

Erasmus
School of
Economics

Sparse Robust Regression and Model Selection

Andreas Alfons

ICORS Workshop, September 20, 2021

Content

- ① Motivation
- ② LARS and lasso
- ③ Robust least angle regression
- ④ Sparse least trimmed squares (trimmed lasso)
- ⑤ Penalized S- and MM-estimator
- ⑥ Hands-on part with R
- ⑦ Discussion and conclusions

Motivation

Motivation

Many empirical applications typically have data with $p > n$ or $p \gg n$

- Gene expression
- fMRI
- Chemometrics
- Financial or macroeconomic time series

Two common strategies for model selection:

- Add penalty on coefficient estimates to objective function
→ Certain penalties allow for sparse model estimates
- Sequentially add variables according to their importance

Motivation

Many empirical applications typically have data with $p > n$ or $p \gg n$

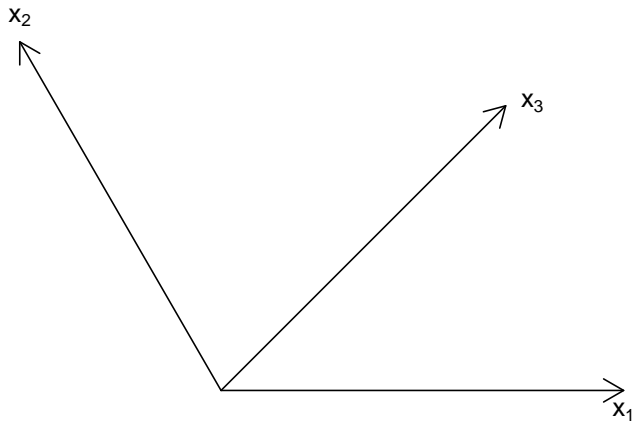
- Gene expression
- fMRI
- Chemometrics
- Financial or macroeconomic time series

Two common strategies for model selection:

- Add **penalty on coefficient estimates** to objective function
→ Certain penalties allow for **sparse model estimates**
- **Sequentially add variables** according to their importance

