

Robust Statistical Methods for Outlier Detection with Application to Household Expenditure Data

Johannes Gussenbauer ¹ Peter Filzmoser ¹
Matthias Templ ¹ Oliver Dupriez ²

¹Vienna University of Technology

²World Bank

October 27, 2015

Outlier in household expenditure data

- ▶ Household expenditure data usually provided through household surveys
 - ▶ Data subject to human error
 - ▶ participants don't want to share every information
- ▶ The Gini coefficient plays an important role in connection with household expenditure data
 - ▶ Measures the inequality of the household spendings among the surveyed households

Impact of Outliers

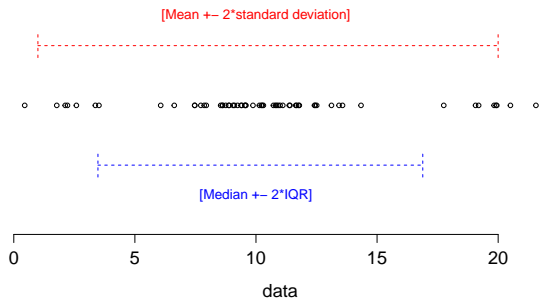
- ▶ Outliers may have a huge impact on non-robust estimators
- ▶ Ranking between countries may completely change only because of outliers
- ▶ From then on World Bank used simple univariate outlier detection and replacement of outliers
- ▶ Projekt with World Bank to improve outlier detection and replacement

Provided data and data structure

- ▶ Household expenditure data from Albania(2008), Mexico(2010), India(2009), Malawi(2010) and Tajikistan(2007)
- ▶ Product of large household surveys containing value of goods or services consumed in local currency for each household over a period of time
- ▶ World Bank started to harmonize the resulting data into a common framework
- ▶ Household consumption categorized by
 - ▶ ICP basic headings / ICP class / ICP group / ICP category

Robust statistical methods

- ▶ Use robust statistical methods to detect potential outliers
- ▶ Due to the structure of the data univariate and multivariate methods were tested



Univariate methods

- ▶ Data points which are "far enough" away from the main bulk of the data
- ▶ The following methods were used:
 - ▶ Estimate location and scale in a robust way to determine interval for "good" observations
 - ▶ $[med - c \cdot S_{IQR}, med + c \cdot S_{IQR}]$
 - ▶ $[med - c \cdot S_{MAD}, med + c \cdot S_{MAD}]$
 - ▶ Boxplot
 - ▶ Expenditure data usually skewed to the right
 - ▶ use Box-Cox transformation \Rightarrow estimate interval \Rightarrow transform back interval boundaries
 - ▶ use skewness-adjusted Boxplot
 - ▶ Pareto tail modeling

Replace univariate potential outliers

- ▶ Place potential outliers onto the lower/upper ends of the calculated intervals
- ▶ For Pareto tail modeling, values larger than a certain quantile of the fitted distribution
 - ▶ are replaced by values drawn from the fitted distribution
 - ▶ their sample weights are set to 1 and recalibrate for the rest of the data

Applying univariate methods to multivariate data

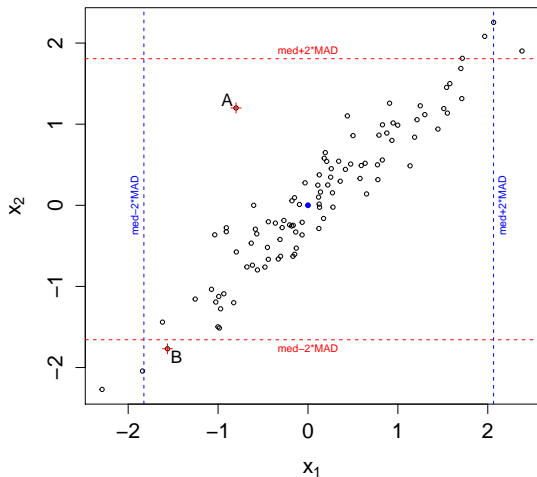


Figure: Simulated data from multivariate standard normal

Mahalanobis distance

- ▶ Use distance measure which takes into account the multidimensional structure of the data \Rightarrow Squared Mahalanobis distance MD_i^2

$$MD_i^2 = (\mathbf{x}_i - \bar{\mathbf{x}})^t S^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}) \quad ,$$

- ▶ Estimate center and covariance in a robust way to gain squared robust distances, RD_i^2
- ▶ If data follows a multivariate normal distribution \Rightarrow
 $MD_i^2 \sim \chi_p^2$
- ▶ Declare data points as potential outliers if they exceed $\chi_{p;0.975}^2$

