

# Avoiding the pitfalls of S-estimators with categorical predictors

Manuel Koller

Icors 2012, University of Vermont



# 1 Outline

We want to compute MM-estimates on data with many factors.

## 2 Example: NO<sub>x</sub>Emissions

*A typical medium sized environmental data set with hourly measurements of NO<sub>x</sub> pollution content in the ambient air.*

The dataset consists of **8088** observations on the following 4 variables.

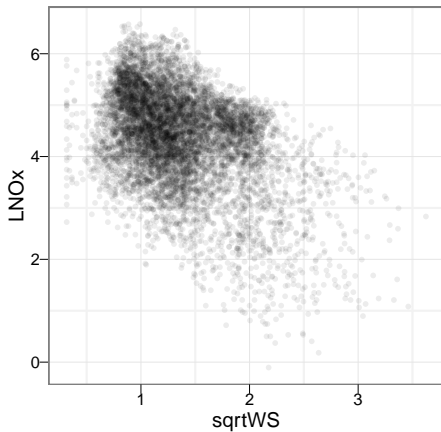
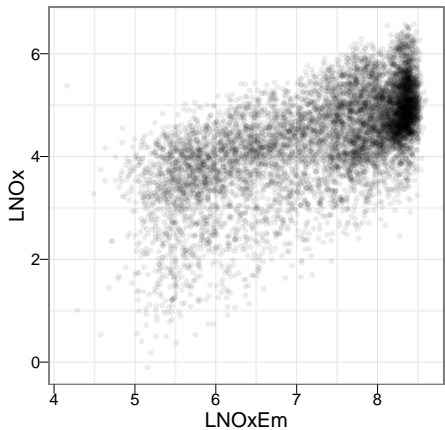
**LNO<sub>x</sub>** log of hourly mean of NO<sub>x</sub> concentration in ambient air [ppb] next to a highly frequented motorway.

**LNO<sub>x</sub>Em** log of hourly sum of NO<sub>x</sub> emission of cars on this motorway in arbitrary units.

**sqrtWS** Square root of wind speed [m/s].

**julday** day number, a factor with levels '373' ... '730', typically with 24 hourly measurements.

(The data set comes with the R package `robustbase`.)



When using (an older version) of `lmrob` on such a dataset, it usually failed to compute the initial S-estimate:

```
> lmrob(LNOx ~ ., data = NOxEmissions)
```

```
Too many singular resamples
```

```
Aborting fast_s_w_mem()
```

```
Error in lmrob.S(x, y, control = control) :
```

```
  C function R_lmrob_S() exited prematurely
```

(`lmrob` is a function in the R package `robustbase`. It computes MM-estimates.)

### 3 The problem

MM-estimates consist of

- an initial S-estimate with high breakdown point,
- a final M-estimate with high efficiency.

The algorithm for computing the initial S-estimate usually involves sub-sampling, which is problematic for such data.

*S*-estimates are defined as

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \hat{\sigma}(r(\beta)),$$

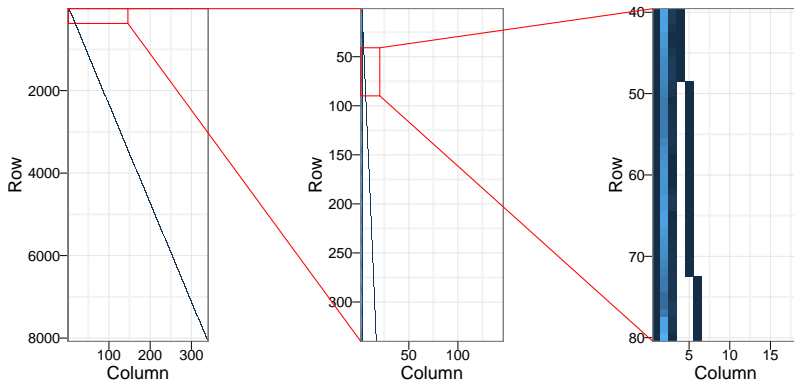
where  $\hat{\sigma}(\cdot)$  is an M-estimate of scale and  $r(\cdot)$  are the residuals.

**The algorithm to compute S-estimates**, simplified:

1. Take a random subsample of size  $p$ , the number of parameters; solve the problem on the subsample.
2. Find local minimum starting from the solution found in 1.
3. Repeat for a fixed number of times.

The S-estimate is then the result of 2. with the smallest scale estimate.

Zooming in the design matrix from the NO<sub>x</sub>Emissions example.