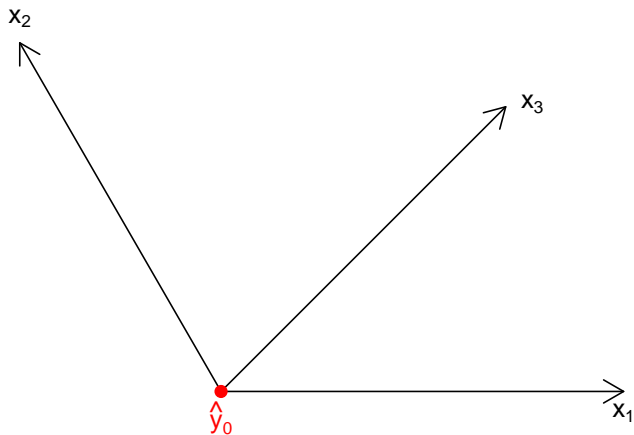
LARS and lasso

Least angle regression (LARS): Idea



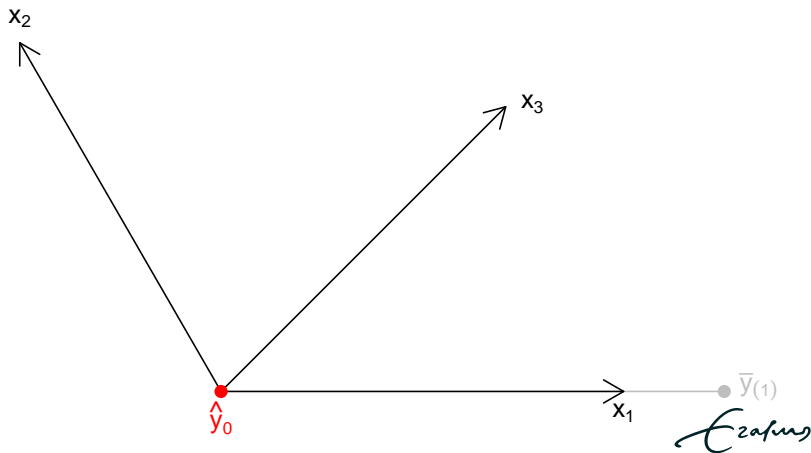
Ezra

Least angle regression (LARS): Idea

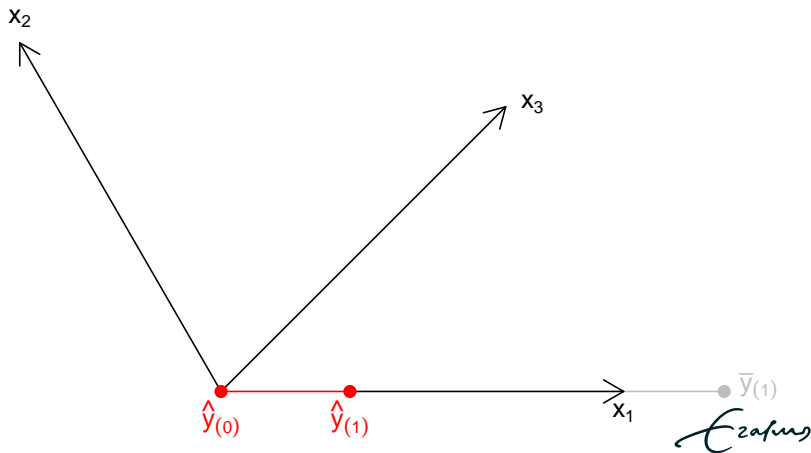


Ezra

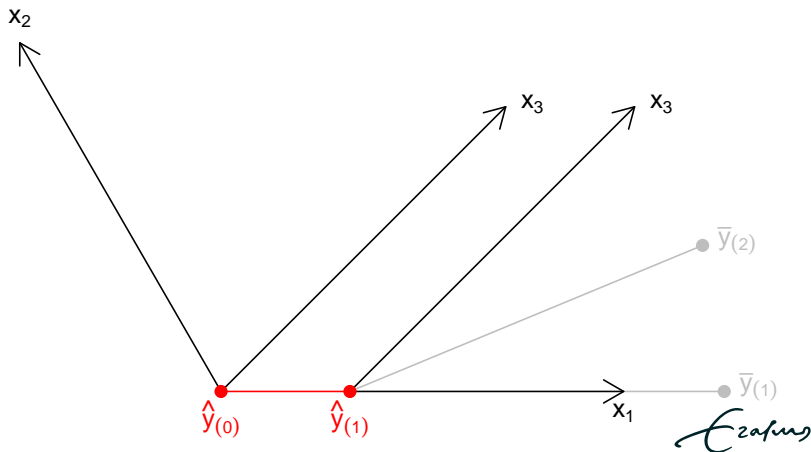
Least angle regression (LARS): Idea



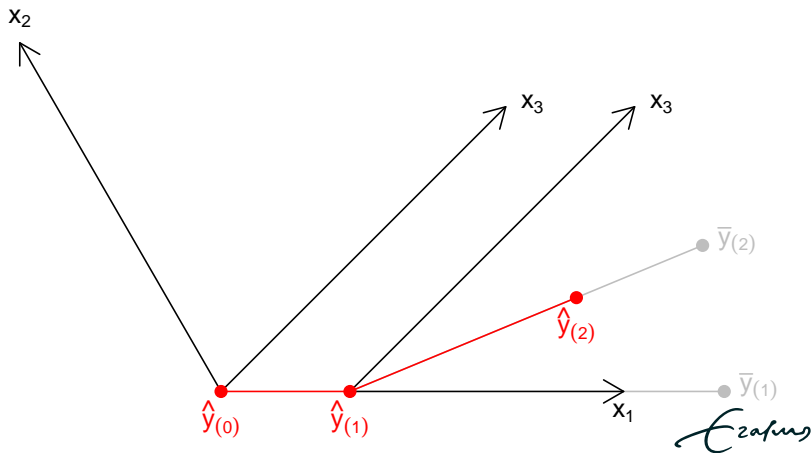
Least angle regression (LARS): Idea



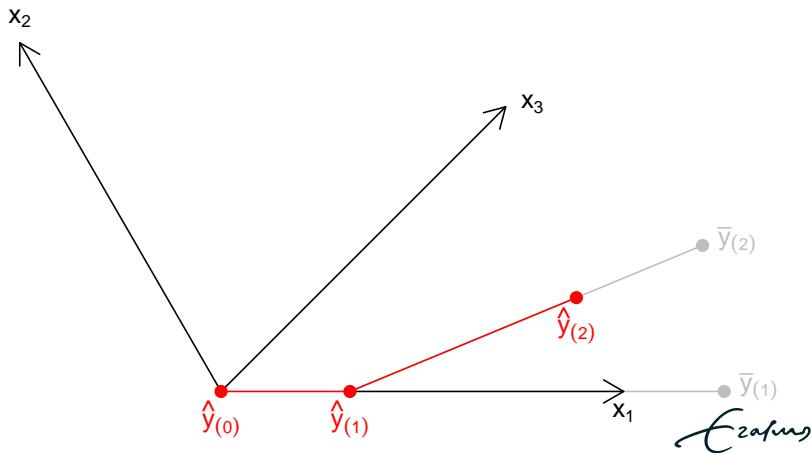
Least angle regression (LARS): Idea



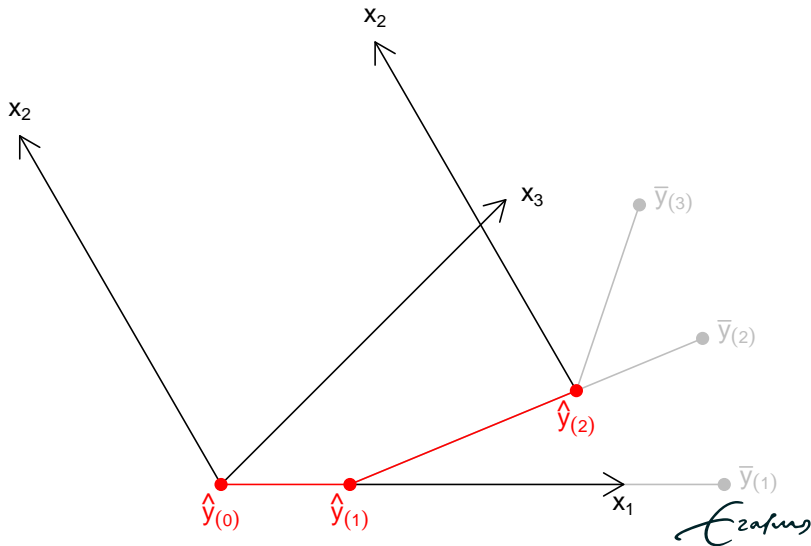
Least angle regression (LARS): Idea



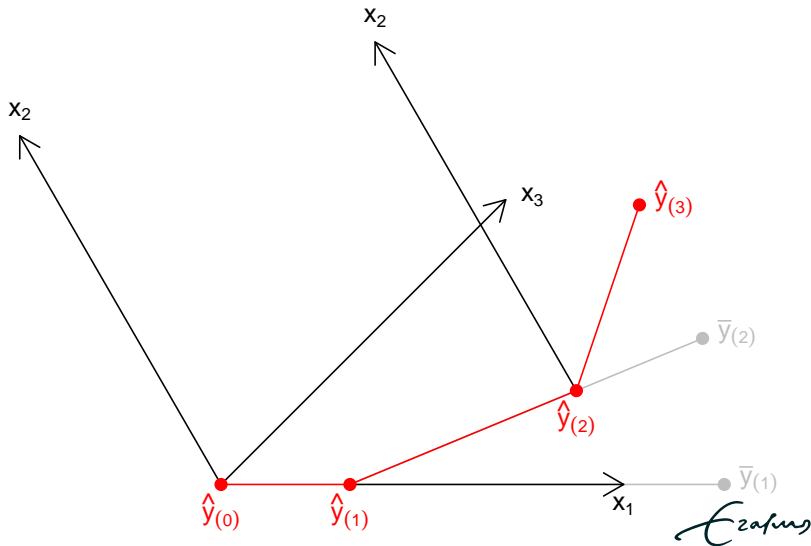
Least angle regression (LARS): Idea



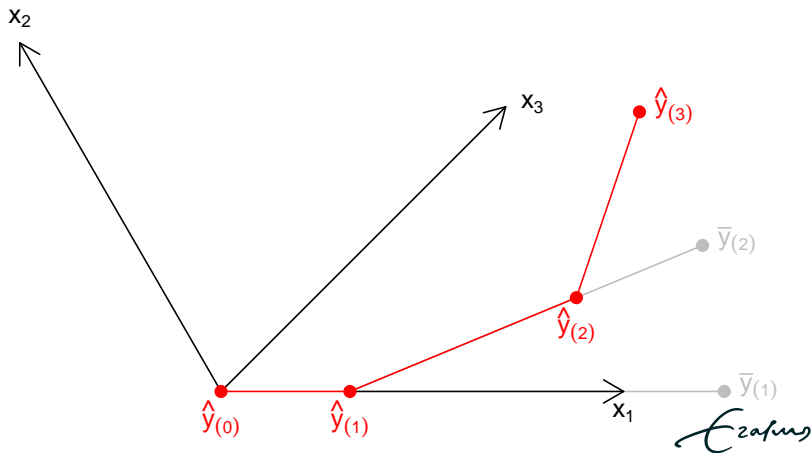
Least angle regression (LARS): Idea



Least angle regression (LARS): Idea



Least angle regression (LARS): Idea



Least angle regression (LARS): Algorithm

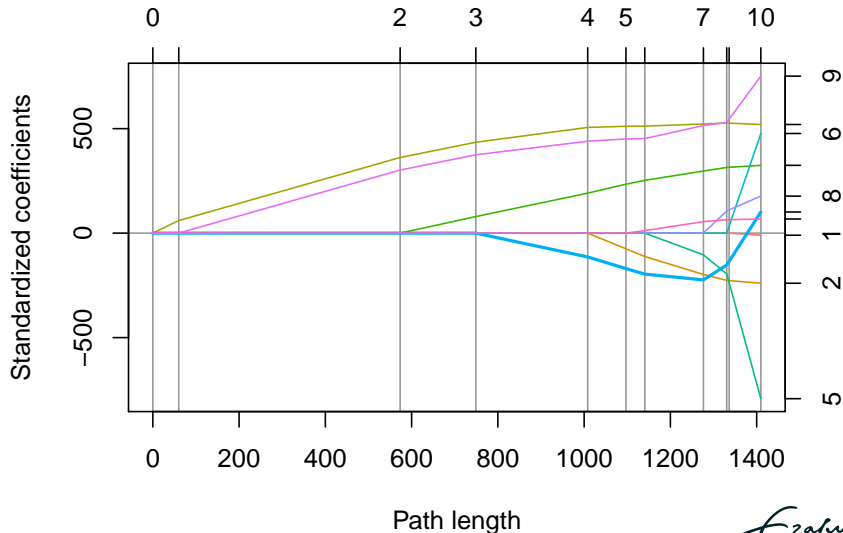
Model selection algorithm based on **forward selection** approach
(Efron et al., 2004)

- Start with **most correlated predictor**
