

Probabilistic algorithms for computing the LTS estimate.



**FACULTY
OF INFORMATION
TECHNOLOGY
CTU IN PRAGUE**

Bachelor's thesis

Probabilistic algorithms for computing the LTS estimate

Martin Jenč

Department of Applied Mathematics
Supervisor: Ing. Karel Klouda, Ph.D.

May 13, 2019

Acknowledgements

THANKS to everybody

Declaration

I hereby declare that the presented thesis is my own work and that I have cited all sources of information in accordance with the Guideline for adhering to ethical principles when elaborating an academic final thesis.

I acknowledge that my thesis is subject to the rights and obligations stipulated by the Act No. 121/2000 Coll., the Copyright Act, as amended, in particular that the Czech Technical University in Prague has the right to conclude a license agreement on the utilization of this thesis as school work under the provisions of Article 60(1) of the Act.

In Prague on May 13, 2019

.....

Czech Technical University in Prague

Faculty of Information Technology

© 2019 Martin Jenč. All rights reserved.

This thesis is school work as defined by Copyright Act of the Czech Republic. It has been submitted at Czech Technical University in Prague, Faculty of Information Technology. The thesis is protected by the Copyright Act and its usage without author's permission is prohibited (with exceptions defined by the Copyright Act).

Citation of this thesis

Jenč, Martin. *Probabilistic algorithms for computing the LTS estimate*. Bachelor's thesis. Czech Technical University in Prague, Faculty of Information Technology, 2019.

Abstrakt

Metoda nejmenších usekaných čtverců je robustní verzí známé metody nejmenších čtverců, jednou ze základních metod regresní analýzy, používané k odhadování koeficientů lineárního regresního modelu. Výpočet odhadu pomocí metody nejmenších usekaných čtverců je znám jako NP-težký a proto jsou v praxi nejčastěji používány pouze suboptimální pravděpodobnostní algoritmy. Mimo popisu těchto algoritmů navrhujeme několik způsobů jak je zkombinovat za účelem dosažení lepších výsledků.

Klíčová slova nejmenší usekané čtverce, LTS odhad, lineární regrese, robustní statistika, nejmenší čtverce, chybná pozorování

Abstract

The least trimmed squares method is a robust version of the method of least squares which is an essential tool of regression analysis used to find an estimate of coefficients in the linear regression model. Computing the least trimmed squared estimate is known to be NP-hard, and hence only suboptimal probabilistic algorithms are usually used in practice. Besides describing those algorithms, we propose a few ways of combining those algorithms to obtain better performance.

Keywords least trimmed squares, LTS estimate , linear regression, robust statistics, ordinary least squares, outliers

Contents

Citation of this thesis	vi
Introduction	1
1 Least trimmed squares	3
1.1 Linear regression model	3
1.1.1 Prediction with the linear regression model	4
1.2 Ordinary least squares	5
1.2.1 Properties of the OLS estimate	7
1.3 Robust statistics	7
1.3.1 Outliers	7
1.3.2 Measuring robustness	8
1.4 Least trimmed squares	9
1.4.1 Discrete objective function	9
2 Algorithms	13
2.0.1 First attempts	13
2.0.2 Strong and weak necessary conditions	14
2.1 Computing OLS	15
2.1.1 Computation using a matrix inversion	15
2.1.2 Computation using the Cholesky decomposition	16
2.1.3 Computation using the QR decomposition	17
2.2 FAST-LTS	22
2.2.1 C-step	22
2.2.2 Choosing an initial \mathbf{m}_1 subset	26
2.2.3 Speed-up of the algorithm	28
2.2.4 Putting all together	29
2.3 Feasible solution	29
2.4 OEA, MOEA, MMEA	33
2.4.1 Multiplicative formula	33

2.4.2	The OEA and its properties	37
2.4.3	Minimum-maximum exchange algorithm	39
2.4.4	Different method of computation	40
2.5	Combined algorithm	45
2.6	BAB algorithm	45
2.6.1	Improvements	47
2.7	BSA algorithm	47
2.7.0.1	Domain of OF-LTS	48
2.7.1	One-dimensional version of the algorithm	50
2.7.2	Multidimensional BSA	51
3	Experiments	53
3.1	Data set generator	53
3.1.1	Generating outliers	54
3.2	Data sets	55
3.3	Implementation of the algorithms	56
3.4	Results	57
	Conclusion	59
	Bibliography	61
	A Data sets	65
	B Contents of enclosed CD	67

List of Figures

1.1	Change of the regression hyperplane given by coefficients estimated with OLS method when one of the four observations (highlighted with red color) starts to deviate from the linear pattern.	9
2.1	Illustration of the C-step algorithm. (1) represents the value of $OF^{(OLS, \mathbf{M}_1 \mathbf{X}, \mathbf{M}_1 \mathbf{y})}(\hat{\mathbf{w}}_1)$ (which is equal to the value $OF_D^{LTS}(\mathbf{m}_1)$). (2) represents the value of $OF^{(OLS, \mathbf{M}_2 \mathbf{X}, \mathbf{M}_2 \mathbf{y})}(\hat{\mathbf{w}}_1)$ and (3) represents the value of $OF^{(OLS, \mathbf{M}_2 \mathbf{X}, \mathbf{M}_2 \mathbf{y})}(\hat{\mathbf{w}}_2)$ (which is equal to the value $OF_D^{LTS}(\mathbf{m}_2)$).	24
2.2	Value of the residual sum of squares (normalized) based on the number of the step of C-step algorithm for 100 different starting subsets. Dataset D3 was used with configuration $n = 500, p = 20$ and 30% of the the outliers (see Section 3 for more details about the dataset).	26
2.3	Tree consisting of all element subsets for $n = 4$ and $h = 3$. Leaves with blue border color represents h -element subsets.	46
2.4	Step in the algorithms when due to the sorted siblings node j_k can be trimmed	48
3.1	Different types of outliers.	54

List of Tables

3.1	Table to test captions and labels	58
-----	---	----

Introduction

Least trimmed squares (LTS) is one of many modifications very well known method ordinary least squares (OLS). Both of these methods are tools of regression analysis, which is a group of processes used to estimate the dependence of variables. Specifically in the case when we try to estimate dependence one variable on one or multiple others. The regression analysis uses the regression models, and in the case of OLS and LTS methods, that model is the linear regression model.

To OLS estimate produces reliable results, many strong assumptions about the data have to be fulfilled. It is assumed that data is generated specific way as well as that the data does not contain measurement errors called outliers. Those assumptions are in practice hardly fulfilled, because outliers in the data are very common. The OLS estimate is in such cases unreliable.

Problems of classical statistic methods try to solve robust statistic. Its methods are usually emulations of classical statistic methods but try to provide good performance even if data contains a large number of outliers. That means such methods does not rely so much on assumptions which are difficult to achieve in practice. Because the OLS method is one of the essential tools of linear regression analysis, multiple its alternatives have been designed to fulfill the assumptions of a robust estimator.

One of the methods is LTS. Its idea is simple, but unlike the OLS method, the exact formula for calculation is not known, so that only to date known exact algorithms for computing the LTS estimate are combinatorial. For this reason beside exact algorithms, several probabilistic algorithms have been proposed.

This work is divided into three chapters. In the first chapter, we introduce theory required to understand linear regression model, OLS and LTS. We mention properties of both methods and also describe field of its usage.

In the second chapter, we cover algorithms for calculating the OLS estimate. It is necessary because most of the algorithm used to compute the LTS estimate relies on those algorithms. Next, we describe all currently used

algorithms for calculating LTS estimate. They consist of several probabilistic and also few exact ones. In this chapter, we also show that those algorithms can be easily combined to obtain higher speed and performance.

In the last chapter, we describe our experimental results. At first, we cover data generator which is used for our experiments and which can provide data affected by various types of outliers. Next, we provide information about our implementation of all algorithms for chapter two. Finally, we present our results for specific data.

Least trimmed squares

In this chapter, we introduce one of the most common regression analysis models which is known as the linear regression model. It aims to model the relationship between one variable which is called *dependent* and one or more variables which are called *explanatory*. The relationship is based on a model function with parameters which are not known in advance and are to be estimated from data. We also describe one of the most common methods for finding those parameters in this model, namely the ordinary least squares method. It is important to note that all vectors in this text are considered as column vectors. On the other hand, we denote a row vector as a transposed vector.

1.1 Linear regression model

Definition 1. The *linear regression model* is given by

$$y = \mathbf{x}^T \mathbf{w} + \varepsilon, \quad (1.1)$$

where $y \in \mathbb{R}$ is a random variable which is called *dependent variable* and $\mathbf{x}^T = (x_1, x_2, \dots, x_p)$ is a vector of *explanatory variables*. Usually we call x_i a regressor. Finally $\varepsilon \in \mathbb{R}$ is a random variable called *noise* or *error*. The vector $\mathbf{w} = (w_1, w_2, \dots, w_p)$ is a vector of parameters called *regression coefficients*.

In regression analysis we aim to estimate the \mathbf{w} using n measurements of y and \mathbf{x} . We can write this in matrix form

$$\mathbf{y} = \mathbf{X} \mathbf{w} + \boldsymbol{\varepsilon} \quad (1.2)$$

where

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1p} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & x_{n3} & \dots & x_{np} \end{bmatrix}, \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix}, \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

This means that we can think of rows of matrix \mathbf{X} as columns vectors \mathbf{x}_i written into the row.

It is assumed that errors are independent and identically distributed so that $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2)$.

Note 2. It is usual refer to given \mathbf{X} and \mathbf{y} as to *data set* and to y_i with corresponding \mathbf{x}_i as to *ith data sample* or *observation*.

1.1.1 Prediction with the linear regression model

The linear regression model contains the vector \mathbf{w} of regression coefficients which are unknown and which we need to estimate in order to be able to use the model for predictions. Let us assume that we already have estimated regression coefficients as $\hat{\mathbf{w}}$. Then the predicted values of y are given by

$$\hat{y} = \hat{\mathbf{w}}^T \mathbf{x}. \quad (1.3)$$

The true value of y is given by

$$y = \mathbf{w}^{*T} \mathbf{x} + \varepsilon \quad (1.4)$$

where \mathbf{w}^* represents actual regression coefficients which we aim to estimate.

