

# Robustification of Statistical and Econometrical Regression Methods

Tomáš Jurczyk

Charles University in Prague,  
Faculty of Mathematics and Physics

Department of Probability and Mathematical Statistics



ROBUST 2016

13. 9. 2016

# Motivation (goal)

To push method in wide usage, it is not sufficient to have good method. You need to have the whole package of methods, diagnostics and tests around it.

We want to build this "package" around **LWS** method (which is method for robust regression).

## Goals:

- We want to find variant of this method in case of multicollinearity
- Investigate and understand the situation in data with contamination together multicollinearity

# Regression task

## Notation:

Linear regression model:

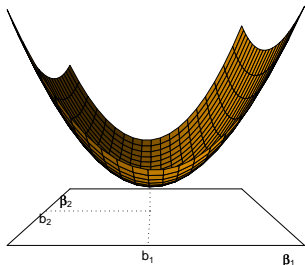
$$Y_i = \sum_{j=1}^p x_{ij} \beta_j^0 + e_i, \quad i = 1, 2, \dots, n.$$

## Classical method:

### Least Squares (LS)

$$\mathbf{b} = \underset{\beta \in \mathcal{R}^p}{\operatorname{argmin}} \sum_{i=1}^n r_i^2(\beta).$$

*Example* (Loss function of the LS estimate, 2 regressors)



- ←  $\sum_{i=1}^n (Y_i - x_{1i}\beta_1 - x_{2i}\beta_2)^2 = \sum_{i=1}^n r_i^2(\beta)$
- ← Loss function of LS is quadratic
- ←  $(b_1, b_2)'$  =  $\mathbf{b}$  denote LS estimate – the minimal value of loss function

# Assumptions

LS estimator is simple and widely used

It is also due to nice properties of this estimator.

Under following assumptions LS is BLUE:

- 1)  $e_i, i = 1, \dots, n$  are independent
- 2)  $Ee_i = 0$  for all  $i$ , which means  $E\mathbf{Y} = \mathbf{X}'_i\beta^0$
- 3) The rank of the matrix  $\mathbf{X}$  is full.
- 4) Variance  $\text{var}e_i = \sigma^2, i = 1, \dots, n$ .
- 5)  $e_i, i = 1, \dots, n$  has normal distribution.

# Assumptions

LS estimator is simple and widely used

It is also due to nice properties of this estimator.

Under following assumptions LS is BLUE:

- 1)  $e_i, i = 1, \dots, n$  are independent
- 2)  $Ee_i = 0$  for all  $i$ , which means  $E\mathbf{Y} = \mathbf{X}'_i\beta^0$
- 3) The rank of the matrix  $\mathbf{X}$  is full.
- 4) Variance  $\text{var}e_i = \sigma^2, i = 1, \dots, n$ .
- 5)  $e_i, i = 1, \dots, n$  has normal distribution.

Typically not all assumptions are fulfilled.

Here are possible data problems:

- ▶ **Observation are not independent**
- ▶ **Heteroscedasticity** (different variances of components of error term)
- ▶ **Multicollinearity** (problem with dependence of regressors)
- ▶ **Not normal distribution** of error term
- ▶ Presence of **outlying observations**

# Problems of data - Multicollinearity

## Multicollinearity

- ▶ situation when regressors are "nearly" linear dependend

# Problems of data - Multicollinearity

## Multicollinearity

- ▶ situation when regressors are “nearly” linear dependend

Consequences for the least squares method (LS)

- ▶ Matrix  $\mathbf{X}'\mathbf{X}$  is almost singular
- ▶ The smallest eigenvalue  $t_p^2$  of the matrix  $\mathbf{X}'\mathbf{X}$  is close to 0.
- ▶ Numerical solution of normal equation is not stable.
- ▶ Multicollinearity induces large expected value of the length of the LS estimate ( $\mathbf{b}$ ).

$$E \|\mathbf{b}\|^2 - \|\beta^0\|^2 = E \|\mathbf{b} - \beta^0\|^2 = \text{tr}(\text{var}(\mathbf{b})) = \sigma^2 \text{tr}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2 \sum_{i=1}^p (1/t_i)^2$$

