

Session 9: Cellwise outliers

Winter course, CMStatistics 2016

Mia Hubert, Peter Rousseeuw, Stefan Van Aelst

*Department of Mathematics
KU Leuven, Belgium*

December 6–7, 2016

KU LEUVEN

Outline of the course

- 1. General notions of robustness
- 2. Robustness for univariate data
- 3. Robust multivariate methods
- 4. Robust regression
- 5. Robust principal component analysis
- 6. Inference
- 7. Multivariate and functional depth
- 8. High dimensional data and sparsity
- 9. Cellwise outliers

Contamination types

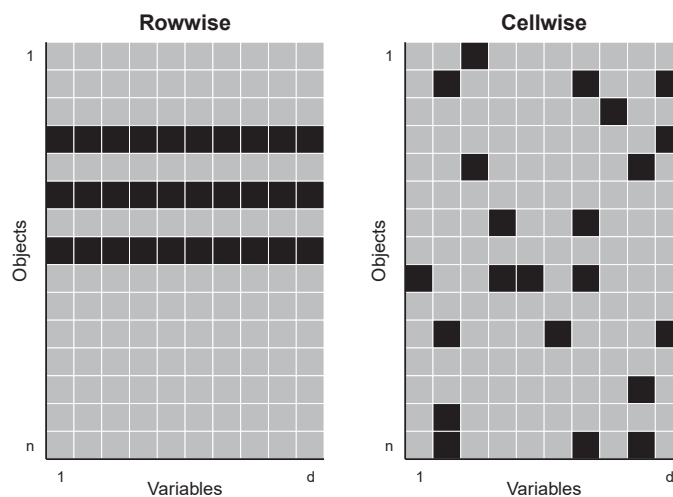
We want to analyze a data matrix \mathbf{X} with n rows (cases) and $d > 1$ columns (variables). But the data may contain outliers.

The usual **rowwise contamination model** of Tukey (1960), Huber (1964),... assumes that some **rows** x_i have been replaced by arbitrary rows.

Such outlying rows may be cases belonging to a different population. Many robust and equivariant estimators of covariance, regression, PCA,... were devised for this situation. They downweight or remove outlying rows. These methods all require at least 50% of clean rows.

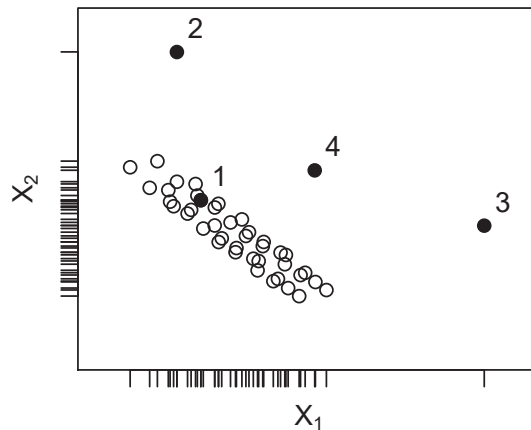
The **cellwise contamination model** of Alqallaf, Van Aelst, Yohai and Zamar (2009 AOS) assumes that some **cells** x_{ij} have been replaced.

In that case downweighting/dropping an entire row loses a lot of information. For high d it is even possible that *every* row contains an outlying cell!



But how can we **detect** the outlying cells? Agostinelli, Leung, Yohai and Zamar (2015 Test) consider each column (variable) separately.

Bivariate example



Which cell is outlying: $x_{4,1}$ or $x_{4,2}$?

Algorithm DetectDeviatingCells

The R-code is available from <http://wis.kuleuven.be/stat/robust/software> as well as a Matlab implementation in LIBRA.

Step 0: preprocessing. Check that the variables are roughly continuous, and temporarily set aside dummy variables.

It is okay when the data contain some missing values, but also set aside columns and rows with over 25% of NA's.

Verify that each remaining variable is approximately gaussian in its center, e.g. with QQ plots. If not, it is recommended to transform that variable so that the bulk of the data becomes roughly gaussian, e.g. by a robust version of the Box-Cox or Yeo-Johnson transformation.