Because we assume linear dependence between dependent variable y and explanatory variables \mathbf{x} which we assume to be non-random, then what makes y random variable is a random variable ε . Because we assume that $\mathbb{E}(\varepsilon) = 0$ we can see that

$$\mathbb{E}(y) = \mathbf{x}^T \mathbf{w} + \mathbb{E}(\varepsilon) = \mathbf{x}^T \mathbf{w} \quad (1.5)$$

so \hat{y} is a point estimation of the expected value of y .

Intercept

In real world situations it is not usual that $\mathbb{E}(\varepsilon) = 0$. Consider this trivial example.

Example 3. Let us consider that y represents price of the room and x represents the number of the windows in such a room. If this room does not have windows thus $x = 0$ and $\mathbb{E}(\varepsilon) = 0$ then $y = wx + \varepsilon$ equals zero. But it is very unlikely that room without windows is free.

Because of that, it is very common to include one constant regressor $x_1 = 1$. The corresponding coefficient w_1 of \mathbf{w} is called an *intercept*. We refer this model as a *model with an intercept*. The intercept then corresponds to expected value of y when all regressors are zero and prevent the problem from Example 3. This means that intercept can be assumed as a shift so that it corresponds to $\mathbb{E}(\varepsilon) = \mu$. With regards to this fact we can still assume that in the model with intercept $\varepsilon \sim \mathcal{N}(0, \sigma^2)$. In this work, we consider the model with the intercept. This means that we consider $\mathbf{x} = (x_1, x_2 \dots x_p)$ where constant $x_1 = 1$ represents intercept.

Note 4. Sometimes in the model with an intercept, the explanatory variable \mathbf{x} is marked as $\mathbf{x} \in \mathbb{R}^{p+1}, \mathbf{x} = (x_0, x_1 \dots x_p)$, which means that actual observation $\mathbf{x} \in \mathbb{R}^p$ and the intercept $x_0 = 1$ is explicitly marked.

1.2 Ordinary least squares

We want to estimate \mathbf{w} so that an error of the model on the whole data set is the least possible. This error is measured by a *loss function* $L : \mathbb{R}^2 \rightarrow \mathbb{R}$, which in case of the *ordinary least squares* (OLS) method is quadratic loss function $L(y, \hat{y}) := (y - \hat{y})^2$. So the idea is to find $\hat{\mathbf{w}}$ so that it minimizes the sum of squared *residuals*

$$r_i(\mathbf{w}) = y_i - \hat{y}_i = y_i - \mathbf{x}_i^T \mathbf{w}, \quad i = 1, 2, \dots, n. \quad (1.6)$$

This is commonly know as residual sum of squares *RSS*

$$RSS(\mathbf{w}) = \sum_{i=1}^n r_i^2(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2. \quad (1.7)$$

Definition 5. The RSS as the function of \mathbf{w} is an *objective function* for OLS

$$\text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 = (\mathbf{Y} - \mathbf{X}\mathbf{w})^T (\mathbf{Y} - \mathbf{X}\mathbf{w}). \quad (1.8)$$

The point of the minimum of this function

$$\hat{\mathbf{w}}^{(OLS, n)} = \arg \min_{\mathbf{w} \in \mathbb{R}^p} \text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}(\mathbf{w}) \quad (1.9)$$

is the ordinary least squares estimate of regression coefficients.

To find the minimum of this function, we first need to find the gradient by calculating all partial derivatives

$$\frac{\partial \text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}}{\partial w_j} = \sum_{i=1}^n 2(y_i - \mathbf{w}^T \mathbf{x}_i)(-x_{ij}), \quad j \in \{1, 2, \dots, p\}. \quad (1.10)$$

1. LEAST TRIMMED SQUARES

By this we obtain the gradient

$$\nabla \text{OF}^{(OLS, \mathbf{X}, \mathbf{y})} = - \sum_{i=1}^n 2(y_i - \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i. \quad (1.11)$$

Putting gradient equal to zero we get the so called *normal equations*

$$- \sum_{i=1}^n 2(y_i - \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i = 0 \quad (1.12)$$

that can be rewritten in a matrix form as

$$\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \mathbf{w} = 0. \quad (1.13)$$

Let us now construct the hessian matrix using second-order partial derivatives:

$$\frac{\partial^2 \text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}}{\partial w_h \partial w_j} = \sum_{i=1}^n 2(-x_{ik})(-x_{ij}), \quad h \in \{1, 2, \dots, p\}. \quad (1.14)$$

We get

$$\mathbf{H}_{\text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}} = 2\mathbf{X}^T \mathbf{X}. \quad (1.15)$$

We can see that hessian $\mathbf{H}_{\text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}}$ is always positive semi-definite because for all $\mathbf{s} \in \mathbb{R}^p$

$$\mathbf{s}^T (2\mathbf{X}^T \mathbf{X}) \mathbf{s} = 2(\mathbf{X} \mathbf{s})^T (\mathbf{X} \mathbf{s}) = 2 \|\mathbf{X} \mathbf{s}\|^2 \quad (1.16)$$

It is easy to prove that twice differentiable function is convex if and only if the hessian of such function is positive semi-definite. Moreover, any local minimum of the convex function is also the global one. Hence the solution of (1.13) gives us the global minimum.

Assuming that $\mathbf{X}^T \mathbf{X}$ is a regular matrix, its inverse exists, and the solution can be explicitly written as

$$\hat{\mathbf{w}}^{(OLS, n)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad (1.17)$$

where (OLS, n) denotes that $\hat{\mathbf{w}}^{(OLS, n)}$ is the estimate of true regression coefficients given by the OLS method for n observations.

Moreover, we can see that if $\mathbf{X}^T \mathbf{X}$ is a regular matrix, then the hessian $\mathbf{H}_{\text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}}$ is positive definite and $\text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}$ is strictly convex and $\hat{\mathbf{w}}^{(OLS, n)}$ is the unique strict global minimum.

1.2.1 Properties of the OLS estimate

Gauss-Markov theorem [1] tells us that if particular assumptions about the regression model are fulfilled then the OLS estimate is unbiased and efficient. Gauss-Markov theorem states that OLS is the best linear unbiased estimator (BLUE). Being an efficient estimate means that any other linear unbiased estimate has the same or higher variance. The most important conditions are:

- Expected value of errors is zero.
- Errors are independently distributed and uncorrelated thus
 $cov(\varepsilon_i, \varepsilon_j), i, j = 1, 2 \dots, n, i \neq j$
- All errors have same finite variance. This is known as *homoscedasticity*.

There are also other theorems which describe properties of OLS under specific conditions, but they are out of the scope of this work.

1.3 Robust statistics

Standard statistics methods rely on multiple assumptions and fail if those assumptions are not met. The goal of the robust statistics is to produce acceptable results even when the data are from some unconventional distributions or if data contains outliers or errors which are not normally distributed.

Such assumptions about the OLS method are described in Section 1.2.1. Before we explain what happens if those conditions are not met or met only partially let us describe one of the most common reasons why assumptions are false.

1.3.1 Outliers

We stated many assumptions that are required for the OLS method to produce a acceptable estimate of $\hat{\boldsymbol{w}}$. Unfortunately, in real conditions, these assumptions are often false so that the ordinary least squares do not guarantee to return reasonable results. One of the most common reasons for the assumptions being not met are observations called outliers.

Outliers are for instance erroneous measurements such as transmission errors or noise. Another common reason for outliers is that nowadays the data are automatically processed by computers. Sometimes we are also given data which are heterogeneous in the sense that they contain data from multiple regression models. In some sense outliers are inevitable. One would say that we should be able to eliminate them by precise examination, repair or removal of such data. That is possible in some cases, but often the data we are dealing with are too big and highly dimensional to check.

The robust methods are sometimes not only useful to create models that are not being unduly affected by the presence of outliers but also capable of identifying data which seems to be outliers.

We use terminology from [2] to describe certain types of outliers. Let us have observation (y_i, \mathbf{x}_i) . If the observation is not outlying in any direction we call it *regular observation*. If it is outlying in direction of the explanatory variable \mathbf{x}_i we call it *leverage point*. We have two types of leverage points. If \mathbf{x}_i is outlying but (y_i, \mathbf{x}_i) follows the liner pattern we call it a *good leverage point*. If it does not follow such a pattern we call it *bad leverage point*. Finally if (y_i, \mathbf{x}_i) is outlying only in direction of y_i , we call it a *vertical outlier*.

1.3.2 Measuring robustness

There are a couple of tools to measure the robustness of an estimate. One of the most popular one is called *breakdown point*. Others are *empirical influence function* and *influence function and sensitivity curve*. Here we describe only breakdown point right now. More on robustness measures can be found in [3].

Definition 6. Let T be a statistics, $\mathbf{x} = (x_1, x_2, \dots, x_n)$ be an n -element random sample and $T_n(\mathbf{x})$ value of this statistics. The breakdown point of T at sample \mathbf{x} is defined using sample $\bar{\mathbf{x}}^{(k)}$, that arose by replacing k points from the original sample \mathbf{x} with random values x_i . Then the *breakdown point* is

$$\text{bdpoint}(T, \mathbf{x}_n) = \frac{1}{n} \min S_{T, \mathbf{x}_n}, \quad (1.18)$$

where

$$S_{T, \mathbf{x}_n, D} = \left\{ k \in \{1, 2, \dots, n\} : \sup_{\bar{\mathbf{x}}^{(k)}} \|T_n(\mathbf{x}) - T_n(\bar{\mathbf{x}}^{(k)})\| = \infty \right\}. \quad (1.19)$$

This definition is says that the breakdown point is the function of the minimal number of observations needed to be changed so that the estimator gives arbitrarily biased results.

Intuitively, a reasonable breakdown point should not be higher than 0.5 [4]; if more than 50% of the data is exchanged, the model of exchanged data should override the model of the original data.

In the case of the OLS estimator, one outlier is enough to increase the value of $\hat{\mathbf{w}}^{(OLS, n)}$ to any desired value [5] thus

$$\text{bdpoint}(\hat{\mathbf{w}}^{(OLS, n)}, \mathbf{x}_n) = \frac{1}{n}. \quad (1.20)$$

Figure 1.1 gives us an idea of how one outlier may change the hyperplane given by the OLS estimator of regression coefficients.