- ▶ It may cause large variance of  $b_j$ .
- $$\text{var} \mathbf{b} = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} = \sigma^2 \sum_{i=1}^p t_i^{-2} \mathbf{q}_i \mathbf{q}_i'$$

where:  $\mathbf{X} = \mathbf{P}\mathbf{T}\mathbf{Q}'$ ,  $\mathbf{P}'\mathbf{P} = \mathbf{Q}\mathbf{Q}' = \mathbf{I}$ ,  $\mathbf{T} = \text{diag}(t_1, t_2, \dots, t_p)$ ,  $t_i^2$  eigen value of  $\mathbf{X}'\mathbf{X}$

# Problems of data - Multicollinearity

## Multicollinearity

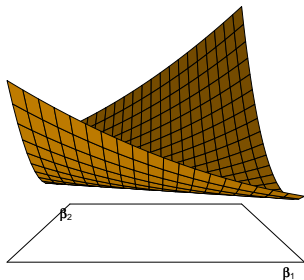
- ▶ situation when regressors are “nearly” linear dependent

Consequences for the least squares method (LS)

- ▶ Matrix  $\mathbf{X}'\mathbf{X}$  is almost singular
- ▶ The smallest eigenvalue  $t_p^2$  of the matrix  $\mathbf{X}'\mathbf{X}$  is close to 0.
- ▶ Numerical solution of normal equation is not stable.
- ▶ Multicollinearity induces large expected value of the length of the LS estimate ( $\mathbf{b}$ ).
- ▶ It may cause large variance of  $b_j$ .

**Example** (Loss function of the LS estimate, 2 regressors, data with multicollinearity)

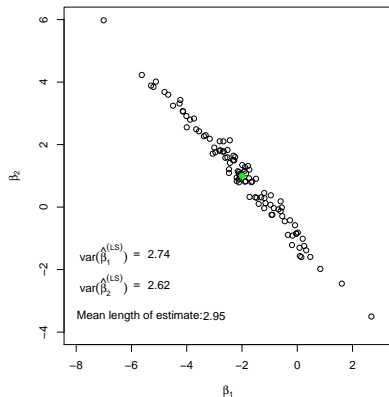
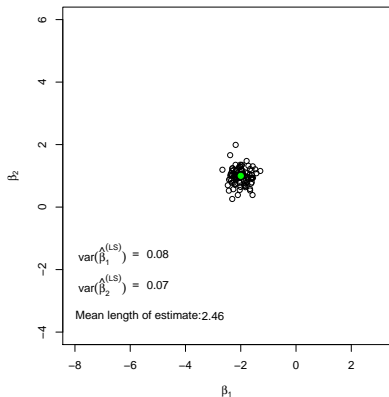
- Typical shape for data with strong multicollinearity
- There are almost the same values of loss function along some line in parameter space  $\Rightarrow$  unstable behaviour  $\Rightarrow$  large variance and large expected length of  $\mathbf{b}$





# Demonstration of multicollinearity for LS

LS estimate on independent regressors (left graph) and LS estimate for multicollinear regressors (right graph). Each point is LS estimate for one run of simulated dataset (correlation of regressors in right graph is around 0.99). Green dot is theoretical value of  $\beta^0$ .



# Ridge regression

We show estimate used instead of the least squares when multicollinearity is present.

## Ridge Regression (RR)

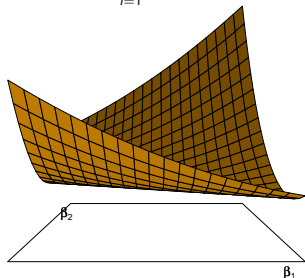
$$\mathbf{b}_\delta = \underset{\boldsymbol{\beta} \in \mathcal{R}^p}{\operatorname{argmin}} \left( \sum_{i=1}^n r_i^2(\boldsymbol{\beta}) + \delta \sum_{j=1}^p \beta_j^2 \right).$$

*Example*

(2 regressors, data with multicollinearity)