Multivariate methods

- ▶ Robust methods to estimate center and covariance
 - ▶ M-estimate
 - ▶ Generalization of Maximum Likelihood estimate
 - ▶ S-estimate
 - ▶ MM-estimate
 - ▶ Uses high breakdown preliminary S-estimate
 - ▶ MCD-& MVE-estimate
 - ▶ Minimum covariance determinant estimate
 - ▶ Minimum volume ellipsoid estimate
 - ▶ Stahel-Donoho estimate
 - ▶ Incorporates multivariate measure of outlyingness
 - ▶ OGK estimate
 - ▶
$$\text{Cov}(X, Y) = \frac{1}{4}(\text{Var}(X + Y) - \text{Var}(X - Y))$$

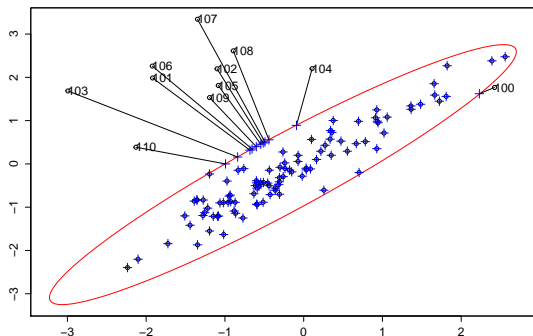
Multivariate methods

- ▶ BACON-EEM
 - ▶ Combines BACON algorithm and EEM algorithm
 - ▶ Uses EEM-algorithm to estimate center and covariance during BACON-procedure
 - ▶ EEM-algorithm able to handle missing values in the data

- ▶ Epidemic Algorithm
 - ▶ Simulate an epidemic, starting from the center of the data
 - ▶ Data points with high infection times are declared potential outliers

Replace potential outliers

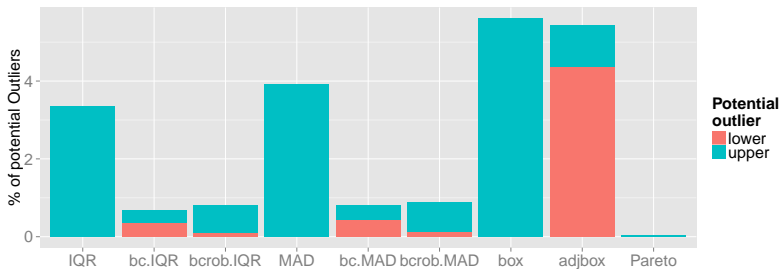
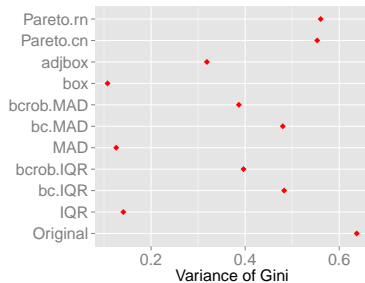
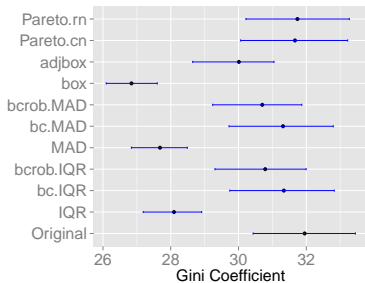
- ▶ Multivariate potential outliers are winsorised onto the boundaries of the 97.5% tolerance ellipse.



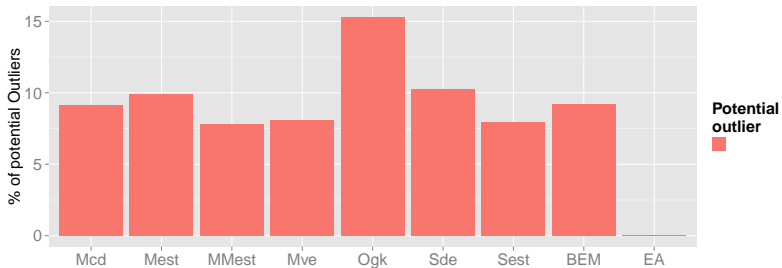
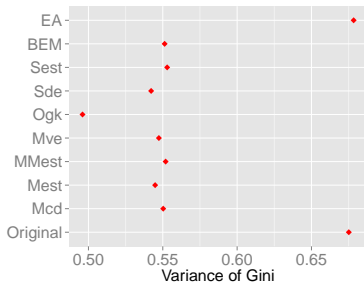
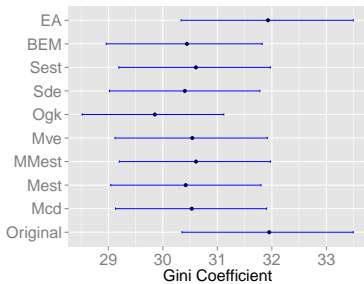
Applying outlier detection methods