For an increasing number of the data samples n the breakdown point of $\hat{\mathbf{w}}^{(OLS, n)}$ tends to zero. We can see that ordinary least squares estimator is

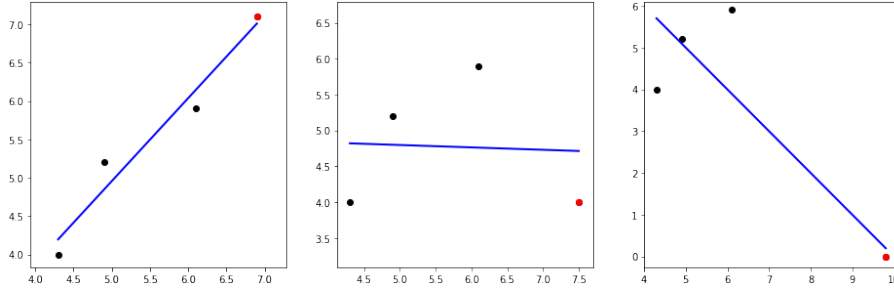


Figure 1.1: Change of the regression hyperplane given by coefficients estimated with OLS method when one of the four observations (highlighted with red color) starts to deviate from the linear pattern.

not resistant to outliers at all. Due to this fact, multiple robust estimators alternatives to the OLS have been proposed.

1.4 Least trimmed squares

The least trimmed squares (LTS) estimator is a robust version of the OLS estimator. In this section, we give a definition and show that its breakdown point is variable and can go up to the maximum possible value of breakdown point, thus 0.5.

Definition 7. Let us have $\mathbf{X} \in \mathbb{R}^{n,p}$, $\mathbf{y} \in \mathbb{R}^{n,1}$, $\mathbf{w} \in \mathbb{R}^p$ and h , $n/2 \leq h \leq n$. The objective function of LTS for data \mathbf{X} and \mathbf{y} is

$$\text{OF}^{(LTS,h,n)}(\mathbf{w}) = \sum_{i=1}^h r_{i:n}^2(\mathbf{w}) \quad (1.21)$$

where $r_{i:n}^2(\mathbf{w})$ denotes the i th smallest squared residuum at \mathbf{w} , i.e.

$$r_{1:n}^2(\mathbf{w}) \leq r_{2:n}^2(\mathbf{w}) \leq \dots \leq r_{n:n}^2(\mathbf{w}). \quad (1.22)$$

Even though that objective function of LTS seems similar to the OLS objective function, finding the minimum is far more complex because the order of the least squared residuals depends on \mathbf{w} . Moreover, $r_{i:n}^2(\mathbf{w})$ residuum is not uniquely determined if more squared residuals have same value. This makes finding the LTS estimate non-convex optimization problem and, in fact, finding the global minimum is NP-hard [6].

1.4.1 Discrete objective function

The LTS objective function from Definition 7 is not differentiable and not convex, so we are unable to use the same approach as with the OLS objective

function. Let us transform this objective function to a discrete version which is easier to use by algorithms to minimize it.

Let us assume for now that we know the vector $\hat{\mathbf{w}}^{(LTS,h,n)}$ of estimated regression coefficients minimizing the LTS objective function. Let π be the permutation of $\hat{n} = \{1, 2, \dots, n\}$ such that

$$r_{i:n}(\hat{\mathbf{w}}^{(LTS,h,n)}) = r_{\pi(j)}(\hat{\mathbf{w}}^{(LTS,h,n)}), \quad j \in \hat{n}. \quad (1.23)$$

Put

$$Q^{(n,h)} = \left\{ \mathbf{m} \in \mathbb{R}^n \mid m_i \in \{0, 1\}, i \in \hat{n}, \sum_{i=1}^n m_i = h \right\}, \quad (1.24)$$

which is simply the set of all vectors $\mathbf{m} \in \mathbb{R}^n$ which contain h ones and $n - h$ zeros. Let $\mathbf{m}^{(LTS)} \in Q^{(n,h)}$ such that $m_j^{(LTS)} = 1$ when $\pi(j) \leq h$ and $m_j^{(LTS)} = 0$ otherwise. Then

$$\hat{\mathbf{w}}^{(LTS,h,n)} = \arg \min_{\mathbf{w} \in \mathbb{R}^p} \sum_{i=1}^h r_{i:n}^2(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^p} \sum_{i=1}^n m_i^{(LTS)} r_i^2(\mathbf{w}). \quad (1.25)$$

This means that if we know the vector \mathbf{m}_{LTS} than we can compute the LTS estimate as the OLS estimate with \mathbf{X} and \mathbf{Y} multiplied by the diagonal matrix $\mathbf{M}_{LTS} = \text{diag}(\mathbf{m}^{(LTS)})$:

$$\hat{\mathbf{w}}^{(LTS,h,n)} = (\mathbf{X}^T \mathbf{M}_{LTS} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{M}_{LTS} \mathbf{y}. \quad (1.26)$$

In other words, finding the minimum of the LTS objective function can be done by finding the OLS estimates (1.26) for all vectors $\mathbf{m} \in Q^{(n,h)}$. Thus if we denote $\mathbf{X}_M = \mathbf{M} \mathbf{X}$ and $\mathbf{y}_M = \mathbf{M} \mathbf{y}$ which corresponds to the h -element subsets of \mathbf{X} and \mathbf{Y} , then as described in [7],

$$\min_{\mathbf{w} \in \mathbb{R}^p} \text{OF}^{(LTS,h,n)}(\mathbf{w}) = \min_{\mathbf{w} \in \mathbb{R}^p} \sum_{i=1}^h r_{i:n}^2(\mathbf{w}) \quad (1.27)$$

$$= \min_{\mathbf{w} \in \mathbb{R}^p, \mathbf{m} \in Q^{(n,h)}} \sum_{i=1}^n m_i r_i^2(\mathbf{w}) \quad (1.28)$$

$$= \min_{\mathbf{m} \in Q^{(n,h)}} \left(\min_{\mathbf{w} \in \mathbb{R}^p} \text{OF}^{(OLS, \mathbf{M} \mathbf{X}, \mathbf{M} \mathbf{y})}(\mathbf{w}) \right) \quad (1.29)$$

$$= \min_{\mathbf{m} \in Q^{(n,h)}} \left(\min_{\mathbf{w} \in \mathbb{R}^p} \|\mathbf{M} \mathbf{y} - \mathbf{M} \mathbf{X} \mathbf{w}\|^2 \right). \quad (1.30)$$

Substituting \mathbf{w} with the OLS estimate as in (1.25) we get

$$\hat{\mathbf{w}}^{(LTS,h,n)} = \min_{\mathbf{m} \in Q^{(n,h)}} \left(\left\| \mathbf{M} \mathbf{y} - \mathbf{M} \mathbf{X} (\mathbf{X}^T \mathbf{M} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{M} \mathbf{y} \right\|^2 \right)$$

We get the discrete objective function with domain $\mathbf{m} \in Q^{(n,h)}$

$$\text{OF}_D^{\text{LTS}}(\mathbf{m}) = \|\mathbf{M}\mathbf{y} - \mathbf{M}\mathbf{X}(\mathbf{X}^T\mathbf{M}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{M}\mathbf{y}\|^2. \quad (1.31)$$

Minimizing this OF could be done straightforwardly by iterating over the $Q^{(n,h)}$ set. Unfortunately, this set has cardinality equal to $\binom{n}{h}$, which is huge, so this approach is infeasible for bigger data sets. Multiple algorithms were proposed to overcome this problem. Majority of them are probabilistic algorithms, but besides those, some exact algorithms were proposed.

Finally, let us point out some fact about the number h of non-trimmed residuals and how it makes least trimmed squares robust. The LTS reaches maximum breakdown point 0.5 at $h = \lfloor (n/2) \rfloor + \lfloor (p+1)/2 \rfloor$ [5]. This means that up to 50% of the data can be outliers. In practice, the portion of outliers is usually lower; if an upper bound on the percentage of outliers is known, h should be set to match this percentage.

Algorithms

In the previous chapter, we have covered the theoretical background required to implement algorithms that are presented in this chapter. We have introduced the discrete version of the LTS objective function whose minimum is equivalent to the continuous one. Finding the minimum of this function requires to find the particular h -element subset and then calculate the estimate of the corresponding regression coefficients \mathbf{w} . To achieve this, we need to examine all h -element subsets. The exhaustive approach fails due to the exponential size of $Q^{(n,h)}$.

2.0.1 First attempts

First known algorithms were based on iterative removal data samples whose residuum had the highest value based on the OLS estimate on the whole dataset. Such attempts are known to be flawed [8] because the initial OLS fit can be already profoundly affected by outliers and the algorithm may remove data samples which represent the actual model.

Other algorithms were based purely on a random approach. One such algorithm is the Random solution algorithm [9] which randomly selects k h -element subsets and subsequently compute the OLS estimate on each of them and chooses the estimate with a minimum value of the objective function. Such approach is straightforward, but the probability of selecting at least one h -element out of k subsets which does not contain outliers thus has a chance of producing good result tends to zero for an increasing number of data samples n as we describe in detail in Section 2.2.2.

Another very similar algorithm called Resampling algorithm introduced in [10]. It selects vectors from $Q^{(n,p+1)}$ instead of $Q^{(n,h)}$. This minor tweak has a higher chance to succeed. Mainly because the probability of selecting k h -element subsets of size $p + 1$ gives the nonzero probability of selecting at least one h -element subset not containing outliers (see Section 2.2.2 for more details). Besides that, the number of vectors in this set is significantly lower

than in $Q^{(n,h)}$ (at least if h is conservatively chosen so that $h = \lceil (n/2) \rceil + \lceil (p+1)/2 \rceil$).

2.0.2 Strong and weak necessary conditions

Generating all possible h -element subsets is computationally exhaustive and relying on randomly selecting “good” h -element subsets does not lead to reliable results. So what are our options? In [11] two criteria called a *weak necessary condition* and a *strong necessary condition* are introduced. They state necessary properties which some h -element subset must satisfy to be a subset which leads to the global optimum of the LTS objective function. Let us introduce those two necessary conditions. For that, it is convenient not only to label a h -element subset of *non-trimmed* observations but also the complementary subset of trimmed observations. We refer to this complementary subset as to *trimmed* subset.

