

Robust Tools for Statistical Data Editing and Imputation

WADA Kazumi (Tsuda University) and TSUBAKI Hiroe (The Institute of Statistical Mathematics (ISM))

- I. M-estimators for regression (with two weight functions and scale parameters)
- II. M-estimators for generalised ratio model (with Tukey's biweight function and average absolute deviation (AAD) for scale parameter)
- III. Modified Stahel-Donoho estimators

I. M-estimators for regression:

$$y = \alpha + \beta x + \varepsilon, \quad \varepsilon \sim N(0, \sigma), \quad \beta = (\beta_1, \dots, \beta_p)$$

Features

- IRLS (Iteratively Reweighted Least Squares) based on Bienias et al. (1997) in *Statistical Data Editing 2 – Methods and Techniques* by Tukey's biweight function with AAD scale.
- R functions by Huber's weight function, and also with MAD scale are implemented for comparison
- Tukey's biweight function is more suitable with AAD scale
- AAD scale converge faster than MAD scale

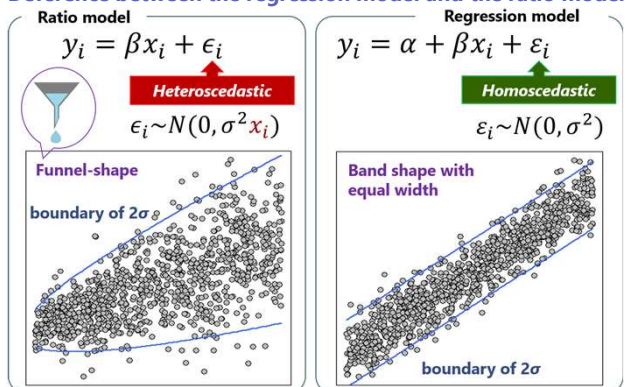
R functions available at <https://github.com/kazwd2008/IRLS>

File name	R function	Weight function	Scale parameter	
Tirls.r	Tirls.aad	Tukey	AAD	* AAD: average absolute deviation $\text{mean}(r_i - \text{mean}(r_i))$
	Tirls.mad		MAD	* MAD: median absolute deviation $\text{median}(r_i - \text{median}(r_i))$
Hirls.r	Hirls.aad	Huber	AAD	* r_i : residual
	Hirls.mad		MAD	

II. M-estimators for generalised ratio model:

$$y = \beta x + \varepsilon x^\gamma, \quad \varepsilon \sim N(0, \sigma)$$

Difference between the regression model and the ratio model



Robustification and generalisation

- Making the error term homoscedastic**
 $\varepsilon \sim N(0, \sigma^2): y_i = \beta x_i + \varepsilon_i$
 $\varepsilon \sim N(0, \sigma^2): y_i = \beta x_i + \varepsilon_i \sqrt{x_i} \quad \because \varepsilon_i = \frac{\varepsilon_i}{\sqrt{x_i}}$
- Robustification**

$$\hat{\beta}_{rob} = \frac{\sum w_i y_i}{\sum w_i x_i}$$

Quasi-residual: $\tilde{\varepsilon}_i = \frac{y_i - \hat{\beta}_{rob} x_i}{\sqrt{x_i}} - \hat{\beta}_{rob} \sqrt{x_i}$

Weight function: Tukey's biweight ($c=8$)

Scale parameter: $\hat{\sigma}_{AAD} = \frac{1}{n} \sum_{i=1}^n |\tilde{\varepsilon}_i|$ * γ : an arbitral constant

* AAD: average absolute deviation
- Generalisation**

Error term proportional to $x_i^{1/2} \Rightarrow x_i^\gamma$

model: $\frac{y_i}{x_i^\gamma} = \beta x_i^{(1-\gamma)} + \varepsilon_i$

estimator: $\hat{\beta} = \frac{\sum x_i y_i^{1-2\gamma}}{\sum x_i^{2(1-\gamma)}}$

robustification: $\hat{\beta}_{rob} = \frac{\sum w_i y_i x_i^{1-2\gamma}}{\sum w_i x_i^{2(1-\gamma)}}$

R functions available at <https://github.com/kazwd2008/IRLS>

File name	R function	Feature	Weight function	Scale parameter
RrT.r*	RrTa.aad	$\gamma=1$	Tukey	AAD
	RrTb.aad	$\gamma=1/2$		
	RrTc.aad	$\gamma=0$		
	RrTa.mad	$\gamma=1$		MAD
	RrTb.mad	$\gamma=1/2$		
	RrTc.mad	$\gamma=0$		
RrH.r*	RrHa.aad	$\gamma=1$	Huber	AAD
	RrHb.aad	$\gamma=1/2$		
	RrHc.aad	$\gamma=0$		
	RrHa.mad	$\gamma=1$		MAD
	RrHb.mad	$\gamma=1/2$		
	RrHc.mad	$\gamma=0$		
RBreds.r	Rbred	Robust estimation for generalized ratio model (γ and β are simultaneously estimated)	Tukey	AAD
	Bred	Non robust estimation for generalized ratio model (γ and β are simultaneously estimated)		