- ▶ Apply univariate outlier detection methods on total annual household expenditures
 - ▶ Exclude missing values/zeros from calculations
- ▶ Apply multivariate outlier detection methods by
 - ▶ Log transforming the data
 - ▶ Impute zeros/missing values if necessary with kNN algorithm
 - ▶ BACON-EEM & EA have an internal imputation mechanism
- ▶ Estimate weighted Gini coefficient of total annual expenditures

Results for Albanian data set



Results for Albanian data set



Simulation

- ▶ Simulate such kind of data sets which can be comparable, regarding the data on household expenditure, with the ones provided by the World Bank

- ▶ Know the number and position of "true" outliers beforehand

Simulation setup

- ▶ Use Albanian data set and
 - ▶ split data into "clean" and "contaminated" data set
 - ▶ data point never flagged \Rightarrow "clean" data
 - ▶ data point flagged by at least 5 univariate outlier detection methods OR at least 6 multivariate outlier detection methods \Rightarrow "contaminated" data
- ▶ estimate location and covariance for "clean" and "contaminated" data set in a classical way
 $\Rightarrow (\boldsymbol{\mu}_{cl}, \boldsymbol{\Sigma}_{cl}), (\boldsymbol{\mu}_{co}, \boldsymbol{\Sigma}_{co})$
- ▶ Simulate data from $MVN(\boldsymbol{\mu}_{cl}, \boldsymbol{\Sigma}_{cl})$

Simulation setup

- ▶ swap observations with contaminated values generated from $MVN(\boldsymbol{\mu}_{co}, \boldsymbol{\Sigma}_{co})$
 - ▶ swap only a single cell for share of contaminated data
- ▶ Simulated data set \mathbf{X} follows the following distribution

$$\mathbf{X} \sim (1 - \epsilon)MVN(\boldsymbol{\mu}_{cl}, \boldsymbol{\Sigma}_{cl}) + \epsilon MVN(\boldsymbol{\mu}_{co}, \boldsymbol{\Sigma}_{co}) \quad ,$$

with $\epsilon \in (0, 1)$ determining the share of contaminated data points.

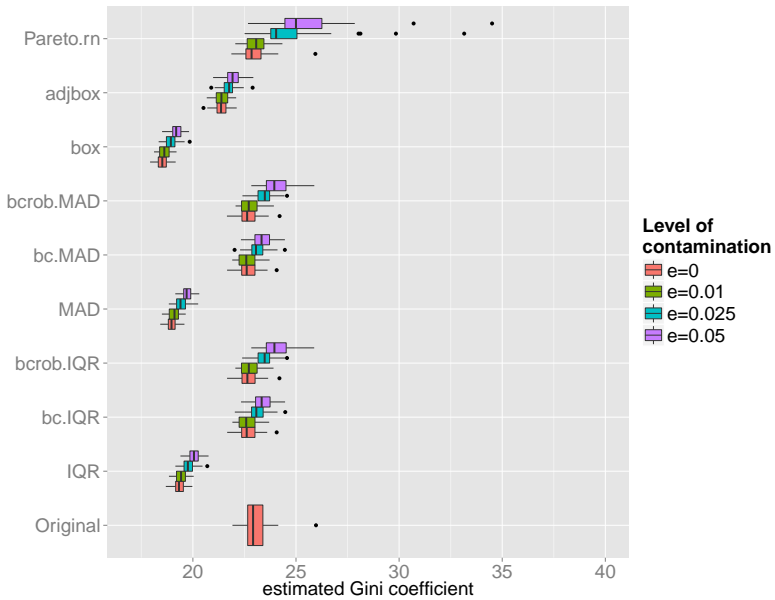
- ▶ Include missing values and sample weights from the Albanian data set