Definition 8. A h -element subset corresponding to $\mathbf{m} \in Q^{(n,h)}$ satisfies the *strong necessary condition* if for any vector \mathbf{m}_{swap} which differs from \mathbf{m} only by swapping one 1 with one 0, we get $\text{OF}_D^{\text{LTS}}(\mathbf{m}) \leq \text{OF}_D^{\text{LTS}}(\mathbf{m}_{\text{swap}})$. In words, the value of the discrete LTS objective function cannot be reduced by swapping one non-trimmed observation with one trimmed observation.

Based on this fact an algorithm can be created. We’ll discuss it in detail in Section 2.3.

Definition 9. An h -element subset corresponding to $\mathbf{m} \in Q^{(n,h)}$ satisfies the *weak necessary condition* if $r_i^2(\hat{\mathbf{w}}^{(\text{OLS}, \mathbf{MX}, \mathbf{My})})$ for all trimmed observation is greater or equal to the greatest non-trimmed squared residuum $r_j^2(\hat{\mathbf{w}}^{(\text{OLS}, \mathbf{MX}, \mathbf{My})})$.

Again, based on this criteria an algorithm can be created. Interesting consequence which we use later gives us the following lemma.

Lemma 10. The strong necessary condition is not satisfied unless the weak necessary condition is satisfied. Thus, if a strong condition is satisfied then weak is too.

Proof. Let us assume that we have some $\hat{\mathbf{w}}^{(\text{OLS}, \mathbf{MX}, \mathbf{My})}$ with $\mathbf{m} \in Q^{(n,h)}$ for which the strong necessary condition is satisfied, but the weak necessary condition is not. That means there exists \mathbf{x}_i with y_i from the non-trimmed subset and \mathbf{x}_j and y_j from the trimmed subset such that

$$r_j^2(\hat{\mathbf{w}}^{(\text{OLS}, \mathbf{MX}, \mathbf{My})}) < r_i^2(\hat{\mathbf{w}}^{(\text{OLS}, \mathbf{MX}, \mathbf{My})}).$$

Thus

$$\text{OF}_D^{\text{LTS}}(\mathbf{m}) > \text{OF}_D^{\text{LTS}}(\mathbf{m}) + r_j^2(\hat{\mathbf{w}}^{(\text{OLS}, \mathbf{MX}, \mathbf{My})}) - r_i^2(\hat{\mathbf{w}}^{(\text{OLS}, \mathbf{MX}, \mathbf{My})}) \quad (2.1)$$

Now we need to show that \mathbf{m}_{swap} vector that is created by swapping j th observation from the trimmed subset with i th observation from the non-trimmed subset leads to

$$\text{OF}_D^{\text{LTS}}(\mathbf{m}) + r_j^2(\hat{\mathbf{w}}^{(OLS, \mathbf{M}\mathbf{X}, \mathbf{M}\mathbf{y})}) - r_i^2(\hat{\mathbf{w}}^{(OLS, \mathbf{M}\mathbf{X}, \mathbf{M}\mathbf{y})}) \geq \text{OF}_D^{\text{LTS}}(\mathbf{m}_{swap}). \quad (2.2)$$

That is true because the value of $\text{OF}_D^{\text{LTS}}(\mathbf{m}_{swap})$ is the minimum on the given subset of observations. That is, of course, a contradiction with our assumption which says that the strong necessary condition is satisfied. \square

2.1 Computing OLS

In this section, we describe a few of many methods that can be used to obtain $\hat{\mathbf{w}}^{(OLS, n)}$. Those methods are parts of the algorithms used to calculate the LTS estimate.

In the following algorithms, we describe its time complexity. Because the matrix multiplication is the fundamental part of all the following algorithms, it is, therefore, appropriate to mention a few facts about the time complexity of the matrix multiplication.

There are multiple algorithms for the multiplication of $n \times n$ matrices, for example:

- Naive algorithm; its time complexity is $\mathcal{O}(n^3)$.
- Strassen algorithm; its time complexity is $\mathcal{O}(n^{2.8074})$ [12].
- Coppersmith-Winograd; its time complexity is $\mathcal{O}(n^{2.375477})$ [13]

In practice, however, the naive algorithm is usually used. Even though some of those algorithms have been efficiently implemented and are known to be numerically stable (primarily variations of Strassen algorithm), their error bound is weaker than in the case of the naive algorithm [14]. For this reason, they are not used even when they might improve performance.

Moreover, currently widely used linear algebra software libraries such as LAPACK uses Basic Linear Algebra Subprograms (BLAS) low-level routines which implement naive matrix-matrix multiplication [15].

For that reason we assume in this work that the matrix multiplication is $\mathcal{O}(n^3)$. For not square matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times p}$ it is then $\mathcal{O}(mnp)$.

2.1.1 Computation using a matrix inversion

Calculating the OLS estimate using the matrix inversion can be done as follows:

1. Compute $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ and $\mathbf{B} = \mathbf{X}\mathbf{y}$.

2. Find the inversion of \mathbf{A} .
3. Multiply $\mathbf{A}^{-1}\mathbf{B}$.

Observation 11. Time complexity of computing the OLS estimate on $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^{n \times 1}$ using the matrix inversion is $\mathcal{O}(p^2n)$.

Proof. First, we compute the $\mathbf{A} = \mathbf{X}^T\mathbf{X}$ and $\mathbf{B} = \mathbf{X}^T\mathbf{y}$. That gives us $\mathcal{O}(p^2n + pn)$. Next, we compute inversion $\mathbf{C} = \mathbf{A}^{-1}$ which gives us $\mathcal{O}(p^3)$. Finally, we compute $\hat{\mathbf{w}}^{(OLS,n)} = \mathbf{CB}$ which is $\mathcal{O}(p^2)$. Altogether, we get $\mathcal{O}(p^2n + pn + p^3 + p^2)$; p^2n and p^3 asymptotically dominates over the rest so

$$\mathcal{O}(p^2n + pn + p^3 + p^2) \sim \mathcal{O}(p^2n + p^3). \quad (2.3)$$

Moreover, if we assume that $n \geq p$, we get $\mathcal{O}(p^2n + p^3) \sim \mathcal{O}(p^2n)$. \square

2.1.2 Computation using the Cholesky decomposition

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be a symmetric positive definite matrix. Then it is possible to find lower triangular matrix \mathbf{L} so that

$$\mathbf{A} = \mathbf{LL}^T \quad (2.4)$$

We call this decomposition *Cholesky factorization* (sometimes also *Cholesky decomposition*)

If we look at our problem of finding a solution of

$$\mathbf{X}^T\mathbf{X}\mathbf{w} = \mathbf{X}^T\mathbf{y}, \quad (2.5)$$

we can easily rewrite it as

$$\mathbf{Z}\mathbf{w} = \mathbf{b} \quad (2.6)$$

where $\mathbf{Z} = \mathbf{X}^T\mathbf{X}$ is a symmetric positive definite matrix and $\mathbf{b} = \mathbf{X}^T\mathbf{y}$. Because \mathbf{Z} can be factorized as \mathbf{LL}^T the solution can be found easily by substitution

$$\mathbf{L}\mathbf{d} = \mathbf{b} \quad (2.7)$$

$$\mathbf{L}^T\mathbf{w} = \mathbf{d}, \quad (2.8)$$

where \mathbf{d} and \mathbf{w} can be obtained by forward and backward substitution. Solution of \mathbf{w} then represents $\hat{\mathbf{w}}^{(OLS,n)}$ estimate.

So the algorithm for solving OLS using Cholesky factorization goes as follows:

1. Compute $\mathbf{X}^T\mathbf{X}$ and $\mathbf{X}^T\mathbf{y}$.
2. Compute Cholesky factorization $\mathbf{Z} = \mathbf{LL}^T$ where $\mathbf{Z} = \mathbf{X}^T\mathbf{X}$.

3. Solve lower triangular system $\mathbf{L}\mathbf{d} = \mathbf{b}$ where $\mathbf{b} = \mathbf{X}^T\mathbf{y}$ for \mathbf{d} using forward substitution.
4. Solve upper triangular system $\mathbf{L}^T\mathbf{w} = \mathbf{d}$ for \mathbf{w} .

Observation 12. The time complexity of solving OLS using Cholesky factorization is $\mathcal{O}(p^2n)$

Proof. First step requires two matrix multiplication $\mathbf{X}^T\mathbf{X}$ which is $\mathcal{O}(np^2)$ and $\mathbf{X}^T\mathbf{y}$ which is $\mathcal{O}(np)$. Second step represents computing Cholesky factorization. This can be done in $\mathcal{O}(\frac{1}{3}p^3)$ [16]. In the third and fourth step, we are solving triangular linear systems; both of them require $\mathcal{O}(\frac{1}{2}p^2)$.

Putting all the steps together we get

$$\mathcal{O}(np^2 + np + \frac{1}{3}p^3 + p^2) \quad (2.9)$$

and because we assume $n \geq p$ then the multiplication $\mathbf{X}^T\mathbf{X}$ asymptotically dominates over the rest so we get $\mathcal{O}(p^2n)$. \square

Computing the OLS estimate using the Cholesky factorization is more numerically stable than using matrix inversion and the time complexity is asymptotically similar.

2.1.3 Computation using the QR decomposition

Let us now look on a similar method of computing the OLS estimate which does not require multiplying $\mathbf{X}^T\mathbf{X}$.

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be square matrix. Then exists matrices \mathbf{Q} and \mathbf{R} so that

$$\mathbf{A} = \mathbf{Q}\mathbf{R}, \quad (2.10)$$

where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is an orthogonal matrix and $\mathbf{R} \in \mathbb{R}^{n \times n}$ is an upper triangular matrix. This decomposition is known as *QR decomposition*.

If matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is not square, then decomposition can be found as

$$\mathbf{A} = \mathbf{Q}\mathbf{R} = \mathbf{Q} \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{Q}_2 \end{bmatrix} \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix} = \mathbf{Q}_1\mathbf{R}_1 \quad (2.11)$$

Where $\mathbf{Q} \in \mathbb{R}^{m \times m}$ is an orthogonal matrix, $\mathbf{R}_1 \in \mathbb{R}^{n \times n}$ is upper an triangular matrix, $\mathbf{0} \in \mathbb{R}^{(m-n) \times n}$ is a zero matrix and $\mathbf{Q}_1 \in \mathbb{R}^{n \times n}$ and $\mathbf{Q}_2 \in \mathbb{R}^{(m-n) \times n}$ are matrices with orthogonal columns.

Given a QR decomposition of \mathbf{X}

$$\mathbf{X} = \mathbf{Q}\mathbf{R} \quad (2.12)$$

we get

$$\mathbf{X}^T\mathbf{X} = (\mathbf{Q}\mathbf{R})^T\mathbf{Q}\mathbf{R} = \mathbf{R}^T\mathbf{Q}^T\mathbf{Q}\mathbf{R} \quad (2.13)$$

2. ALGORITHMS

and because \mathbf{Q} is orthogonal, $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ and so

$$\mathbf{X}^T \mathbf{X} = \mathbf{R}^T \mathbf{R}. \quad (2.14)$$

Because $\mathbf{X} \in \mathbb{R}^{n \times p}$ is not a square matrix then $\mathbf{R} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}$ and

$$\mathbf{X}^T \mathbf{X} = \mathbf{R}_1^T \mathbf{R}_1, \quad (2.15)$$

where \mathbf{R}_1^T is a lower triangular matrix.

Solution to

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (2.16)$$

is then given by

$$\mathbf{R}_1^T \mathbf{R}_1 \mathbf{w} = \mathbf{R}_1^T \mathbf{Q}_1^T \mathbf{y}. \quad (2.17)$$

If we assume \mathbf{X} have full column rank, \mathbf{R}_1 must be invertible and this equation can be simplified to

$$\mathbf{R}_1 \mathbf{w} = \mathbf{b}_1 \quad (2.18)$$

where $\mathbf{b}_1 = \mathbf{Q}_1^T \mathbf{y}$. Because \mathbf{R}_1 is upper an triangular matrix, finding the is trivial using the backward substitution. Resulting \mathbf{w} is then the OLS estimate $\hat{\mathbf{w}}^{(OLS,n)}$.

Note 13. We can see that Cholesky factorization $\mathbf{X}^T \mathbf{X} = \mathbf{L} \mathbf{L}^T$ is closely connected to the QR decomposition $\mathbf{X} = \mathbf{Q}_1 \mathbf{R}_1$. Indeed, putting

$$\mathbf{L} = \mathbf{R}_1^T \quad (2.19)$$

we can get Cholesky decomposition directly from the QR decomposition.

The algorithm for solving the OLS using QR decomposition can go as follows:

1. Calculate a QR decomposition $\mathbf{X} = \mathbf{Q} \mathbf{R}^T = \mathbf{Q}_1 \mathbf{R}_1^T$.
2. Calculate $\mathbf{b}_1 = \mathbf{Q}_1^T \mathbf{y}$.
3. Solve upper triangular system $\mathbf{R}_1 \mathbf{w} = \mathbf{b}_1$.

The QR factorization can be calculated in multiple ways. The most basic method is applying the Gram-Schmidt process to columns of the matrix \mathbf{X} . This approach is not numerically stable so in practice it is not used as much as two following methods: *Householder transformations* and *Givens rotations*. The time complexity of both algorithms is similar. QR decomposition of a matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ is [17]

$$\mathcal{O}(2p^2n - \frac{2}{3}p^3) \sim \mathcal{O}(p^2n). \quad (2.20)$$

On the other hand, using Givens rotation for QR decomposition of the same matrix is

$$\mathcal{O}(3np^2 - p^3) \sim \mathcal{O}(p^2n). \quad (2.21)$$

Although Householder transformations are about 50% faster, both have asymptotically equal time complexity. Moreover Givens rotations are more numerically stable. We speak about Givens rotations in detail in Section 2.1.3 where we also make proof of the (2.21). Moreover, Givens rotations are suitable for sparse matrices. We use this property of Givens rotations in Section 2.4.4.

For now, let us look at the time complexity of solving the OLS using Householder transformations.

Observation 14. The time complexity of finding the OLS estimate using the QR decomposition by Householder transformations is $\mathcal{O}(2p^2n - \frac{2}{3}p^3) \sim \mathcal{O}(p^2n)$.

Proof. In the first step, we calculate a QR decomposition. Householder transformations can be done in $\mathcal{O}(2np^2 - \frac{2}{3}p^3)$. The second step consists of the matrix multiplication $\mathbf{Q}_1^T \mathbf{y}$ which is $\mathcal{O}(np)$. In the last step we solve upper triangular linear system which is $\mathcal{O}(\frac{1}{2}p^2)$. Putting all steps together we get

$$\mathcal{O}(2np^2 - \frac{2}{3}p^3 + np + \frac{1}{2}p^2) \quad (2.22)$$

We can see that p^2n asymptotically dominates (this is the same as in the case of Givens rotations). \square

The QR decomposition is considered as a standard way of computing OLS estimate because of its high numerical stability. Let us mention that a little slower, but a more stable method of computing the OLS estimate exists, and that is the one using the singular value decomposition (SVD). On the other hand, the QR decomposition is sufficiently stable for most cases so describing SVD decomposition is out of the scope of this work.

Givens Rotation

In this section, we describe Givens rotations in detail. Moreover, we also show how to update QR decomposition when adding or deleting a from the matrix \mathbf{X} . We use this in the algorithm described in Section 2.4.4.

As described in [18], we compute the QR factorization of $\mathbf{A} \in \mathbb{R}^{m \times n}$ so that we apply an orthogonal transformation using a matrix \mathbf{Q}^T as

$$\mathbf{Q}^T \mathbf{A} = \mathbf{R} \quad (2.23)$$

2. ALGORITHMS

where \mathbf{Q} is a product of orthogonal matrices. These matrices have the following matrix as sub-matrix:

$$\mathbf{Q}_\varphi = \begin{bmatrix} \cos(\varphi) & \sin(\varphi) \\ -\sin(\varphi) & \cos(\varphi) \end{bmatrix} \quad (2.24)$$

This matrix is indeed orthogonal since

$$\mathbf{Q}_\varphi \mathbf{Q}_\varphi^T = \mathbf{Q}_\varphi^T \mathbf{Q}_\varphi = \mathbf{I}. \quad (2.25)$$

Moreover, if we multiply this orthogonal matrix with some vector $\mathbf{x} \in \mathbb{R}^2$, as a result $\mathbf{Q}\mathbf{x}$ we get a vector of the same length as \mathbf{x} which is rotated clockwise by the angle φ . When we say that vector has the same length we mean that L^2 norm of such vector stays the same. This is an important property of all orthogonal matrices (operations that preserve L^2 norm are known as *unitary transformations*). Let us verify this claim. Let us have an orthogonal matrix $\mathbf{Q} \in \mathbb{R}^{m \times m}$ and any vector $\mathbf{x} \in \mathbb{R}^m$ then

$$\|\mathbf{Q}\mathbf{x}\|^2 = (\mathbf{Q}\mathbf{x})^T(\mathbf{Q}\mathbf{x}) = \mathbf{x}^T \mathbf{Q}^T \mathbf{Q} \mathbf{x} = \mathbf{x}^T \mathbf{I} \mathbf{x} = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|^2. \quad (2.26)$$

So as we said, the idea is to create a series of orthogonal matrices which gradually rotates two-element column sub-vectors of \mathbf{A} , so that we obtain zeros under diagonal. One such method - *Givens rotation* uses *Givens matrices* that erase one element under diagonal at a time.

The example of (2.24) is not random. Let us look at the orthogonal matrix \mathbf{Q}_φ one more time and let us multiply this matrix with some vector.

$$\begin{bmatrix} \cos(\varphi) & \sin(\varphi) \\ -\sin(\varphi) & \cos(\varphi) \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} r \\ z \end{bmatrix}. \quad (2.27)$$

To obtain the zeroing effect, we need to rotate this vector so that it is parallel to $(1, 0)^T$. So we only need to find φ so that $\cos(\varphi)a + \sin(\varphi)b = r$. As we know, rotation with \mathbf{Q}_φ preserves L^2 norm. Hence if we want $z = 0$, then r must be equal to L^2 norm of the vector $(a, b)^T$ i.e. $r = \sqrt{a^2 + b^2}$. This leads to the solution

$$\cos(\varphi) = \frac{a}{\sqrt{a^2 + b^2}} \quad (2.28)$$

and

$$\sin(\varphi) = \frac{b}{\sqrt{a^2 + b^2}}. \quad (2.29)$$

In practice, the algorithm of computing $\cos(\varphi)$ and $\sin(\varphi)$ is slightly different because we want to prevent arithmetic overflow. This algorithm is described by the following pseudocode:

Algorithm 1: Rotate

Input: a, b
Output: \cos, \sin

```

1  $\sin \leftarrow \emptyset$ ;
2  $\cos \leftarrow \emptyset$ ;
3 if  $b == 0$  then
4    $\sin \leftarrow 0$ ;
5    $\cos \leftarrow 1$ ;
6 else if  $\text{abs}(b) \geq \text{abs}(a)$  then
7    $\cot g \leftarrow \frac{a}{b}$ ;
8    $\sin \leftarrow \frac{1}{\sqrt{1+(\cot g)^2}}$ ;
9    $\cos \leftarrow \sin \cot g$ ;
10 else
11    $\tan \leftarrow \frac{b}{a}$ ;
12    $\cos \leftarrow \frac{1}{\sqrt{1+(\tan)^2}}$ ;
13    $\sin \leftarrow \cos \tan$ ;
14 end
15 return  $\cos, \sin$ ;