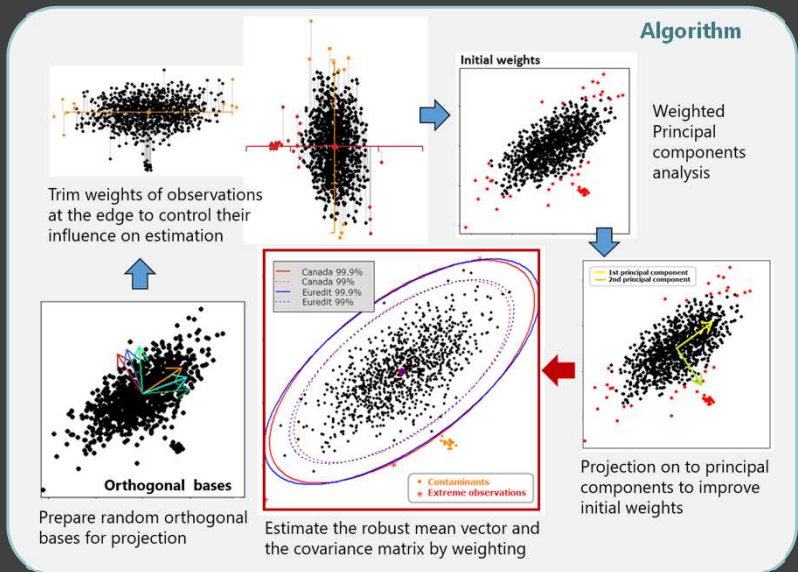
III. Modified Stahel-Donoho (MSD) estimator

Features

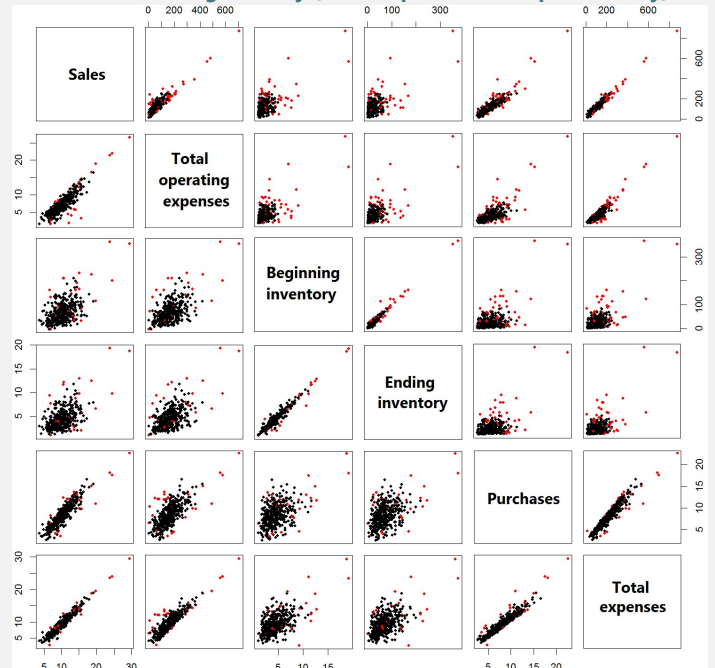
- Multivariate outlier detection method estimating mean vector and covariance matrix
- Adapted for Annual Wholesale and Retail Trade Survey (AWRTS) in Statistics Canada (Franklin & Brodeur, 1997)
- Suggestions by Béguin & Hulliger (2003) improve the performance while having a problem of the curse of dimensionality due to the increase of orthogonal bases (Wada, 2010)
- Higher performance compared with BACON, NNVE and MCD (Wada, Kawano & Tsubaki, 2020), while computationally expensive with high dimensional data

Parallelisation

- Paralleled version to cope with higher dimensional data is implemented (Wada & Tsubaki, 2013)



Manufacturing industry (Unincorporated Enterprise Survey)



Upper triangular matrix: square root transformation
 Lower triangular matrix: biquadratic root transformation

Practical application

- 2016 Economic Census for Business Activities**
 - Robust estimator for generalised ratio model ($\gamma=1/2$)
- 2019 Unincorporated Enterprise Survey**
 - Robust estimator for generalised ratio model ($\gamma=1/2$)
 - MSD estimators (single core version)