Step 1: standardization. For each column j of \mathbf{X} we estimate

$$m_j = \text{robLoc}_i(x_{ij}) \quad \text{and} \quad s_j = \text{robScale}_i(x_{ij} - m_j)$$

where *robLoc* is a robust estimator of location (such as the sample median) and *robScale* is a robust estimator of scale about zero. Next, we standardize \mathbf{X} to \mathbf{Z} by $z_{ij} = (x_{ij} - m_j)/s_j$.

Step 2: univariate outlier detection. We define a new matrix \mathbf{U} with entries $u_{ij} = z_{ij}$ except when

$$|z_{ij}| > c$$

in which case we set $u_{ij} = \text{NA}$ (missing). The cutoff value c is taken as

$$c = \sqrt{\chi_{1,p}^2}$$

where the probability p is 99% by default.

Step 3: bivariate relations. For any two columns $h \neq j$ we compute

$$\text{cor}_{jh} = \text{robCorr}_i(u_{ij}, u_{ih})$$

where **robCorr** is a robust correlation measure.

We only use the relation between variables j and h when

$$|\text{cor}_{jh}| \geq \text{corlim}$$

in which $\text{corlim} = 0.5$ by default. Variables j satisfying this for some $h \neq j$ will be called **connected**. The others are called **standalone** variables.

For the pairs (j, h) with $|\text{cor}_{jh}| \geq \text{corlim}$ we also compute

$$b_{jh} = \text{robSlope}_i(u_{ij} | u_{ih})$$

where **robSlope** computes the slope of a robust no-intercept regression line that predicts variable j from variable h .

Step 4: predicted values. Next we compute predicted values \hat{z}_{ij} for all cells. For each variable j we consider the set H_j consisting of all variables h with $|cor_{jh}| \geq \text{corrlim}$, including j itself. For all $i = 1, \dots, n$ we then set

$$\hat{z}_{ij} = \frac{\sum_h w_{jh} b_{jh} u_{ih}}{\sum_h w_{jh}}$$

where $w_{jh} = |cor_{jh}|$. Other choices are possible, such as a weighted median.

Step 5: deshrinkage. Computing \hat{z}_{ij} often shrinks the scale of the entries, which is undesirable. To counteract the shrinkage we replace \hat{z}_{ij} by

$$\hat{z}_{ij} \text{ robSlope}_{i'}(z_{i'j} | \hat{z}_{i'j})$$

for all connected variables j .

Step 6: flagging cellwise outliers. Compute the standardized cell residuals

$$r_{ij} = \frac{z_{ij} - \hat{z}_{ij}}{\text{robScale}_{i'}(z_{i'j} - \hat{z}_{i'j})}.$$

In each column j we then flag all cells with $|r_{ij}| > c$ as cellwise outliers.

Next we assemble the ‘imputed’ matrix Z_{imp} given by

$$(Z_{imp})_{ij} = \begin{cases} \hat{z}_{ij} & \text{if } z_{ij} \text{ was flagged or NA} \\ z_{ij} & \text{otherwise.} \end{cases}$$

Step 7: flagging rowwise outliers. The method can also flag some outlying rows i based on the standardized cell residuals r_{ij} .

For multivariate gaussian data without outliers we have $r_{ij} \approx N(0, 1)$ so the cdf of r_{ij}^2 is approximately the cdf F of χ_1^2 . This leads us to the criterion

$$T_i = \text{ave}_{j=1}^d F(r_{ij}^2) .$$

We then robustly standardize the T_i and flag the rows i for which the standardized T_i exceeds the cutoff $\sqrt{\chi_{1,p}^2}$.

Note that when row i has an unusually large T_i this doesn't 'prove' that i is a member of a different population, but at least it is worth looking into.

Also note that although the T_i can flag **some** rowwise outliers, there are types of rowwise outliers that it may not detect. Therefore, it is recommended to use **rowwise robust** methods in subsequent analyses of the data.