```

We can scale the same idea to higher dimensions. Let us denote matrix $\mathbf{Q}_\varphi(i, j) \in \mathbb{R}^{m \times m}$ defined as

$$\mathbf{Q}_\varphi(i, j) = \begin{matrix} & & i & & j & & \\ & & & & & & \\ i & & \begin{bmatrix} 1 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \dots & c & \dots & s & \dots & 0 \\ \vdots & & \vdots & \ddots & \vdots & & \vdots \end{bmatrix} & & \\ j & & \begin{bmatrix} 0 & \dots & -s & \dots & c & \dots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 1 \end{bmatrix} & & \end{matrix}$$

where $c = \cos(\varphi)$ and $s = \sin(\varphi)$ for some φ . That means this matrix is orthogonal. Now if we have some column vector $\mathbf{x} \in \mathbb{R}^{m \times 1}$ and we calculate c and s by (2.28) and (2.29) for $a := x_i$ and $b := x_j$. Finally if we multiply $\mathbf{Q}_\varphi(i, j)\mathbf{x} = \mathbf{p}$ we can see that \mathbf{p} is the same vector as \mathbf{x} except for p_i and p_j so that:

$$\mathbf{Q}_\varphi(i, j)\mathbf{x} = \mathbf{Q}_\varphi(i, j) \begin{bmatrix} x_1 \\ \vdots \\ x_i \\ \vdots \\ x_j \\ \vdots \\ x_m \end{bmatrix} = \mathbf{p} = \begin{bmatrix} x_1 \\ \vdots \\ r \\ \vdots \\ 0 \\ \vdots \\ x_m \end{bmatrix}$$

where $r = \sqrt{x_i^2 + x_j^2}$. If we have matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$ instead of only one column \mathbf{x} , it works the same way. We need to create matrix $\mathbf{Q}_\varphi(i, j)$ for each a_{ij} under the diagonal in order to create upper triangular matrix. Usually we are zeroing columns so that we start with a_{12} and continue with $a_{13} \dots a_{1n}$. Then we start with the second column with element a_{23} and so on. That means we need to create in total $e = \frac{p^2 - p}{2} + np - p^2$ matrices $\mathbf{Q}_\varphi(i, j)$. We can denote this sequence of matrices as $\mathbf{Q}_{\varphi_1}, \mathbf{Q}_{\varphi_2}, \dots, \mathbf{Q}_{\varphi_e}$. The QR decomposition then looks like

$$\mathbf{Q}_{\varphi_e} \dots \mathbf{Q}_{\varphi_2} \mathbf{Q}_{\varphi_1} \mathbf{A} = \mathbf{R} \quad (2.30)$$

where $\mathbf{Q}_{\varphi_e} \dots \mathbf{Q}_{\varphi_2} \mathbf{Q}_{\varphi_1}$ is actually \mathbf{Q}^T and \mathbf{Q} is obtained by

$$\mathbf{Q} = \mathbf{Q}_{\varphi_1}^T \mathbf{Q}_{\varphi_2}^T \dots \mathbf{Q}_{\varphi_e}^T. \quad (2.31)$$

We can see that for the each non-zero under diagonal element of the matrix \mathbf{A} , we need one additional matrix \mathbf{Q}_{φ_i} , and with such matrix, we are multiply only two rows. Moreover, the rows are getting shorter after each finished column. Hence the total number of operations is at most

$$\sum_{i=1}^n \sum_{j=i+1}^p 6(n-i+1) \approx 6np^2 - 3np^2 - 3p^3 + 2p^3 = 3np^2 - p^3. \quad (2.32)$$

2.2 FAST-LTS

In this section, we introduce the FAST-LTS algorithm from [19]. It is, as well as other algorithms we introduce an iterative algorithm. We discuss all main components of the algorithm starting with its core idea called a concentration step which authors call a *C-step*.

2.2.1 C-step

We show that from an existing LTS estimate $\hat{\mathbf{w}}$ we can construct a new LTS estimate $\hat{\mathbf{w}}_{new}$ so that value the objective function at $\hat{\mathbf{w}}_{new}$ is less or equal to

the value at $\hat{\mathbf{w}}$. Based on this property, the algorithm creates a sequence of LTS estimates leading to better results.

Theorem 15. Consider $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^{n \times 1}$. Let us also have $\mathbf{w}_0 = (m_1^{(0)}, \dots, m_n^{(0)}) \in \mathbb{R}^p$ and $\mathbf{m}_0 \in Q^{(n,h)}$. Let us put $L_0 = \sum_{i=1}^n m_i^{(0)} r_i^2(\mathbf{w}_0)$. Let $\pi : \hat{n} \rightarrow \hat{n}$ be permutation of \hat{n} such that

$$|r_{\pi(1)}(\mathbf{w}_0)| \leq \dots \leq |r_{\pi(n)}(\mathbf{w}_0)| \quad (2.33)$$

and mark $\mathbf{m}_1 \in Q^{(n,h)}$ such that $m_i^{(1)} = 1$ for $i \in \{\pi(1), \pi(2), \dots, \pi(h)\}$ and $m_i^{(1)} = 0$ otherwise. This means that \mathbf{m}_1 corresponds to h -element subset with smallest squared residuals $r_i^2(\mathbf{w}_0)$.

Finally let \mathbf{w}_1 be the least squares fit on the \mathbf{m}_1 h -element subset of observations and $L_1 = \sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_1)$, then

$$L_1 \leq L_0. \quad (2.34)$$

Proof. Because \mathbf{m}_1 represents h observations with the smallest squared residuals $r_i^2(\mathbf{w}_0)$ at point \mathbf{w}_0 , then $\sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_0) \leq \sum_{i=1}^n m_i^{(0)} r_i^2(\mathbf{w}_0) = L_0$. When we take into account that the OLS estimate minimizes the sum of squared residuals for the \mathbf{m}_1 subset of observations, then

$$L_1 = \sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_1) \leq \sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_0). \quad (2.35)$$

Together we get

$$L_1 = \sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_1) \leq \sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_0) \leq \sum_{i=1}^n m_i^{(0)} r_i^2(\mathbf{w}_0) = L_0 \quad (2.36)$$

□

Let us denote $\hat{\mathbf{w}}_{old}^{(LTS,h,n)} = \text{OF}_D^{\text{LTS}}(\mathbf{m}_{old})$ and $L_{old} = \sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_{old})$ for some vector \mathbf{m}_{old} . Then based on the previous theorem, using some \mathbf{m}_{old} with corresponding $\hat{\mathbf{w}}_{old}^{(LTS,h,n)}$ we can construct \mathbf{m}_{new} with corresponding $\hat{\mathbf{w}}_{new}^{(LTS,h,n)}$ such that $L_{new} \leq L_{old}$.

Applying the theorem repetitively leads to the iterative sequence of $L_1 \leq L_2 \leq \dots$. One step, called the *C-step*, of this process is described by the

2. ALGORITHMS

following pseudocode.

Algorithm 2: C-step

Input: dataset consisting of $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^{n \times 1}$, $\hat{\mathbf{w}}_{old} \in \mathbb{R}^{p \times 1}$

Output: $\hat{\mathbf{w}}_{new}$, \mathbf{m}_{new}

- 1 $R \leftarrow \emptyset$;
 - 2 **for** $i \leftarrow 1$ **to** n **do**
 - 3 $R \leftarrow R \cup \{|y_i - \hat{\mathbf{w}}_{old} \mathbf{x}_i^T|\}$;
 - 4 **end**
 - 5 $\mathbf{m}_{new} \leftarrow$ select set of h smallest absolute residuals from R ;
 - 6 $\hat{\mathbf{w}}_{new} \leftarrow$ OLS fit the on \mathbf{m}_{new} subset;
 - 7 **return** $\hat{\mathbf{w}}_{new}$, \mathbf{m}_{new} ;
-

C-step algorithm is visualized on Figure 2.1 where we start with h -element subset \mathbf{m}_1 and corresponding $\hat{\mathbf{w}}_1$ OLS fit on the \mathbf{m}_1 subset and $L_1 = \text{OF}_D^{\text{LTS}}(\mathbf{m}_1)$. By sorting absolute the residuals and selecting h smallest we obtain \mathbf{m}_2 h -element subset. Its value $\text{OF}^{(OLS, \mathbf{M}_2 \mathbf{X}, \mathbf{M}_2 \mathbf{y})}(\hat{\mathbf{w}}_1)$ is highlighted with red dot. Finally we calculate OLS fit on \mathbf{m}_2 and obtain $\hat{\mathbf{w}}_2$ estimate.

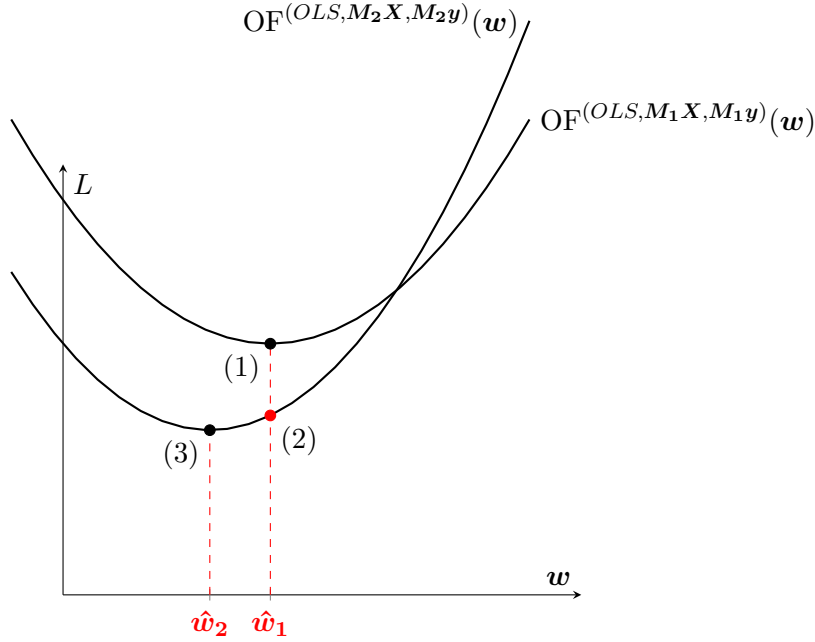


Figure 2.1: Illustration of the C-step algorithm. (1) represents the value of $\text{OF}^{(OLS, \mathbf{M}_1 \mathbf{X}, \mathbf{M}_1 \mathbf{y})}(\hat{\mathbf{w}}_1)$ (which is equal to the value $\text{OF}_D^{\text{LTS}}(\mathbf{m}_1)$). (2) represents the value of $\text{OF}^{(OLS, \mathbf{M}_2 \mathbf{X}, \mathbf{M}_2 \mathbf{y})}(\hat{\mathbf{w}}_1)$ and (3) represents the value of $\text{OF}^{(OLS, \mathbf{M}_2 \mathbf{X}, \mathbf{M}_2 \mathbf{y})}(\hat{\mathbf{w}}_2)$ (which is equal to the value $\text{OF}_D^{\text{LTS}}(\mathbf{m}_2)$).