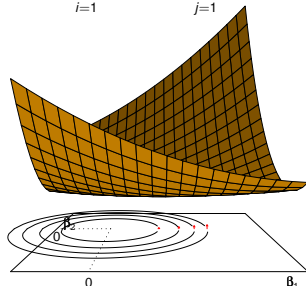
Loss function of the LS estimate

$$\sum_{i=1}^n r_i^2(\boldsymbol{\beta})$$



Loss function of the ridge estimate

$$\sum_{i=1}^n r_i^2(\boldsymbol{\beta}) + \delta \sum_{j=1}^p \beta_j^2$$



# Problem of data - contamination

## Basic goal of robust statistics

- ▶ Finding models which correspond to the structure of the majority of the data.
- ▶ The outlier detection is closely connected with this objective.

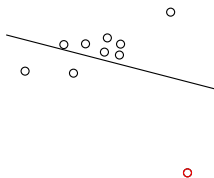
Outlier in our context - **observation which does not follow the regression model.**

## Problems with outlier presence

- ▶ **Already one outlier is able to change the value of LS estimate essentially.**

---

**Example** (Influence of one outlier to LS estimate)



# Approaches to overcome contamination

Reduction of influence of contaminating points:

- ▶ Different loss function (M-estimates, Least Absolute Deviation regression,...)
- ▶ Implicit residual weighting (LTS, LMS, LWS)

## Representative of robust methods

### Least Trimmed Squares (LTS)

Let  $n/2 \leq h \leq n$ . Then

$$\mathbf{b}^{(LTS,n,h)} = \underset{\beta \in \mathcal{R}^p}{\operatorname{argmin}} \sum_{i=1}^h r_{(i)}^2(\beta)$$

is called *Least Trimmed Squares Estimator*,  $r_{(j)}^2(\beta)$  is the  $j$ -th order statistic among the squared residuals.

Generalization of LTS:

### Least Weighted Squares (LWS)

Let  $1 = w_1 \geq w_2 \geq \dots \geq w_n \geq 0$  are weights. Then

$$\mathbf{b}^{(LWS,w)} = \underset{\beta \in \mathcal{R}^p}{\operatorname{argmin}} \sum_{i=1}^n w_i r_{(i)}^2(\beta)$$

is called *Least Weighted Squares estimate*.

# Multicollinearity in company with outliers

We need to stress that **contamination and multicollinearity are different in their essence**. Multicollinearity is a problem of regressors, whereas the presence of outliers is the problem with non-compliance of the regression model.

This induces lot of problematic combined situations.

Known problems with outliers in connection with multicollinearity

- ▶ Outlier can affect value of estimate of classical methods (like RR).
- ▶ In addition already one outlier can **hide or create multicollinearity** for classical methods for multicollinearity detection.
- ▶ Therefore we need some robust detector of multicollinearity.

# Multicollinearity in company with outliers

We need to stress that **contamination and multicollinearity are different in their essence**. Multicollinearity is a problem of regressors, whereas the presence of outliers is the problem with non-compliance of the regression model.

This induces lot of problematic combined situations.

Known problems with outliers in connection with multicollinearity

- ▶ Outlier can affect value of estimate of classical methods (like RR).
- ▶ In addition already one outlier can **hide or create multicollinearity** for classical methods for multicollinearity detection.
- ▶ Therefore we need some robust detector of multicollinearity.

---

First idea and approach also presented in literature is to use high breakdown methods for outlier detection and after revelation of outliers use some classical multicollinearity diagnostics on non-contaminated data.

*Will such approach work?*

---

# Multicollinearity in company with outliers

We need to stress that **contamination and multicollinearity are different in their essence**. Multicollinearity is a problem of regressors, whereas the presence of outliers is the problem with non-compliance of the regression model.

This induces lot of problematic combined situations.

Known problems with outliers in connection with multicollinearity

- ▶ Outlier can affect value of estimate of classical methods (like RR).
- ▶ In addition already one outlier can **hide or create multicollinearity** for classical methods for multicollinearity detection.
- ▶ Therefore we need some robust detector of multicollinearity.