- Move along equi-angular vector **until a new predictor is equally correlated** and add that predictor to the **active set**
- **Update coefficients** of active predictors along solution path

Least angle regression (LARS): Properties

- **Simple formula** for the step size when next predictor is added
- Solution path is **piecewise linear**
 - **Efficient computation**
- Applicable to **high-dimensional data** by limiting the number of steps

LARS: piecewise linear solution path vs path length



Erasmus

Least absolute shrinkage and selection operator (lasso)

Different parametrization than proposed by Tibshirani (1996):

$$\hat{\beta}_{\text{lasso}} = \arg \min_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}'_i \beta)^2 + n\lambda \|\beta\|_1$$

- Can be computed through LARS framework (Efron et al., 2004)
- But modern implementations use a coordinate descent algorithm (Friedman et al., 2010) or an ADMM algorithm (Boyd et al., 2010)

The logo for Erasmus University, featuring the word "Erasmus" in a stylized, handwritten-style script.

Least absolute shrinkage and selection operator (lasso)

Different parametrization than proposed by Tibshirani (1996):

$$\hat{\beta}_{\text{lasso}} = \arg \min_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}'_i \beta)^2 + n\lambda \|\beta\|_1$$

- Can be computed through LARS framework (Efron et al., 2004)
- But modern implementations use a coordinate descent algorithm (Friedman et al., 2010) or an ADMM algorithm (Boyd et al., 2010)

The logo for Erasmus University, featuring a stylized signature of the name 'Erasmus' in a cursive font.

Relationship between LARS and lasso

Modification of the LARS algorithm:

- If the coefficient of an active predictor reaches 0, drop that predictor from the active set
- Continue algorithm with reduced active set

→ Lasso solution path

→ If no coefficient changes signs, LARS solution path is identical to lasso solution path

Relationship between LARS and lasso

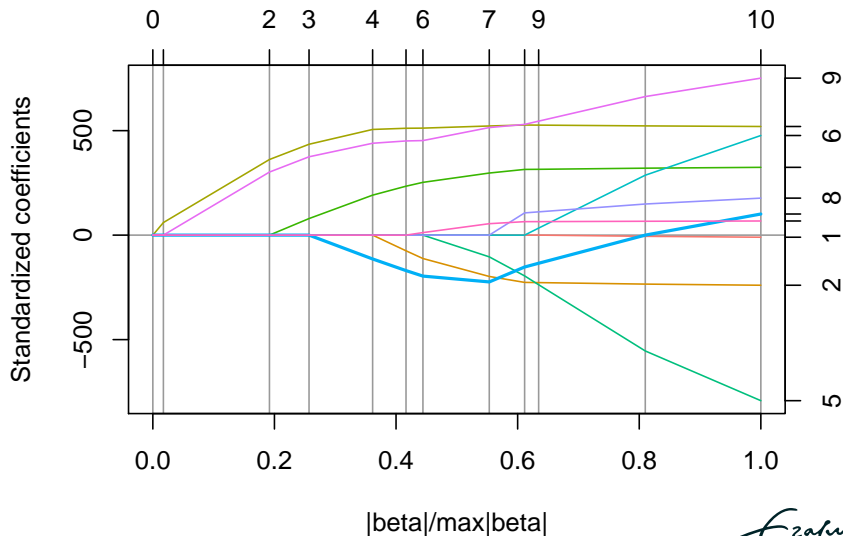
Modification of the LARS algorithm:

- If the coefficient of an active predictor reaches 0, drop that predictor from the active set
- Continue algorithm with reduced active set

→ Lasso solution path

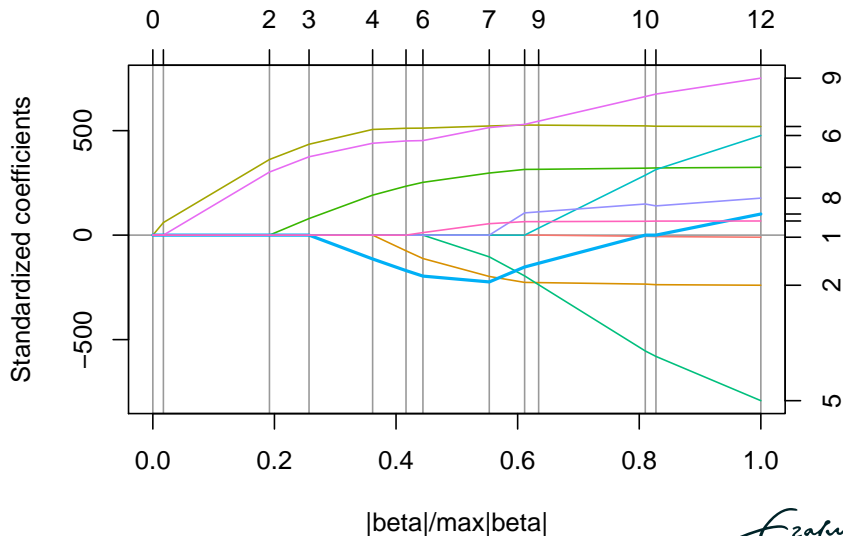
→ If no coefficient changes signs, LARS solution path is identical to lasso solution path

LARS: piecewise linear solution path vs L_1 norm



Ezra

Lasso: piecewise linear solution path vs L_1 norm



Ezra

Robust least angle regression

Robust least angle regression (RLARS)

Hybrid procedure (Khan et al., 2007):

- ① Sequence predictors based on robust correlations
- ② Fit robust regression models along the sequence

→ Applicable to high-dimensional data by limiting the number of steps

→ Implemented in function `rlars()` of R package `robustHD`

The logo for Erasmus University, featuring the word "Erasmus" in a stylized, handwritten-style script.

Robust least angle regression (RLARS)

Hybrid procedure (Khan et al., 2007):

- ① Sequence predictors based on robust correlations
- ② Fit robust regression models along the sequence

→ Applicable to **high-dimensional data** by limiting the number of steps

→ Implemented in function `rlars()` of R package `robustHD`

Robust groupwise least angle regression (RGrpLARS)

Robust extension of LARS to groupwise variable selection
(Alfons et al., 2016):

- ➊ Sequence of predictor groups based on R^2 after initial data cleaning
- ➋ Fit robust regression models along the sequence using original data

→ Applicable to high-dimensional data for some of the proposed data cleaning approaches

→ Implemented in function `rgrplars()` of R package `robustHD`



Robust groupwise least angle regression (RGrpLARS)

Robust extension of LARS to groupwise variable selection

(Alfons et al., 2016):

- ① Sequence of predictor groups based on R^2 after initial data cleaning
- ② Fit robust regression models along the sequence using original data

→ Applicable to **high-dimensional data** for some of the proposed data cleaning approaches

→ Implemented in function `rgrplars()` of R package `robustHD`



Sparse least trimmed squares (trimmed lasso)

Sparse least trimmed squares regression

Objective function:

$$\hat{\beta}_{\text{sparseLTS}} = \arg \min_{\beta} \sum_{i=1}^h (r^2(\beta))_{i:n} + h\lambda \|\beta\|_1$$

with

$$h \leq n$$

$$r^2(\beta) = (r_1^2, \dots, r_n^2)' \quad \text{squared residuals}$$

$$(r^2(\beta))_{1:n} \leq \dots \leq (r^2(\beta))_{n:n} \quad \text{order statistics}$$



Sparse least trimmed squares regression

Combining...

- Least trimmed squares regression for **robustness** (Rousseeuw and Van Driessen, 2006)
- Lasso for **sparsity** (Tibshirani, 1996)

→ C-step algorithm for computation

→ Reweighting step to increase efficiency

→ Details and theory in Alfons et al. (2013) and Öllerer et al. (2015)

The logo for Erasmus University, featuring the word "Erasmus" in a stylized, handwritten script.

Sparse least trimmed squares regression

Combining...

- Least trimmed squares regression for **robustness** (Rousseeuw and Van Driessen, 2006)
- Lasso for **sparsity** (Tibshirani, 1996)

→ **C-step algorithm** for computation

→ **Reweighting step** to increase efficiency

→ Details and theory in Alfons et al. (2013) and Öllerer et al. (2015)

The logo for Erasmus University, featuring the word "Erasmus" in a stylized, handwritten script.