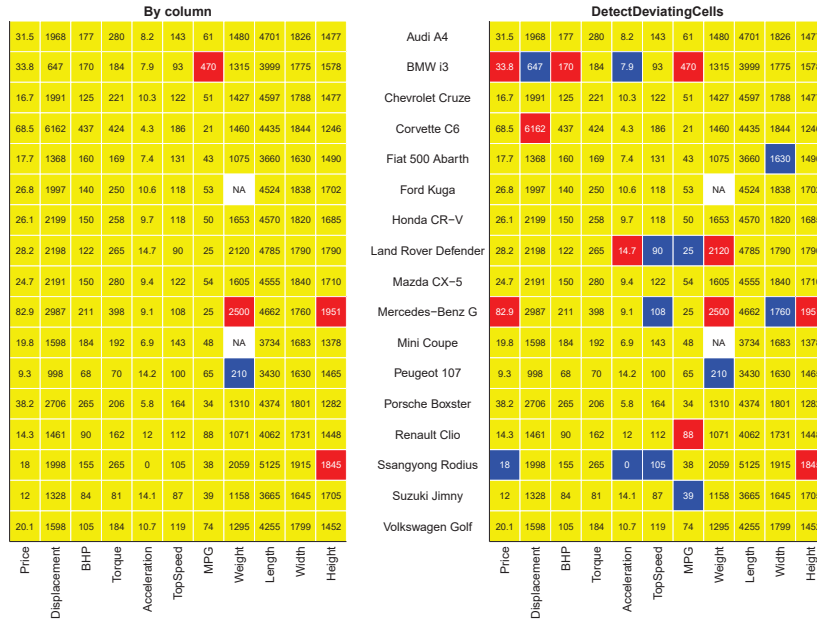
Step 8: unstandardize. Turn the imputed matrix Z_{imp} into an imputed matrix X_{imp} by undoing the standardization. The output of DetectDeviatingCells is X_{imp} together with the list of cells and rows flagged as outlying.

DetectDeviatingCells does **not** require over 50% of clean rows!
It is equivariant for translations, for diagonal linear transformations, and for permuting rows and columns, but not for general linear transformations.

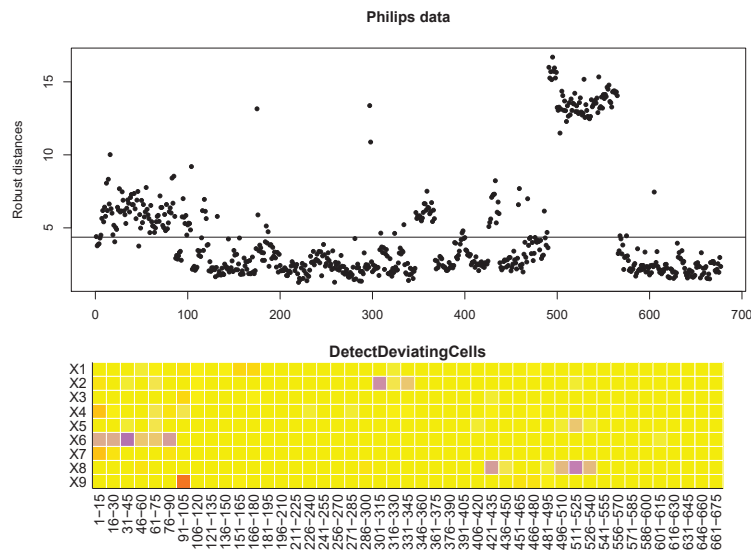
As a byproduct, DetectDeviatingCells **imputes** all NA's in the data.
This is far less efficient than the EM algorithm when the data are outlier-free, but it is more robust against cellwise outliers.

Note that DetectDeviatingCells starts from relations between 2 variables.
If instead we would fit a q -dimensional model to each set of $q > 2$ variables by a rowwise robust method, the computation time would explode and those fits would be less robust due to cellwise outlier propagation. The current algorithm runs in $O(nd^2)$ time and $O(nd)$ space, and speedups are possible.