---

First idea and approach also presented in literature is to use high breakdown methods for outlier detection and after revelation of outliers use some classical multicollinearity diagnostics on non-contaminated data.

*Will such approach work?*

---

You will see in the poster...

# Poster content

Outline what you will see in the poster:

- ▶ Simple examples for understanding multicollinearity and outliers separately
- ▶ Investigation of functionality of recent proposals for regression methods suitable for combined outlier-multicollinearity problem – some important results have been done in this area
- ▶ Proposal of new regression method called Ridge Least Weighted Squares
- ▶ Properties of this estimator
- ▶ Using new estimate for diagnostics



# Strategy

*Least  
squares*



*Ridge  
regression*



*Least  
weighted  
squares*

# Strategy

*Least  
squares*



*Ridge  
regression*



*Least  
weighted  
squares*



*Robust  
ridge  
regression*

# Strategy

$$\operatorname{argmin}_{\beta \in \mathcal{R}^p} \sum_{i=1}^n r_i^2(\beta)$$

→

$$\operatorname{argmin}_{\beta \in \mathcal{R}^p} \left( \sum_{i=1}^n r_i^2(\beta) + \delta \sum_{j=1}^p \beta_j^2 \right)$$

↓

↓

$$\operatorname{argmin}_{\beta \in \mathcal{R}^p} \sum_{i=1}^n w_i r_{(i)}^2(\beta)$$

→


$$\operatorname{argmin}_{\beta \in \mathcal{R}^p} \left( \sum_{i=1}^n w_i r_{(i)}^2(\beta) + \delta \sum_{j=1}^p \beta_j^2 \right)$$


# How to find poster

Follow the same color as this presentation...

## Robustification of Statistical and Econometrical Regression Methods (Problems of Combined Outlier-Multicollinearity Presence)

Tomáš Jurczyk  
Faculty of Mathematics and Physics, CHARLES UNIVERSITY IN PRAGUE, Department of Probability and Mathematical Statistics,  
Sokolovská 83, 18675 Prague 8, Czech Republic





### Abstract

This poster is an illustration of problems caused by multicollinearity and outlier presence in the data. Through simple examples we can see behavior of classical least squares method and also robust least trimmed squares (LTS) method in different situations. The most interesting in multicollinearity of LTS (and other methods based on implicit residual weighting) in revealing outliers in the situation where majority of the data outliers from multicollinearity. Method ridge least trimmed squares (RLTS) is presented as a remedy. Apart from showing that this method could be a robust multicollinearity detector, derived properties and diagnostic plots are briefly recalled.

### Notation

We consider the linear regression model

$$Y_i = \sum_{j=1}^p x_{ij}\beta_j + \varepsilon_i, \quad i = 1, 2, \dots, n.$$

$\varepsilon_i(\beta) = Y_i - \sum_{j=1}^p x_{ij}\beta_j$  denotes the  $i$ -th residual and  $r_{(j)}^2(\beta)$  the  $j$ -th order statistic among the squared residuals.

### Classical methods

#### Definition of Least Squares

We define the least squares estimate of the vector  $\beta$  as

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n \varepsilon_i^2(\beta).$$

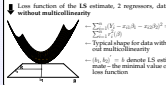
**Loss function**

Loss function of the LS estimate, 2 regressors, data without multicollinearity

$$\sum_{i=1}^n (Y_i - a_{1i}\hat{b}_1 - a_{2i}\hat{b}_2)^2$$

Typical shape for data without multicollinearity

$(\hat{b}_1, \hat{b}_2) = \hat{b}$  denote LS estimate – the minimal value of loss function



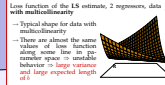
### Multicollinearity

**Multicollinearity**  
situation when regressors are "nearly" linear dependent  
 $\Rightarrow X'X$  is nearly singular with some eigen values close to 0  
 $\Rightarrow$  Large variance and expected length of estimate

**Loss function of the LS estimate, 2 regressors, data with multicollinearity**

Typical shape for data with multicollinearity

There are almost the same values of loss function along some line in parameter space  $\Rightarrow$  unstable behavior  $\Rightarrow$  large variance and large expected lengths of  $\hat{b}$ .