Observation 16. The time complexity of the C-step 2 algorithm is asymptotically similar to the time complexity of the computation of the OLS fit,

namely $\mathcal{O}(p^2n)$.

Proof. In the C-step we must compute n absolute residuals. Computation of one absolute residual consists of matrix multiplication of shapes $1 \times p$ and $p \times 1$ that gives us $\mathcal{O}(p)$. Hence, the time of computing n residuals is $\mathcal{O}(np)$. Next, we must select a set of h smallest residuals which can be done in $\mathcal{O}(n)$ using a modification of the QuickSelect algorithm [20]. Finally, we must compute the $\hat{\mathbf{w}}_{new}$ OLS estimate on an h -element subset of the data. Because h is proportional to n , we can say that this is $\mathcal{O}(p^2n + p^3)$ which is asymptotically dominant over the previous steps which are $\mathcal{O}(np + n)$. Because we assume $n \geq p$, we get $\mathcal{O}(p^2n + p^3) \sim \mathcal{O}(p^2n)$. \square

As we stated above, repeating C-step leads to a sequence of $\hat{\mathbf{w}}_1, \hat{\mathbf{w}}_2 \dots$ on subsets $\mathbf{m}_1, \mathbf{m}_2 \dots$ with corresponding sum of squared residuals $L_1 \geq L_2 \geq \dots$. One could ask if this sequence converges, so that $L_i = L_{i+1}$. Answer to this question is presented by the following theorem.

Theorem 17. The sequence of the estimates $\hat{\mathbf{w}}_1, \hat{\mathbf{w}}_2 \dots$ obtained by the C-step becomes constant after at most $k = \binom{n}{h}$, i.e. $\hat{\mathbf{w}}_k = \hat{\mathbf{w}}_{k+1}$.

Proof. $\hat{\mathbf{w}}_i$ is uniquely given by the h -element subset $\mathbf{m}_i \in Q^{(n,h)}$ and since $Q^{(n,h)}$ is finite, namely its size is $\binom{n}{h}$, the sequence becomes constant at the latest after this number of steps. \square

The theorem gives us a clue to create algorithm described by the following pseudocode.

Algorithm 3: Repeat-C-step

Input: dataset consisting of $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^{n \times 1}$, $\hat{\mathbf{w}}_{old} \in \mathbb{R}^{p \times 1}$, \mathbf{m}_0
Output: $\hat{\mathbf{w}}_{final}$, \mathbf{m}_{final}

```

1  $\hat{\mathbf{w}}_{new} \leftarrow \emptyset$ ;
2  $\mathbf{m}_{new} \leftarrow \emptyset$ ;
3  $L_{new} \leftarrow \infty$ ;
4 while True do
5    $L_{old} \leftarrow \text{OF}_D^{\text{LTS}}(\hat{\mathbf{w}}_{old})$ ;
6    $\hat{\mathbf{w}}_{new}, \mathbf{m}_{new} \leftarrow \text{C-step}(\mathbf{X}, \mathbf{y}, \hat{\mathbf{w}}_{old})$ ;
7    $L_{new} \leftarrow \text{OF}_D^{\text{LTS}}(\hat{\mathbf{w}}_{new})$ ;
8   if  $L_{old} == L_{new}$  then
9     break
10  end
11   $\hat{\mathbf{w}}_{old} \leftarrow \hat{\mathbf{w}}_{new}$ 
12 end
13 return  $\hat{\mathbf{w}}_{new}, \mathbf{m}_{new}$ ;
```

It is important to note that although the maximum number of steps of this algorithm is $\binom{n}{h}$, in practice, it is most often under 20 steps as can be seen on Figure 2.2.

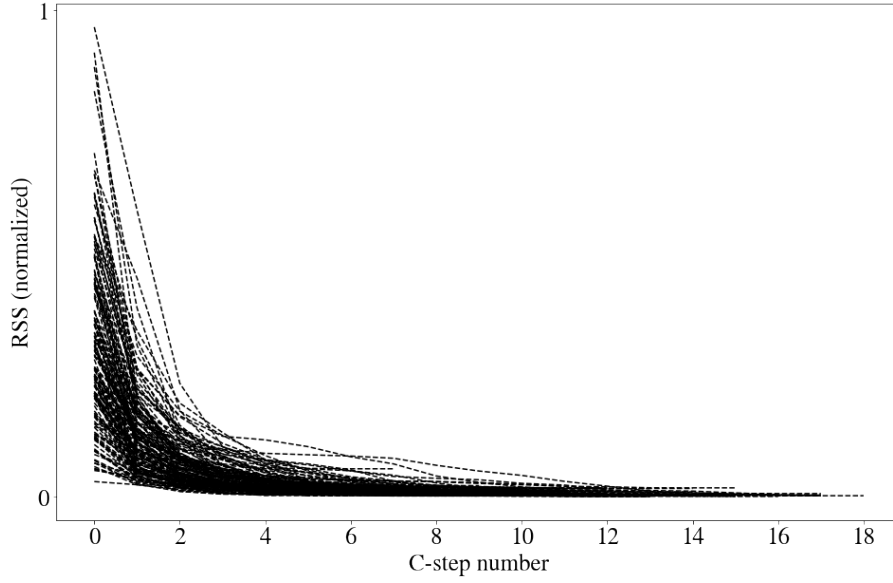


Figure 2.2: Value of the residual sum of squares (normalized) based on the number of the step of C-step algorithm for 100 different starting subsets. Dataset D3 was used with configuration $n = 500, p = 20$ and 30% of the the outliers (see Section 3 for more details about the dataset).

Now we can describe the final algorithm from [19]: Choose a lot of initial subsets \mathbf{m}_1 and for each of them apply the Repeat-C-step algorithm. From all resulting subsets with the corresponding $\hat{\mathbf{w}}$ estimates choose the one with the least value of $\text{OF}_D^{\text{LTS}}(\hat{\mathbf{w}})$.

Before we can construct final the algorithm, we must decide how to choose the initial subset \mathbf{m}_1 and how many of them means “a lot”.

2.2.2 Choosing an initial \mathbf{m}_1 subset

It is important to note, that when we choose \mathbf{m}_1 subset such that it contains outliers, then iteration of C-steps usually does not converge to a good results, so we should focus on methods with non zero probability of selecting \mathbf{m}_1 such that it does not contain outliers. There are many possibilities of how to create an initial \mathbf{m}_1 subset. Let us start with the most trivial one.

Random selection

Most basic way of creating the \mathbf{m}_1 subset is simply to choose random $\mathbf{m}_1 \in Q^{(n,h)}$. The following observation shows that it is not the best way.

Observation 18. Assume that the n -element dataset contains outliers whose number is proportional to ϵ/n with $\epsilon > 0$. Let $\mathbf{m}_{1_1}, \dots, \mathbf{m}_{1_k}, \mathbf{m}_{1_i} \in Q^{(n,h)}$

be k randomly selected h -element subset of data. Then the probability that at least one of these h -element subsets does not contain an outlier tends to 0 as n goes to infinity.

Proof. It follows from observation above and the fact that $h > n/2$ that

$$\begin{aligned} P(\text{one random data sample not an outliers}) &= (1 - \epsilon) \\ P(\text{one random } h \text{ element subset without outliers}) &= (1 - \epsilon)^h \\ P(\text{one subset with at least one outlier}) &= 1 - (1 - \epsilon)^h \\ P(k \text{ subsets with at least one outlier in each}) &= (1 - (1 - \epsilon)^h)^k \\ P(k \text{ subsets with at least one subset without outliers}) &= 1 - (1 - (1 - \epsilon)^h)^k \end{aligned}$$

Because $n \rightarrow \infty$, then $(1 - \epsilon)^h \rightarrow 0$, $1 - (1 - \epsilon)^h \rightarrow 1$, $(1 - (1 - \epsilon)^h)^k \rightarrow 1$, and $1 - (1 - (1 - \epsilon)^h)^k \rightarrow 0$ \square

That means that we should consider other options for selecting initial \mathbf{m}_1 subset. Authors of the algorithm came with the following solution.

P-subset selection

Let us choose a vector $\mathbf{c} \in Q^{(n,p)}$ and compute the rank of the matrix $\mathbf{X}_C = \mathbf{C}\mathbf{X}$, where $C = \text{diag}(\mathbf{c})$. If $\text{rank}(\mathbf{X}_C) < p$ add randomly selected rows of \mathbf{X} to \mathbf{X}_C without repetition until $\text{rank}(\mathbf{X}_C) = p$. Next let us denote $\hat{\mathbf{w}}_0 = \text{OF}_D^{\text{LTS}}(\mathbf{c})$. Let $\pi : \hat{n} \rightarrow \hat{n}$ be the permutation of \hat{n} such that $|r_{\pi(1)}(\mathbf{c})| \leq \dots \leq |r_{\pi(n)}(\mathbf{c})|$.

Finally, let $\mathbf{m}_1 \in Q^{(n,h)}$ be initial h -element subset such that $m_i^{(1)} = 1$ for $i \in \{\pi(1), \pi(2), \dots, \pi(h)\}$ and $m_i^{(1)} = 0$ otherwise.

Observation 19. Assume that the n -element dataset contains outliers whose number is proportional to ϵ/n with $\epsilon > 0$. Let $\mathbf{c}_{1_1}, \dots, \mathbf{c}_{1_k}, \mathbf{c}_{1_i} \in Q^{(n,p)}$ be k randomly selected p -element subset of data. Then the probability that at least one of these p -element subsets does not contain an outlier tends to

$$1 - (1 - (1 - \epsilon)^h)^k > 0. \quad (2.37)$$

Proof. Similarly as in previous observation. \square

The last missing piece of the algorithm is determining the number k of initial subsets \mathbf{m}_1 , which maximize the probability to at least one of them leads to a sequence of estimates ending up in the global minimum. Simply put, the more, the better. So before we answer this question accurately, let us discuss some key observations about the algorithm.