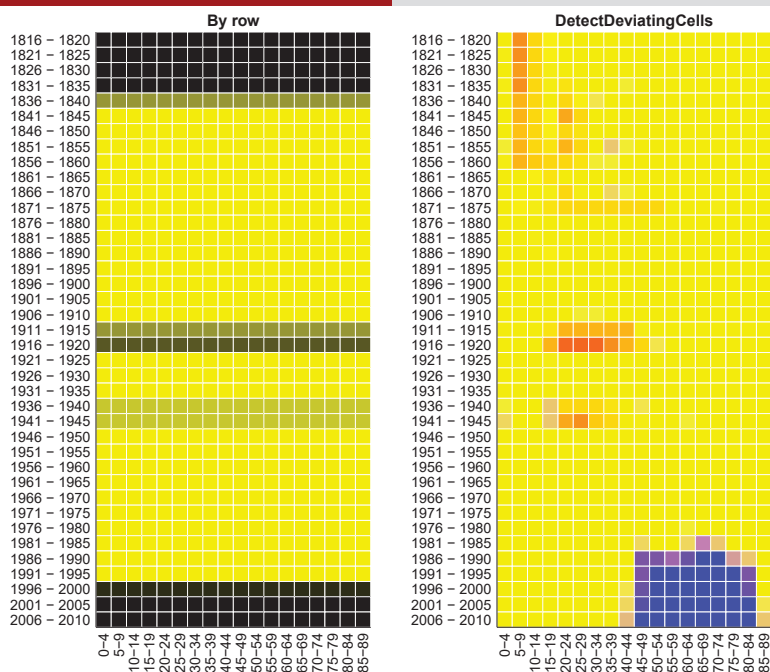
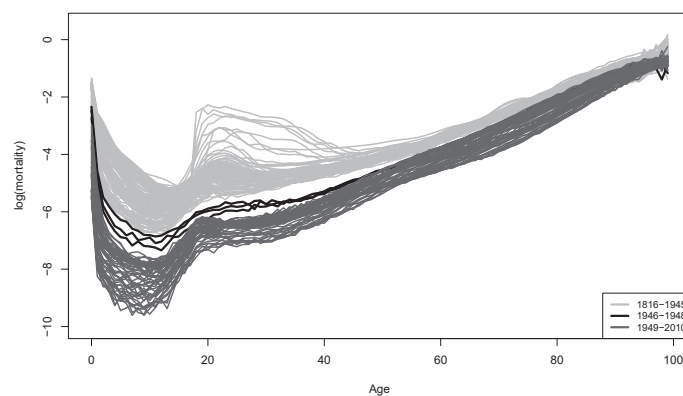
Example: `library(robustHD); data(TopGear)`



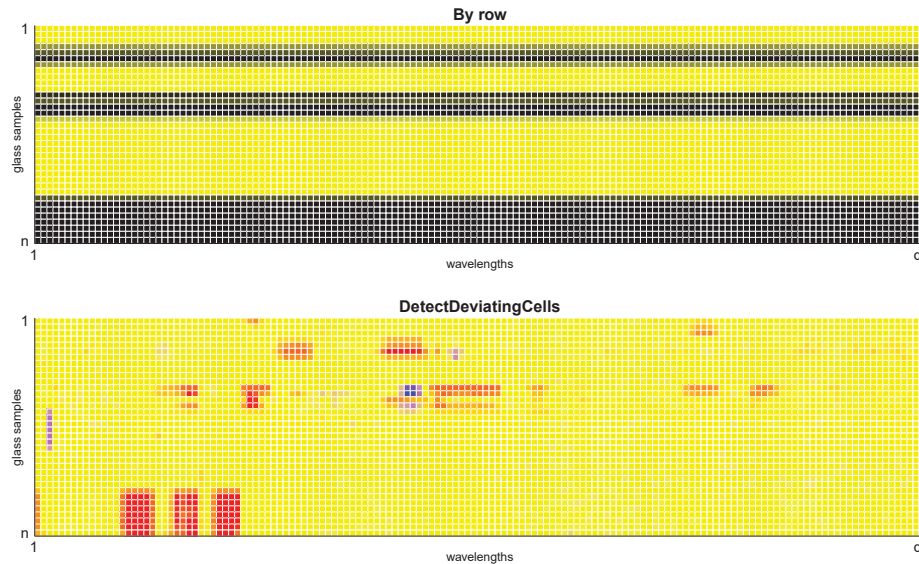
Example: Philips production line, $n=677$ parts, $d=9$ characteristics. Robust distances from FastMCD, and transposed cell map with blocks of 15 products.