A valid subsample in this example must not contain any 0 only columns.



## 4 Strategies before robustbase version 0.9

The user had the following options:

- increase the allowed number of singular subsamples,
- use `lmRob` (R package `robust`) that uses an M/S-estimate as initial estimate for such datasets,
- switch to another estimator to solve the problem, e.g., M or L1.

In other words: wait a long time or ditch `lmrob`.

So... we improved `lmrob`!

As of version 0.9 of `robustbase`:

- Support for M/S-estimates.
- Improved algorithm to compute S-estimates to deal with such datasets as well (nonsingular subsampling).

## 5 Nonsingular subsampling

“Algorithm”:

Build up the subsample observation by observation, checking for collinearities each time. If an observation introduces collinearities, then discard it and continue with another one.

This always works, except if:

- such a subsample does not exist, i.e., the design matrix is not of full rank, or,
- the design matrix is ill-conditioned, causing numerical problems.

And the best thing about it:

Checking for collinearities comes (almost) for free.

## 6 LU-factorization

The LU-decomposition of a nonsingular matrix  $\mathbf{A}$  consists of

- a lower triangular matrix  $\mathbf{L}$ ,
- an upper triangular matrix  $\mathbf{U}$ , and,
- a permutation matrix  $\mathbf{P}$  (to avoid divisions by 0),

such that

$$\mathbf{PA} = \mathbf{LU}.$$

This is used, e.g., for solving linear systems of equations  $\mathbf{A}\beta = \mathbf{b}$ , since

$$\beta = \mathbf{U}^{-1}\mathbf{L}^{-1}\mathbf{P}^{-1}\mathbf{b}.$$

Basic algorithm to compute the LU-factorization: a series of Gaussian eliminations (Doolittle's algorithm).

Example: Compute LU-factorization of the (singular) matrix  $\mathbf{A}$ .

$$\mathbf{A} = \begin{pmatrix} 6 & 18 & 24 & 12 \\ 3 & 17 & 20 & 2 \\ 2 & 4 & 6 & 15 \\ 3 & 5 & 8 & 8 \end{pmatrix}$$

Note: The third column is the sum of the first two columns.

LU-doolittle step 1:

$$\begin{matrix} & \mathbf{L}_1 & & & & \mathbf{A}^{(1)} \\ \left( \begin{array}{cccc} 1 & 0 & 0 & 0 \\ -0.5 & 1 & 0 & 0 \\ -0.33 & 0 & 1 & 0 \\ -0.5 & 0 & 0 & 1 \end{array} \right) & & & & \left( \begin{array}{cccc} 6 & 18 & 24 & 12 \\ 0 & 8 & 8 & -4 \\ 0 & -2 & -2 & 11 \\ 0 & -4 & -4 & 2 \end{array} \right) \end{matrix}$$

LU-doolittle step 2:

$$\begin{matrix} & L_2 & & \\ \left( \begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0.25 & 1 & 0 \\ 0 & 0.5 & 0 & 1 \end{array} \right) & & & \left( \begin{array}{cccc} 6 & 18 & 24 & 12 \\ 0 & 8 & 8 & -4 \\ 0 & 0 & 0 & 10 \\ 0 & 0 & 0 & 0 \end{array} \right) \\ & & & A^{(2)} \end{matrix}$$

LU-doolittle step 3:

$$\begin{array}{c}
 L_3 \\
 \left( \begin{array}{cccc}
 1 & 0 & 0 & 0 \\
 0 & 1 & 0 & 0 \\
 0 & 0 & 1 & 0 \\
 0 & 0 & NaN & 1
 \end{array} \right)
 \end{array}
 \quad
 \begin{array}{c}
 A^{(3)} \\
 \left( \begin{array}{cccc}
 6 & 18 & 24 & 12 \\
 0 & 8 & 8 & -4 \\
 0 & 0 & 0 & 10 \\
 NaN & NaN & NaN & NaN
 \end{array} \right)
 \end{array}$$



**Better:** Gaxpy variant of LU-factorization algorithm.

This variant of the algorithm does the operations in a different order. It only calculates the elements of  **$U$**  when they are actually needed.