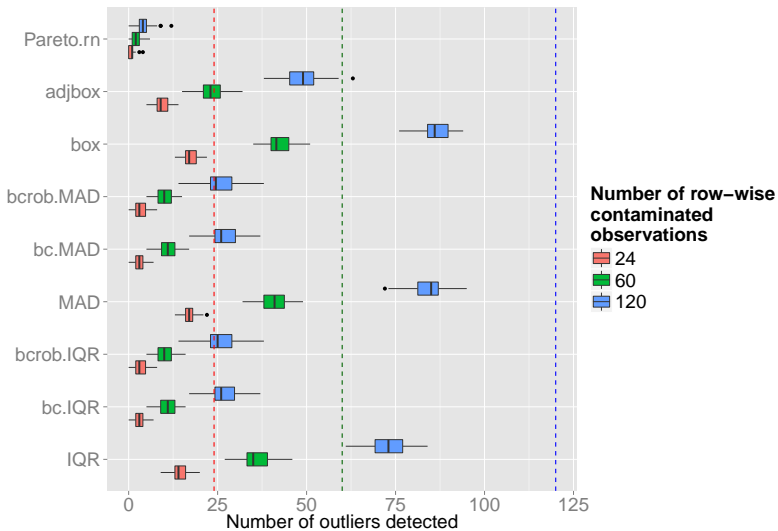
Simulation parameters

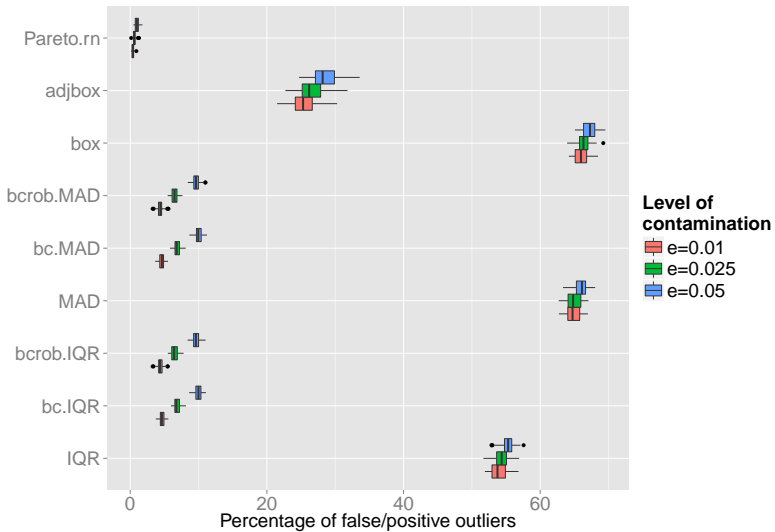
- ▶ Simulation and application of univariate and multivariate outlier detection methods is repeated 50 times
- ▶ $\epsilon \in \{0; 0.01; 0.025; 0.05\}$
- ▶ 1/3 of the contamination is cell-wise

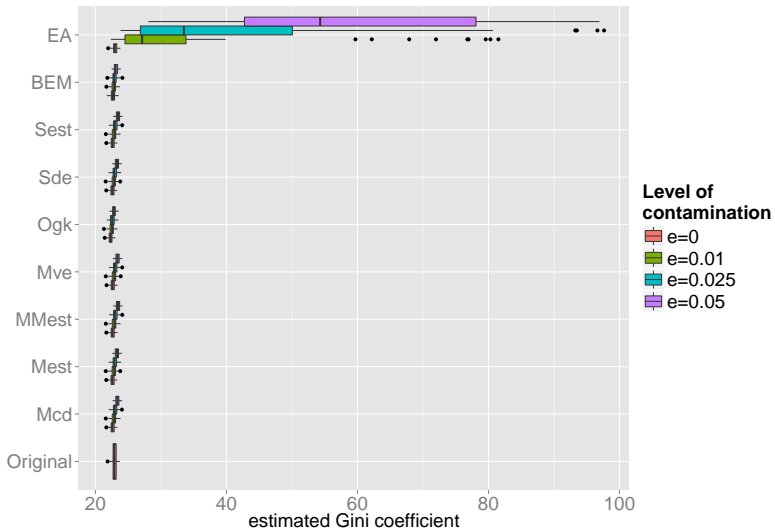
Application of outlier detection methods