Example: mortality by age for males in France, from 1816 to 2010 (from www.mortality.org)



Example: $n=180$ archeological glass samples, spectra with $d=750$ wavelengths, so more dimensions than objects!



After running the algorithm DetectDeviatingCells:

1. Ideally, the user looks at the anomalous cells and whether their values are higher or lower than predicted, and makes sense of what is going on. This may lead to a better understanding of the data pattern, to changes in the way the data are collected/measured, to dropping certain rows or columns, to transforming variables, to changing the model,...
2. If the data set is too large for visual inspection of the results or the analysis is automated, the anomalous cells can be set to missing after which the data set is analyzed by a method appropriate for incomplete data.
3. If no such method is available, one can analyze the imputed data set X_{imp} produced by DetectDeviatingCells, which has no missings.

In 2. and 3. one can drop the flagged rows before taking the step. If that step is carried out by a sparse method such as the Lasso (Tibshirani 1996, Städler-Stekhoven-Bühlmann 2014) or another form of variable selection: look more closely at the deviating cells in the variables that were selected.

Comparison with existing method

We compare with the univariate **Gervini-Yohai (GY)** filter (2002 AOS).

The **uncontaminated** data are gaussian with mean zero and a covariance matrix with unit diagonal. We use two types of correlation matrices:

- **ALYZ**: the random correlation matrices generated by Agostinelli, Leung, Yohai and Zamar (2015 Test). These yield relatively low correlations.
- **A09**: the true correlation matrix is generated as

$$\rho_{jh} = (-0.9)^{|j-h|}.$$

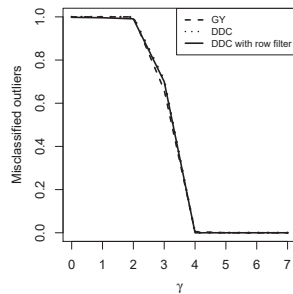
The two types of **contamination** are generated as follows:

- **cells**: A random subset of the nd cells are replaced by the constant γ .
- **rows**: compute the last eigenvector v of the true covariance matrix C . Rescale v to the typical size of a data point, by making

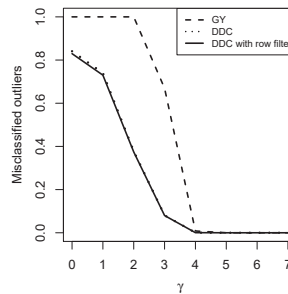
$$MD_C^2(v) = E[Y^2] = d \quad \text{where} \quad Y^2 \sim \chi_d^2.$$

Then add point contamination at γv .