Sparse least trimmed squares regression

Combining...

- Least trimmed squares regression for **robustness** (Rousseeuw and Van Driessen, 2006)
- Lasso for **sparsity** (Tibshirani, 1996)

→ **C-step algorithm** for computation

→ **Reweighting step** to increase efficiency

→ Details and theory in Alfons et al. (2013) and Öllerer et al. (2015)



C-step

Objective function in terms of subset H :

$$Q(H, \beta) = \sum_{i \in H} (y_i - \mathbf{x}'_i \beta)^2 + h\lambda \|\beta\|_1$$

Step k with current subset H_k :

- Obtain lasso solution $\hat{\beta}_{H_k} = \arg \min_{\beta} Q(H_k, \beta)$
- Compute squared residuals $\mathbf{r}_k^2 = (r_{k,1}^2, \dots, r_{k,n}^2)'$
- Construct H_{k+1} from observations with smallest squared residuals:

$$H_{k+1} = \left\{ i \in \{1, \dots, n\} : r_{k,i}^2 \in \{(r_k^2)_{j:n} : j = 1, \dots, h\} \right\}$$



C-step

Objective function in terms of subset H :

$$Q(H, \beta) = \sum_{i \in H} (y_i - \mathbf{x}'_i \beta)^2 + h\lambda \|\beta\|_1$$

Step k with current subset H_k :

- Obtain lasso solution $\hat{\beta}_{H_k} = \arg \min_{\beta} Q(H_k, \beta)$
- Compute squared residuals $\mathbf{r}_k^2 = (r_{k,1}^2, \dots, r_{k,n}^2)'$
- Construct H_{k+1} from observations with smallest squared residuals:

$$H_{k+1} = \left\{ i \in \{1, \dots, n\} : r_{k,i}^2 \in \{(r_k^2)_{j:n} : j = 1, \dots, h\} \right\}$$



Basic C-step algorithm

- ① Obtain m initial subsets of size h from elementary 3-subsets
- ② For $j = 1, \dots, m$ do C-steps until convergence
- ③ Return $\hat{\beta}_{\text{sparseLTS}}$ corresponding to the subset with the lowest value of the objective function

→ Improvements as in FAST-LTS algorithm

→ Implemented in function `sparseLTS()` of R package `robustHD`

Basic C-step algorithm

- ① Obtain m initial subsets of size h from elementary 3-subsets
- ② For $j = 1, \dots, m$ do C-steps until convergence
- ③ Return $\hat{\beta}_{\text{sparseLTS}}$ corresponding to the subset with the lowest value of the objective function

→ Improvements as in FAST-LTS algorithm

→ Implemented in function `sparseLTS()` of R package `robustHD`

Reweighted estimator

Weights from outlier detection via raw estimator:

$$w_i = \begin{cases} 1 & \text{if } |(r_i - \hat{\mu}_{\text{raw}})/\hat{\sigma}_{\text{raw}}| \leq \Phi^{-1}(1 - \delta) \\ 0 & \text{if } |(r_i - \hat{\mu}_{\text{raw}})/\hat{\sigma}_{\text{raw}}| > \Phi^{-1}(1 - \delta) \end{cases} \quad i = 1, \dots, n$$

→ Reweighted estimator given by weighted lasso fit

$$\hat{\beta}_{\text{reweighted}} = \arg \min_{\beta} \sum_{i=1}^n w_i (y_i - \mathbf{x}'_i \beta)^2 + \lambda n_w \|\beta\|_1$$

with

$$n_w = \sum_{i=1}^n w_i \quad \text{number of detected good data points}$$



Reweighted estimator

Weights from outlier detection via raw estimator:

$$w_i = \begin{cases} 1 & \text{if } |(r_i - \hat{\mu}_{\text{raw}})/\hat{\sigma}_{\text{raw}}| \leq \Phi^{-1}(1 - \delta) \\ 0 & \text{if } |(r_i - \hat{\mu}_{\text{raw}})/\hat{\sigma}_{\text{raw}}| > \Phi^{-1}(1 - \delta) \end{cases} \quad i = 1, \dots, n$$

→ Reweighted estimator given by weighted lasso fit

$$\hat{\beta}_{\text{reweighted}} = \arg \min_{\beta} \sum_{i=1}^n w_i (y_i - \mathbf{x}'_i \beta)^2 + \lambda n_w \|\beta\|_1$$

with

$$n_w = \sum_{i=1}^n w_i \quad \text{number of detected good data points}$$



Breakdown point

Finite sample breakdown point (FBP):

$$\epsilon^*(\hat{\beta}; \mathbf{Z}) = \min \left\{ \frac{m}{n} : \sup_{\tilde{\mathbf{Z}}} \|\hat{\beta}(\tilde{\mathbf{Z}})\|_2 = \infty \right\}$$

with

$\mathbf{Z} = (\mathbf{X}, \mathbf{y})$ original sample
 $\tilde{\mathbf{Z}}$ contaminated sample with m points replaced by arbitrary values