LU-gaxpy step 1:

$$\begin{matrix} & \mathbf{L}^{(1)} & & \mathbf{U}^{(1)} \\ \left( \begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0.5 & 1 & 0 & 0 \\ 0.33 & 0 & 1 & 0 \\ 0.5 & 0 & 0 & 1 \end{array} \right) & & \left( \begin{array}{cccc} 6 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right) \end{matrix}$$

LU-gaxpy step 2:

$$\begin{matrix} & \mathbf{L}^{(2)} & & \mathbf{U}^{(2)} \\ \left( \begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0.5 & 1 & 0 & 0 \\ 0.33 & -0.25 & 1 & 0 \\ 0.5 & -0.5 & 0 & 1 \end{array} \right) & & \left( \begin{array}{cccc} 6 & 18 & 0 & 0 \\ 0 & 8 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right) \end{matrix}$$

LU-gaxpy step 3:

$$\begin{matrix} & \mathbf{L}^{(3)} & & & & \mathbf{U}^{(3)} & & & & \\ \left( \begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0.5 & 1 & 0 & 0 \\ 0.33 & -0.25 & 1 & 0 \\ 0.5 & -0.5 & \text{NaN} & 1 \end{array} \right) & & & & \left( \begin{array}{cccc} 6 & 18 & 24 & 0 \\ 0 & 8 & 8 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right) & & & & & \end{matrix}$$

LU-gaxpy step 2:

$$\begin{matrix} & \mathbf{L}^{(2)} & & & & \mathbf{U}^{(2)} & & & & \\ \left( \begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0.5 & 1 & 0 & 0 \\ 0.33 & -0.25 & 1 & 0 \\ 0.5 & -0.5 & 0 & 1 \end{array} \right) & & & & & \left( \begin{array}{cccc} 6 & 18 & 0 & 0 \\ 0 & 8 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right) & & & & \end{matrix}$$

Two key facts about LU-gaxpy:

- Collinearities are detected immediately.
- In the  $i$ -th step, the algorithm only touches columns 1 to  $i$ .

Now consider the transposed design matrix:

columns = observations,  
rows = predictor variables.

Applying the LU-gaxpy, we only need to repeat one step if an observation introduces collinearity. All computations from the previous steps are still valid.

The LU-factorization is needed anyway, so the only extra work comes from repeating steps in case of collinearities.

## 7 Where's the random part?

Permute the observations in the design matrix.



## 8 The nonsingular subsampling algorithm

1. Permute the observations in the design matrix randomly.
2. Run LU-gaxpy step by step on the transposed design matrix, discarding any observation that introduces collinearities.
3. Use computed LU-factorization to solve the least-squares problem.

## 9 Avoiding numerical problems

Sometimes, for ill-conditioned design matrices, the nonsingular subsampling algorithm falsely declares a subsample nonsingular.

Preconditioning the design matrix helps to avoid this problem. `lmrob` uses a technique called matrix equilibration on the whole design matrix.

Instead of

$$\mathbf{A}\beta = y,$$

we solve

$$(\mathbf{D}_{\text{row}}\mathbf{A}\mathbf{D}_{\text{col}})\bar{\beta} = \mathbf{D}_{\text{row}}y, \quad \beta = \mathbf{D}_{\text{col}}\bar{\beta}.$$

The diagonal matrices  $\mathbf{D}_{\text{col}}$  and  $\mathbf{D}_{\text{row}}$  need to be computed only once for the whole design matrix.

## 10 M/S-estimates (Maronna & Yohai, 2000)

Split the design matrix in a categorical part and a continuous part.

- For the categorical part, use an M-estimate (usually L1), while
- for the continuous part, use an S-estimate.

This avoids computational difficulties on the categorical part while keeping the better robustness for the continuous part.

In formulas:

$$y_i = \mathbf{x}_{1i}^T \beta_1 + \mathbf{x}_{2i}^T \beta_2 + \varepsilon_i.$$

Then use an S-estimate for the continuous part  $\beta_2$ :

$$\hat{\beta}_2 = \underset{\beta_2}{\operatorname{argmin}} \hat{\sigma}(r(\beta_1^*(\beta_2), \beta_2)),$$

and an M-estimate for the categorical part  $\beta_1$ :

$$\beta_1^*(\beta_2) = \underset{\beta_1}{\operatorname{argmin}} \sum_{i=1}^n \rho(y_i - \mathbf{x}_{1i}^T \beta_1 - \mathbf{x}_{2i}^T \beta_2).$$

**What about interactions of categorical and continuous variables?**

## 11 Comparison

How long does it take to fit a multiple linear regression model?

$$\text{LNOx} \sim 1 + \text{LNOxE}_m + \text{sqrtWS} + \text{julday}$$

Estimator	Running time [s]
Least Squares (lm)	1.213
L1 (lmrob.lar)	2.680
S, simple subsampling (max tries = $\infty$ )	> 2592000
M/S	919.513
S, nonsingular subsampling	1082.109

Design matrix:  $n = 8088$ ,  $p = 340$ .