### Dealing with multicollinearity

#### Definition of Ridge Regression

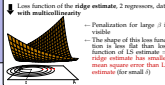
For  $\delta \geq 0$  we define the ridge estimate of the vector  $\beta$  as

$$\hat{b}_\delta = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left( \sum_{i=1}^n \varepsilon_i^2(\beta) + \delta \sum_{j=1}^p \beta_j^2 \right).$$

**Loss function of the ridge estimate, 2 regressors, data with multicollinearity**

Penalization for large  $\beta$  is visible

The shape of this loss function is less flat than loss function of LS estimate  $\Rightarrow$  ridge estimate has smaller mean square error than LS estimate (for small  $\delta$ )



### Outliers

**Basic goal of robust statistics**

- Finding models which correspond to the structure of the majority of the data, therefore also solving outlier presence

**Problems with outlier presence**

- Already one outlier is able to change the value of LS estimate essentially
- Outlier in data point which does not follow regression model. Outlier far away from the model will move the minimum of the LS loss function in direction of its influence because all residuals have the same importance



$\sum_{i=1}^n \varepsilon_i^2(\beta)$

$\sum_{i=1}^n \varepsilon_i^2(\beta) + \delta \sum_{j=1}^p \beta_j^2$

$\sum_{i=1}^k \varepsilon_i^2(\beta)$

$\sum_{i=1}^k \varepsilon_i^2(\beta) + \delta \sum_{j=1}^p \beta_j^2$

**Problems with outlier presence**

- The same story as for LS estimate because also for ridge regression all residuals have the same importance
- Ridge estimate is not robust

### Dealing with outliers

#### Definition of Least Trimmed Squares

Let  $n/2 \leq k \leq n$ , then the least trimmed squares estimate is defined as

$$\hat{b}_{LTS} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^k \varepsilon_i^2(\beta).$$

Compared to LS estimate, the values of the  $n - k$  largest residuals do not affect the LTS estimate  $\Rightarrow$  observe them, which will, keep the  $n - k$  lowest residuals

**Problems with outliers in connection with multicollinearity**

- Already one outlier can hide or create multicollinearity for classical methods for multicollinearity detection.
- Therefore we need some robust detector of multicollinearity

Now, we consider following type of the data:

- The majority of the data suffers from multicollinearity and follows the regression model. Rest of the data are outliers – contamination
- From a good robust method we expect correct detection of such outliers and therefore revelation of the

### Dealing with both problems






#### Definition of Ridge LTS

Let  $n/2 \leq k \leq n$  and  $\delta \geq 0$ , then the ridge least trimmed squares estimate is defined as

$$\hat{b}_{RLTS} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left( \sum_{i=1}^k \varepsilon_i^2(\beta) + \delta \sum_{j=1}^p \beta_j^2 \right).$$

In the middle of the poster we can see (from scheme with loss functions) that method is logical combination of both

# My articles concerning this work

- [1]  T. Jurczyk, High Breakdown Point in Regression, in *WDS'08 Proceedings of Contributed Papers: Part I - Mathematics and Computer Sciences* (eds. J. Safrankova and J. Pavlu), Prague, Matfyzpress, pp. 94–99, 2008.
- [2]  T. Jurczyk, Ridge least weighted squares, *Acta Universitatis Carolinae, Mathematica et Physica*, 52, 1, 15-26, 2011.
- [3]  T. Jurczyk, Weak consistency and weak  $\sqrt{n}$ -consistency of ridge least weighted squares, *Proceedings of the 14th Applied Stochastic Models and Data Analysis (AMSDA 2011) Conference, Faculty of Economics of the University of Rome, Rome, Italy*, pp. 635–643, 2011.
- [4]  T. Jurczyk, Trimmed Estimators in Regression Framework, *Acta Univ. Palacki. Olomuc., Fac.rer. nat., Mathematica 50, 2*, 45–53, 2011.
- [5]  T. Jurczyk, Outlier detection under multicollinearity, *Journal of Statistical Computation and Simulation vol. 82, no. 2*, pp. 261–278, 2012.