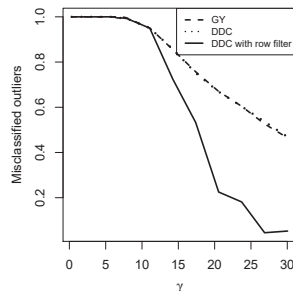
ALYZ model, 10% cells, d=20



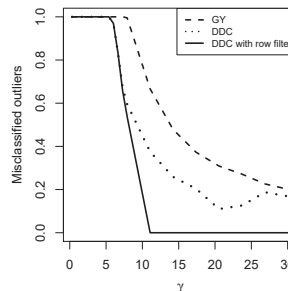
A09 model, 10% cells, d=20



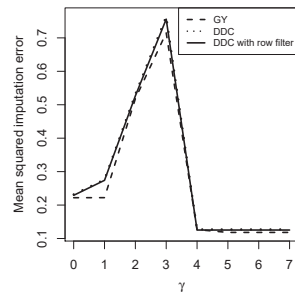
ALYZ model, 10% rows, d=20



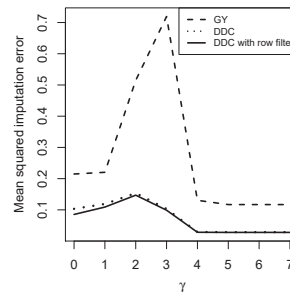
A09 model, 10% rows, d=20



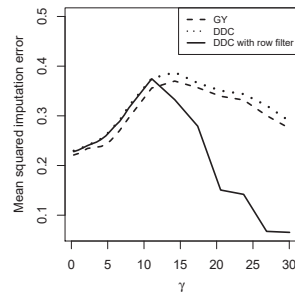
ALYZ model, 10% cells, d=20



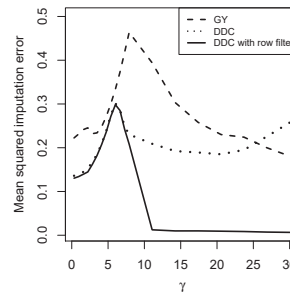
A09 model, 10% cells, d=20



ALYZ model, 10% rows, d=20

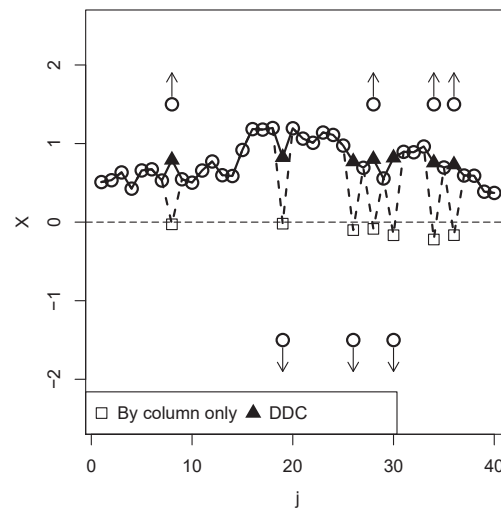


A09 model, 10% rows, d=20



If the true correlations are given by $\rho_{jh} = 0.99^{|j-h|}$ then the rows of \mathbf{X} look like autocorrelated time series.

Compare DetectDeviatingCells imputation to median imputation by column:



Multivariate location and scatter

The 2SGS method (Agostinelli, Leung, Yohai and Zamar 2015) consists of:

- 1 Apply the **GY** filter to the columns of \mathbf{X} and set the flagged cells to NA
- 2 Apply the **GSE** estimator of Danilov, Yohai and Zamar (2012 JASA) to this incomplete data set, yielding $\hat{\boldsymbol{\mu}}$ and $\hat{\boldsymbol{\Sigma}}$.

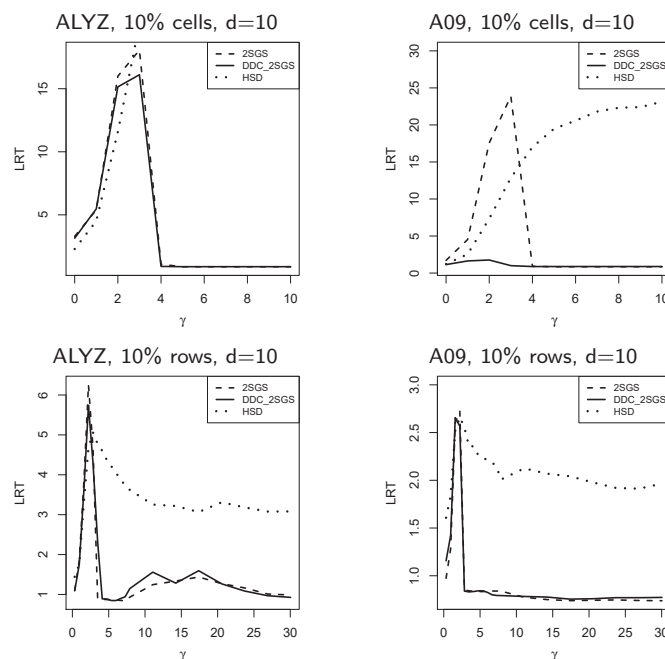
Our version replaces GY in step 1 by DetectDeviatingCells. When an entire row is flagged, we remove it.