## Conclusions

- `lmrob` of the R-package `robustbase` is now suitable also for datasets with many factors, even when some levels have low frequency.
- Nonsingular subsampling allows us to use the regular S-estimate. This does not require extra work for easy problems.

## References

- M. Koller** (to appear). Nonsingular subsampling for S-estimators with categorical predictors.
- M. Salibian-Barrera and V. J. Yohai** (2006). A fast algorithm for S-regression estimates. *Journal of Computational and Graphical Statistics*, 15(2), 414–427.
- R. A. Maronna and V. J. Yohai** (2000). Robust regression with both continuous and categorical predictors. *Journal of Statistical Planning and Inference*, 89, 197–214.

Data:  $p \times p$  matrix  $\mathbf{A}$ .

Result: Matrices  $\mathbf{L}$  and  $\mathbf{U}$ .

$$1 \quad \mathbf{A}^{(0)} \leftarrow \mathbf{A}$$

2 for  $j \leftarrow 1$  to  $p$  do

$$3 \quad \mathbf{L}_j \leftarrow \begin{pmatrix} 1 & & & & \\ & \dots & & & \\ & & 1 & & \\ & & -\frac{a_{j+1,j}^{(j-1)}}{a_{j,j}^{(j-1)}} & \dots & \\ & & \vdots & \dots & \\ & & -\frac{a_{p,j}^{(j-1)}}{a_{j,j}^{(j-1)}} & & \\ & & & & 0 \end{pmatrix}$$
$$4 \quad \mathbf{A}^{(j)} \leftarrow \mathbf{L}_j \mathbf{A}^{(j-1)}$$

$$5 \quad \mathbf{U} \leftarrow \mathbf{A}^{(p)}$$

$$6 \quad \mathbf{L} \leftarrow \mathbf{I} + \sum_{j=1}^p (\mathbf{I} - \mathbf{L}_j)$$

Algorithm 1: LU-doolittle (without pivoting).



**Data:**  $n \times p$  matrix  $\mathbf{X}$ , response vector  $y$ , singularity threshold  $\varepsilon$ .

**Result:** Return code (0 for success, otherwise failing step), initial estimate  $\hat{\beta}$ .

```

1  $\mathbf{U} \leftarrow \mathbf{0}$ ;  $\mathbf{L} \leftarrow \mathbf{I}$ ;  $s \leftarrow 1 : p$ ;  $k \leftarrow 1$ 
2  $t \leftarrow \text{perm}(1 : n)$ ;  $\mathbf{A} \leftarrow \mathbf{X}_{t,1:p}^T$ ;  $y \leftarrow y_t$ 
3 for  $j$  in 1 to  $p$  do
4   if  $j == 1$  then  $v_{1:p} \leftarrow \mathbf{A}_{1:p,k}$ 
5   else
6      $\mathbf{U}_{1:j-1,j} \leftarrow \mathbf{L}_{1:j-1,1:j-1}^{-1} \mathbf{A}_{1:j-1,k}$ 
7      $v_{j:p} \leftarrow \mathbf{A}_{j:p,k} - \mathbf{L}_{j:p,1:j-1} \mathbf{U}_{1:j-1,j}$ 
8   if  $j < p$  then
9     if  $|v_j| \geq \varepsilon$  then
10       $s_j \leftarrow k$ 
11       $\mathbf{L}_{j+1:p,j} \leftarrow v_{j+1:p} / v_j$ 
...
...

```

...  
for  $j$  in 1 to  $p$  do

```

12   |   ...
13   |   if  $|v_j| < \varepsilon$  then
14   |       |   if  $k < n$  then
15   |       |       |    $k \leftarrow k + 1$ 
16   |       |       |   Goto 4
17   |       |   else
18   |       |       |   return  $j$ 
19   |    $U_{jj} \leftarrow v_j$ 
20   |    $k \leftarrow k + 1$ 

```

20  $\hat{\beta} \leftarrow \mathbf{L}^{-T} \mathbf{U}^{-T} y_s$

21 return 0,  $\hat{\beta}$

**Algorithm 2:** Nonsingular subsampling using modified LU-gaxpy (without pivoting).

$1 : p - 1 = (1, 2, \dots, p - 1)$ .