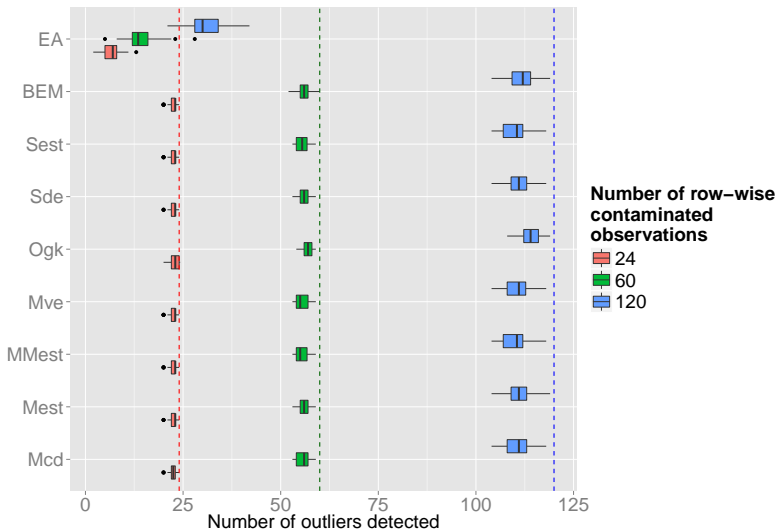
- ▶ Simulate data
- ▶ Apply outlier detection methods
 - ▶ Apply univariate methods on each of the columns of the generated data \Rightarrow results more comparable to multivariate case
- ▶ Detect and impute potential outliers
- ▶ Count correctly identified outliers and false positive outliers
- ▶ Estimate Gini coefficient of the total sum of each observation

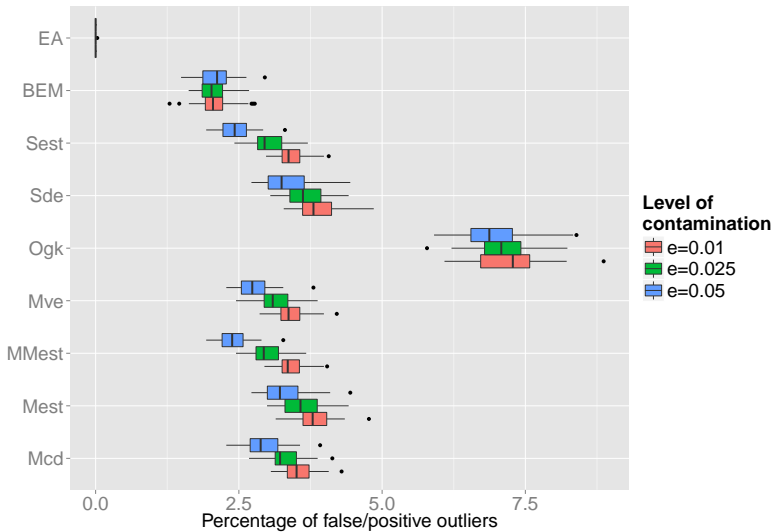












Estimates of Gini for the 5 different countries

Country	Number of households		Original	IQR	BACON EEM
Albania(2008)	3600	Gini	31.95	28.10	30.44
		Number outlier	–	121	332
India(2009)	100852	Gini	39.82	33.56	37.44
		Number outlier	–	9131	9404
Mexico(2010)	27655	Gini	44.20	37.62	42.75
		Number outlier	–	1669	2429
Malawi(2010)	12096	Gini	48.52	36.13	41.22
		Number outlier	–	1003	796
Tajikistan(2007)	4860	Gini	33.11	28.59	30.32
		Number outlier	–	244	505

Summary

- ▶ Simulation study necessary to determine performance of outlier detection methods on household expenditure data
- ▶ The simulation study presented in this work favored the BACON-EEM to be the most suitable method, but
 - ▶ simulation study favored multivariate methods in contrast to univariate methods
 - ▶ did not take into account sociodemographic criteria or household specific information

References



A. Alfons and M. Templ.

Estimation of social exclusion indicators from complex surveys: The R package laeken.

[Journal of Statistical Software](#), 54(15):1–25, 2013.



C. Béguin and B. Hulliger.

The BACON-EEM algorithm for multivariate outlier detection in incomplete survey data.

[Survey Methodology](#), 34(1):91–103, 2008.



N. Billor, A. S. Hadi, and P. F. Velleman.

BACON: Blocked adaptive computationally-efficient outlier nominators.

[Computational Statistics and Data Analysis](#), 34(3):279–298, 2000.



B. Hulliger, A. Alfons, P. Filzmoser, A. Meraner, T. Schoch, and M. Templ.

Robust methodology for laeken indicators.

Research Project Report WP4 – D4.2, FP7-SSH-2007-217322 AMELI, 2011.



T. Todorov and P. Filzmoser.

An object oriented framework for robust multivariate analysis.

[Journal of Statistical Software](#), 32(3):1–47, 2009.

This work has been funded by the World Bank under the project "Improving the quality of sample household expenditure data and the reliability of poverty and inequality measures", selection number 1157976.