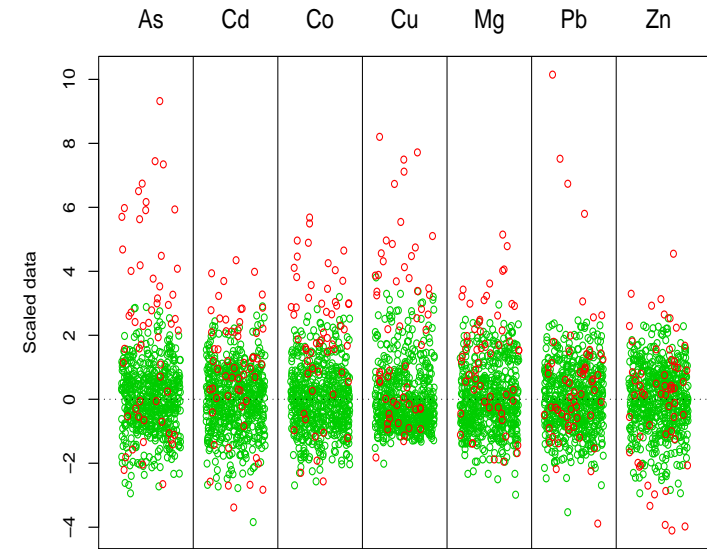
Robust MD with symbols



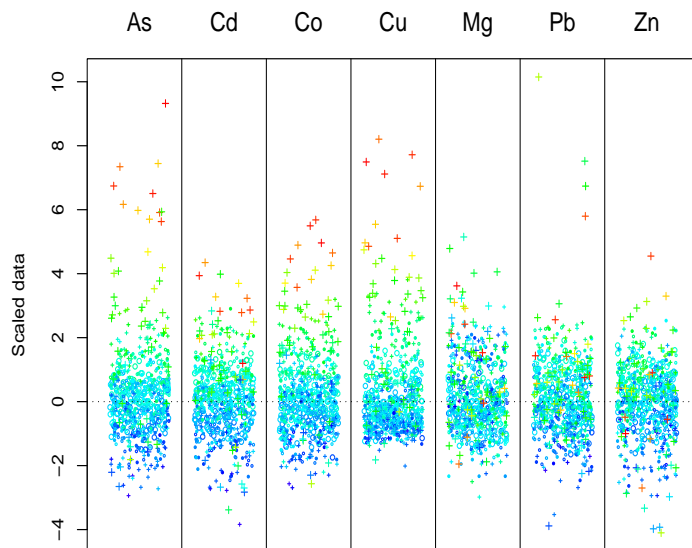
Example: Map showing outliers



Example: From Multivariate to Univariate



Example: Symbols from multivariate plot



Summary

```
library(mvoutlier)
```

includes

- all routines to generate the presented plots
- Kola data and other interesting geochemical data sets