Breakdown point theorem

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^h (\rho(\mathbf{y} - \mathbf{X}\beta))_{i:n} + h\lambda \|\beta\|_1$$

where

$$h \leq n$$

$\rho(x)$ convex, symmetric, $\rho(0) = 0$ and $\rho(x) > 0$ for $x \neq 0$

$\rho(\mathbf{y} - \mathbf{X}\beta) := (\rho(y_1 - \mathbf{x}_1\beta), \dots, \rho(y_n - \mathbf{x}_n\beta))'$ losses

$(\rho(\mathbf{y} - \mathbf{X}\beta))_{1:n} \leq \dots \leq (\rho(\mathbf{y} - \mathbf{X}\beta))_{n:n}$ order statistics

→ Breakdown point of the estimator $\hat{\beta}$:

$$\varepsilon^*(\hat{\beta}; \mathbf{Z}) = \frac{n - h + 1}{n}$$



Breakdown point of selected estimators

Sparse LTS:

$$\varepsilon^*(\hat{\beta}_{\text{sparseLTS}}; \mathbf{Z}) = \frac{n - h + 1}{n}$$

Lasso:

$$\varepsilon^*(\hat{\beta}_{\text{lasso}}; \mathbf{Z}) = \frac{1}{n}$$

→ Note: Breakdown point does not depend on dimension p

Breakdown point of selected estimators

Sparse LTS:

$$\varepsilon^*(\hat{\beta}_{\text{sparseLTS}}; \mathbf{Z}) = \frac{n - h + 1}{n}$$

Lasso:

$$\varepsilon^*(\hat{\beta}_{\text{lasso}}; \mathbf{Z}) = \frac{1}{n}$$

→ Note: Breakdown point does **not** depend on dimension p

Questions

- ① For a small enough h (i.e., a large enough trimming proportion), the sparse LTS has a breakdown point larger than 50%.
 - How is this possible?
 - Does this make sense from the perspective of robust statistics?
- ② The lasso is equivalent to the constrained optimization problem

$$\min_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}_i' \beta)^2 \quad \text{subject to} \quad \|\beta\|_1 \leq t$$

We also have equivalence of the L_1 and L_2 norms.

→ How is it possible that the breakdown point of the lasso is 0%?



Penalized S- and MM-estimator

Penalized Elastic Net S-Estimator (PENSE)

Objective function:

$$\hat{\beta}_{\text{PENSE}} = \arg \min_{\beta} \hat{\sigma}^2(\beta) + \lambda_S \left(\alpha \|\beta\|_1 + \frac{1-\alpha}{2} \|\beta\|_2^2 \right)$$

$$\text{with } \frac{1}{n} \sum_{i=1}^n \rho \left(\frac{r_i(\beta)}{\hat{\sigma}(\beta)} \right) = b$$

where

$$r_i(\beta) = y_i - \mathbf{x}_i' \beta$$

residuals

$$b = \mathbb{E}_Z[\rho(Z)] \text{ with } Z \sim \mathcal{N}(0, 1)$$

consistency parameter

Penalized Elastic Net S-Estimator (PENSE)

Combining...

- S-estimator for **robustness** (Salibián-Barrera and Yohai, 2006)
- Elastic net for **regularization** and **sparsity** (Friedman et al., 2005)

→ Iteratively reweighted elastic net (IRWEN) algorithm based on initial estimator for computation

→ Implemented in function `pense()` of R package `pense`

→ Robust, but not efficient

→ Details and theory in Cohen Freue et al. (2019)

The logo for Erasmus University, featuring the word "Erasmus" in a stylized, handwritten script.

Penalized Elastic Net S-Estimator (PENSE)

Combining...

- S-estimator for **robustness** (Salibián-Barrera and Yohai, 2006)
- Elastic net for **regularization** and **sparsity** (Friedman et al., 2005)

→ **Iteratively reweighted elastic net (IRWEN) algorithm** based on initial estimator for computation

→ Implemented in function `pense()` of R package `pense`

→ Robust, but not efficient

→ Details and theory in Cohen Freue et al. (2019)

The logo for Erasmus University, featuring the word "Erasmus" in a stylized, handwritten script font.

Penalized Elastic Net S-Estimator (PENSE)

Combining...