Removing a row will count as removing all of its d cells.

We also include the HSD estimator of Van Aelst et al (2012).

Measure how far $\hat{\boldsymbol{\Sigma}}$ is from the true $\boldsymbol{\Sigma}$ by

$$\text{LRT} = \text{trace}(\hat{\boldsymbol{\Sigma}}\boldsymbol{\Sigma}^{-1}) - \log(\det(\hat{\boldsymbol{\Sigma}}\boldsymbol{\Sigma}^{-1}) - d) .$$



Three-step versions

Note that the default GSE estimator starts by computing a version of the minimum volume ellipsoid (MVE) suitable for incomplete data, which is rather time-consuming.

Therefore Agostinelli, Leung, Yohai and Zamar (Test 2015, rejoinder) also proposed **Fast2SGS**:

- 1 Apply the GY filter to the columns of \mathbf{X} and set the flagged cells to NA.
- 2 Apply the MVE.S estimator (an S-estimator starting from MVE) to the complete data set obtained by imputing the NA's by the median of their column;
- 3 Apply GSE to the data set with the NA's but now starting from the $\hat{\mu}$ and $\hat{\Sigma}$ obtained in step 2.

Our version again replaces GY in step 1 by DetectDeviatingCells.

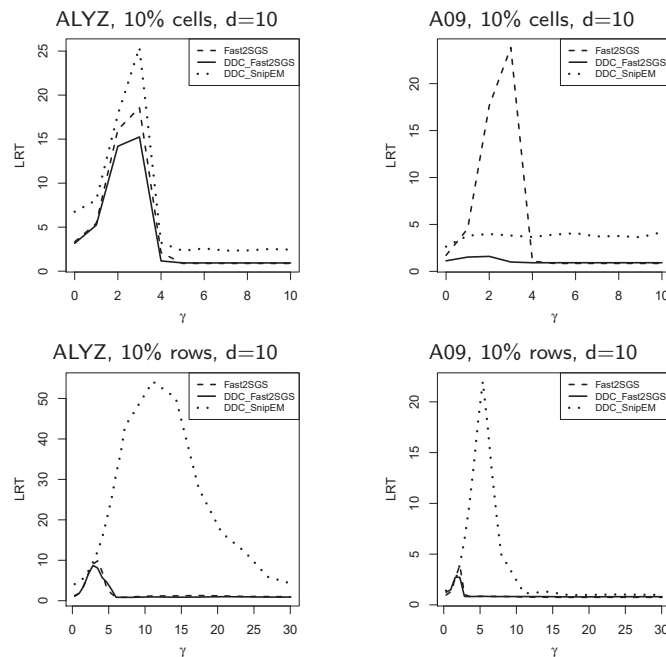
Another approach to estimating μ and Σ is the **SnipEM** method of Farcomeni (2014 Technometrics). This searches for the subset of cells such that 'snipping' them (setting them to NA) and then running the EM algorithm on this incomplete data set yields the highest partial likelihood.

The SnipEM algorithm requires a good initial subset.

Agostinelli et al (2015) found that SnipEM works fairly well for cell contamination but is insufficiently robust for row contamination.

Therefore we take the following steps:

- 1 Apply DetectDeviatingCells to \mathbf{X} with its imputation of the flagged cells and remove the flagged rows.
- 2 Apply the MVE.S estimator to this modified data set and flag the outlying rows it detects too.
- 3 Apply SnipEM to the data set without the rows flagged in steps 1 and 2. As initial subset we use the cells flagged in step 1.