- S-estimator for **robustness** (Salibián-Barrera and Yohai, 2006)
- Elastic net for **regularization** and **sparsity** (Friedman et al., 2005)

→ **Iteratively reweighted elastic net (IRWEN) algorithm** based on initial estimator for computation

→ Implemented in function `pense()` of R package `pense`

→ Robust, but not efficient

→ Details and theory in Cohen Freue et al. (2019)

The logo for Erasmus University, featuring the word "Erasmus" in a stylized, handwritten script font.

PENSE refined via M-estimator (PENSEM)

- To increase efficiency, Cohen Freue et al. (2019) propose a **penalized elastic net M-estimator**, using the initial scale estimate from PENSE
- However, this creates other issues and this approach will not be discussed further

PENSE refined via M-estimator (PENSEM)

- To increase efficiency, Cohen Freue et al. (2019) propose a **penalized elastic net M-estimator**, using the initial scale estimate from PENSE
- However, this creates other issues and this approach will not be discussed further

Hands-on part with R

R packages and script

We will use:

→ R packages `robustHD` (version 0.7.0!) and `pense`

```
R> install.packages(c("robustHD", "pense"))
```

→ R script `ICORS2021_workshop_script.R` available from <https://personal.eur.nl/alfons/ICORS2021.html>

→ Run the commands **in your own R session** along with me



NCI-60 cancer cell panel

- Data on 60 human cancer cell lines
- Available from <http://discover.nci.nih.gov/cellminer/>
- Protein expressions based on 162 antibodies
- Gene expression data with $p = 22\,283$
 - $n = 59$: one observation with all gene expressions missing

→ Use protein expression with largest MAD as response variable

→ Candidate predictors: $d = 100$ most correlated gene expressions

NCI-60 cancer cell panel

- Data on 60 human cancer cell lines
- Available from <http://discover.nci.nih.gov/cellminer/>
- Protein expressions based on 162 antibodies
- Gene expression data with $p = 22\,283$
 - $n = 59$: one observation with all gene expressions missing

→ Use protein expression with largest MAD as response variable

→ Candidate predictors: $d = 100$ most correlated gene expressions

Discussion and conclusions

Some issues to look out for

- Residual scale is typically **underestimated** in high-dimensions
 - Outlier detection via standardized residuals is prone to **false positives**

- BIC for regularization parameter selection can be **unstable** for values of λ close to 0 due to exact fit situations
 - Cross-validation is preferred, but computationally expensive

Conclusions

- Robust regression in high dimensions remains a challenging problem
- R packages `robustHD` and `pense` provide promising functionality
- A trimmed version of the elastic net (Kurnaz et al., 2017) is available in R package `enetLTS`, also for logistic regression

References I

- Alfons, A., Croux, C., and Gelper, S. (2013). Sparse least trimmed squares regression for analyzing high-dimensional large data sets. The Annals of Applied Statistics, 7(1):226–248.
- Alfons, A., Croux, C., and Gelper, S. (2016). Robust groupwise least angle regression. Computational Statistics & Data Analysis, 93:421–435.
- Boyd, S., Parikh, N., Chu, E., Peleato, B., and Eckstein, J. (2010). Distributed optimization and statistical learning via the alternating direction method of multipliers. Foundations and Trends in Machine Learning, 3(1):1–122.
- Cohen Freue, G., Kepplinger, D., Salibian-Barrera, M., and Smucler, E. (2019). Robust elastic net estimators for variable selection and identification of proteomic biomarkers. The Annals of Applied Statistics, 13(4):2065–2090.
- Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. (2004). Least angle regression. The Annals of Statistics, 32(2):407–499.
- Friedman, J., Hastie, T., and Tibshirani, R. (2005). Regularization and variable selection via the elastic net. Journal of the Royal Statistical Society, Series B, 67(2):301–320.



References II

- Friedman, J., Hastie, T., and Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. Journal of Statistical Software, 33(1):1–22.
- Khan, J., Van Aelst, S., and Zamar, R. (2007). Robust linear model selection based on least angle regression. Journal of the American Statistical Association, 102(480):1289–1299.
- Kurnaz, F., Hoffmann, I., and Filzmoser, P. (2017). Robust and sparse estimation methods for high dimensional linear and logistic regression. Chemometrics and Intelligent Laboratory Systems, 172:211–222.
- Öllerer, V., Croux, C., and Alfons, A. (2015). The influence function of penalized regression estimators. Statistics: A Journal of Theoretical and Applied Statistics, 49(4):741–765.
- Rousseeuw, P. and Van Driessen, K. (2006). Computing LTS regression for large data sets. Data Mining and Knowledge Discovery, 12(1):29–45.
- Salibian-Barrera, M. and Yohai, V. (2006). A fast algorithm for S-regression estimates. Journal of Computational and Graphical Statistics, 15(2):414–427.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. Journal of the Royal Statistical Society, Series B, 58(1):267–288.