Regression

Leung, Zhang and Zamar (2016 CSDA) proposed the **3S** regression method:

- 1 Apply a generalization of the GY filter to the columns of \mathbf{X} (not y) and set the flagged cells to NA
- 2 Apply GSE to the incomplete data set $[\mathbf{X}|y]$, yielding $\hat{\boldsymbol{\mu}}$ and $\hat{\boldsymbol{\Sigma}}$
- 3 Partitioning $\hat{\boldsymbol{\mu}}$ and $\hat{\boldsymbol{\Sigma}}$ as

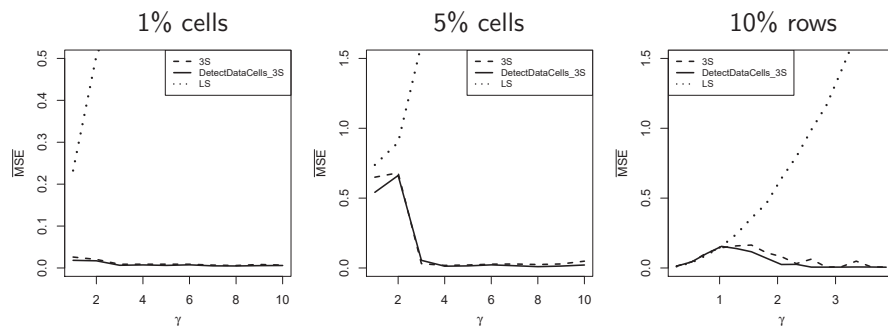
$$\hat{\boldsymbol{\mu}} = \begin{pmatrix} \hat{\mu}_x \\ \hat{\mu}_y \end{pmatrix} \quad \text{and} \quad \hat{\boldsymbol{\Sigma}} = \begin{pmatrix} \hat{\Sigma}_{xx} & \hat{\Sigma}_{xy} \\ \hat{\Sigma}_{yx} & \hat{\Sigma}_{yy} \end{pmatrix}$$

yields the slope and intercept estimates

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \hat{\Sigma}_{xx}^{-1} \hat{\Sigma}_{xy} \\ \hat{\alpha} &= \hat{\mu}_y - \hat{\boldsymbol{\beta}}^t \hat{\mu}_x. \end{aligned}$$

This is again robust to both cellwise and rowwise outliers.
Our version replaces step 1 by DetectDeviatingCells.

Average $\text{MSE}(\hat{\beta}, \hat{\alpha})$ under ALYZ with cell and row contamination ($d=15$, $n=300$):



Under ALYZ, the average lengths of the confidence intervals around the coefficients is similar to those found by Leung-Zhang-Zamar, for the same contamination settings.







The same holds for the average coverage rates of those intervals.

Conclusions

Our approach turns high dimensionality (usually considered a **curse**) into an advantage, as having more variables may improve the accuracy of the predicted cells.

DetectDeviatingCells and methods starting from it perform about as well as the robust methods of Zamar et al (GY filter, 2SGS, Fast2SGS, 3S regression) when the correlations between the variables are small to moderate, and better when there are higher correlations.

References

-  Agostinelli, C., Leung, A., Yohai, V.J., and Zamar, R.H. (2015), "Robust estimation of multivariate location and scatter in the presence of cellwise and casewise contamination," *Test*, 24, 441–461.
-  Alqallaf, F., Van Aelst, S., Yohai, V., and Zamar, R.H. (2009), "Propagation of outliers in multivariate data," *The Annals of Statistics*, 37, 311–331.
-  Gervini, D., and Yohai, V.J. (2002), "A class of robust and fully efficient regression estimators," *The Annals of Statistics*, 30, 583–616.
-  Leung, A., Zhang, H., and Zamar, R. (2016), "Robust regression estimation and inference in the presence of cellwise and casewise contamination," *Computational Statistics and Data Analysis*, 99, 1–11.
-  Rousseeuw, P.J., and Van den Bossche, W. (2016), "Detecting deviating data cells," arXiv 1601.07251 .
-  Van Aelst, S., Vandervieren, E., and Willems, G. (2012), "A Stahel-Donoho estimator based on huberized outlyingness," *Computational Statistics and Data Analysis*, 56, 531–542.