2.2.3 Speed-up of the algorithm

In this section, we describe essential observations which help us to formulate the final algorithm. In two subsections we describe how to optimize the current algorithm.

Selective iteration

The most computationally demanding part of one C-step is computation of the OLS fit for the subset \mathbf{m}_i and then the calculation of n absolute residuals. As we stated above, convergence is usually achieved under 20 steps. So for fast algorithm run, we would like to repeat C-step as little as possible and at the same time do not lose the performance of the algorithm.

Since the convergence of the Repeat-C-step algorithm is very fast, it turns out that we can distinguish between starts that leads to good solutions and those which does not after few C-steps iterations. Based on empiric observation, we can distinguish good or bad solution already after two or three iterations of C-steps based on the values $\text{OF}_D^{\text{LTS}}(\hat{\mathbf{w}}_3)$ and $\text{OF}_D^{\text{LTS}}(\hat{\mathbf{w}}_4)$ respectively (see Figure 2.2).

So even though authors do not specify the size of k explicitly, they propose that after a few C-steps we can choose (say 10) best solutions among all \mathbf{m}_1 starting subsets and continue with iterating the C-steps only using the best solutions. Authors refer to this process as to the *selective iteration*.

Nested extension

C-step computation is usually very fast for small n . It gets slow for very high n , say $n > 10^3$, because we need to compute the OLS estimate for the \mathbf{m}_i subset of size h which proportional to n and then calculate n absolute residuals.

Authors came up with a solution they call a *nested extension*. Let k be number initial subsets \mathbf{m}_1 . Then we can describe the nested extension as follows.

- If n is greater than a given limit l , we create subset L of data $|L| = l$, and divide this subset into s disjoint sets $P_1, P_2, \dots, P_s, |P_i| = \frac{l}{s}, P_i \cap P_j = \emptyset, \bigcup_{i=1}^s P_i = L$.
- For every P_i we set the number of starts $k_{P_i} = \frac{k}{l}$.
- Next in every P_i we create k_{P_i} number of initial $\mathbf{m}_{P_{i_1}}$ subsets and apply C-steps few times for each of them.
- Choose 10 best results from each subsets and merge them together. We get family of sets F_{merged} containing 10 best $\mathbf{m}_{P_{i_3}}$ subsets from each P_i .

- On each subset from F_{merged} family of subsets we again apply 2 C-steps twice and then choose 10 best results.
- Finally we use these 10 best subsets and apply Repeat-C-step algorithm.
- As a result we choose the best of those 10 results.

2.2.4 Putting all together

We have described all major parts of the algorithm FAST-LTS. One last thing we need to mention is that even though C-steps iteration usually converges under 20 steps, it is appropriate to introduce two parameters i_{max} and r which limits the number of C-steps iterations in some rare cases when convergence is too slow. Parameter i_{max} denotes the maximum number of iterations in the final Repeat-C-step iteration. Parameter r denotes the threshold for the stopping criterion because of the rounding errors we use $|\text{OF}_D^{\text{LTS}}(\hat{\mathbf{w}}_i) - \text{OF}_D^{\text{LTS}}(\hat{\mathbf{w}}_{i+1})| \leq r$ instead of $\text{OF}_D^{\text{LTS}}(\hat{\mathbf{w}}_i) = \text{OF}_D^{\text{LTS}}(\hat{\mathbf{w}}_{i+1})$.

2.3 Feasible solution

In this section we introduce feasible solution algorithm from [8]. It is based on the strong necessary condition described at Definition 8. The basic idea can be described as follows.

Let us consider that we have some $\mathbf{m} \in Q^{(n,h)}$, then we can denote $O_m = \{i \in \{1, 2, \dots, n\}; w_i = 1\}$ and $Z_m = \{j \in \{1, 2, \dots, n\}; w_j = 0\}$ thus sets of indexes of positions where is 0 respectively 1 in vector \mathbf{m} . We can think about it as indexes of observations in h set and g set respectively. Then we can mark $\mathbf{m}^{(i,j)}$ as a vector which is constructed by swap of its i th and j th element where $i \in O$ and $j \in Z$. Such vector corresponds to vector \mathbf{m}_{swap} which we marked at Definition 8.

With this in mind, let us mark

$$\Delta S_{i,j}^{(m)} = \text{OF}_D^{\text{LTS}}(\mathbf{m}^{(i,j)}) - \text{OF}_D^{\text{LTS}}(\mathbf{m}) \quad (2.38)$$

thus a change of the LTS objective function by swapping one observation from not trimmed subset with another from trimmed subset. To calculate this we can obviously first calculate the $\text{OF}_D^{\text{LTS}}(\mathbf{m})$ and the $\text{OF}_D^{\text{LTS}}(\mathbf{m}^{(i,j)})$ and finally subtract both results. Although it is a option, it is computationally exhaustive. So the question is if there is easier way of calculating $\Delta S_{i,j}^{(m)}$. The answer is positive and we describe it now.

Let us mark $M = \text{diag}(\mathbf{m})$ and $M^{(i,j)} = \text{diag}(\mathbf{m}^{(i,j)})$ and also $\mathbf{Z}_M = (\mathbf{X}^T \mathbf{M}^T \mathbf{X})$. For now let us assume that we also have \mathbf{Z}_M^{-1} and $\hat{\mathbf{w}} = \mathbf{Z}_M^{-1} \mathbf{X}^T \mathbf{M} \mathbf{y}$. We now want to calculate $\Delta S_{i,j}^{(m)}$. Let us mark vector of residuals $\mathbf{r}^{(m)} =$

2. ALGORITHMS

$\mathbf{Y} - \mathbf{X}\hat{\mathbf{w}}$ and also $d_{r,s} = \mathbf{x}_r \mathbf{Z} \mathbf{M}^{-1} \mathbf{x}_s$ then by equation introduced in [21] we get

$$\Delta S_{i,j}^{(\mathbf{m})} = \frac{(r_j^{(\mathbf{m})})^2(1 - d_{i,i}) - (r_i^{(\mathbf{m})})^2(1 + d_{j,j}) + 2r_i^{(\mathbf{m})}r_j^{(\mathbf{m})}d_{j,j}}{(1 - d_{i,i})(1 + d_{j,j}) + d_{i,j}^2}. \quad (2.39)$$

Let us now describe the core of the algorithm. It is similar to the FAST-LTS algorithm in the sense of iterative refinement of a h -element subset. So let us assume that we have some vector $\mathbf{m} \in Q^{(n,h)}$. Now we compute $\Delta S_{i,j}^{(\mathbf{m})}$ for all $i \in O$ and $j \in Z$. This may lead to several different outcomes:

1. all $S_{i,j}^{(\mathbf{m})}$ are non-negative
2. one $S_{i,j}^{(\mathbf{m})}$ is negative
3. multiple $S_{i,j}^{(\mathbf{m})}$ are negative

In the first case, all $\text{OF}_D^{\text{LTS}}(\mathbf{m}^{(i,j)})$ are higher or the same as the $\text{OF}_D^{\text{LTS}}(\mathbf{m})$ so none swap will lead to an improvement. That also means that strong necessary condition is satisfied and the algorithm ends.

In the second and third case, the strong necessary condition is not satisfied, and we can make the swap. In the second case, it is easy which one to choose because we have only one. In the third case, we have a couple of options again:

1. use the first swap that leads to the improvement
2. from all possible swaps choose one that has highest improvement value $\text{OF}_D^{\text{LTS}}(\mathbf{m}^{(i,j)})$
3. use the first swap that has improvement higher than some given threshold

In terms of complexity, all three options give the same results, because to find a feasible solution you need to evaluate all pair swaps. In practice, the third option is the winner because it leads to the least amount of iterations. On the other hand, as we said this does not improve the complexity of the algorithm, so from now on let us assume that we use the case number two. So if there negative $S_{i,j}^{(\mathbf{m})}$, we choose the one with the lowest value, make the swap and repeat the process.

The algorithm ends when there is no possible improvement, i.e., when all $S_{i,j}^{(\mathbf{m})}$ are non-negative. The number of iterations needed till algorithm stops is usually quite low, but for practical usage, it is still convenient to use some parameter i_{max} after which algorithm stops without finding h -element subset satisfying the strong necessary condition. One step of this algorithm

we call optimal swap additive algorithm (OSAA) is described by the following pseudocode.

Algorithm 4: OSAA

Input: $\mathbf{Z}_M^{-1} \in \mathbb{R}^{p \times p}$, $\mathbf{r}^{(m)} \in \mathbb{R}^{n \times 1}$, O_m , Z_m , $\mathbf{X} \in \mathbb{R}^{n \times p}$
Output: $\hat{\mathbf{w}}_{new}$, m_{new}

```

1  $S \leftarrow 0$ ;
2  $i_{swap} \leftarrow \emptyset$ ;
3  $j_{swap} \leftarrow \emptyset$ ;
4 for  $m_i \in O_m$  do
5   for  $m_j \in Z_m$  do
6      $r_i^{(m)} = \mathbf{r}_{m_i}^{(m)}$ ;
7      $r_j^{(m)} = \mathbf{r}_{m_j}^{(m)}$ ;
8      $d_{i,i} = \mathbf{x}_{m_i} \mathbf{Z}_M^{-1} \mathbf{x}_{m_i}^T$ ;
9      $d_{i,j} = \mathbf{x}_{m_i} \mathbf{Z}_M^{-1} \mathbf{x}_{m_j}^T$ ;
10     $d_{j,j} = \mathbf{x}_{m_j} \mathbf{Z}_M^{-1} \mathbf{x}_{m_j}^T$ ;
11     $S_{tmp} = \text{calculate } \Delta S_{i,j}^{(m)} \text{ by (2.39)}$ ;
12    if  $S_{tmp} < S$  then
13       $S \leftarrow S_{tmp}$ ;
14       $i_{swap} \leftarrow m_i$ ;
15       $j_{swap} \leftarrow m_j$ ;
16    end
17  end
18 end
19 return  $i_{swap}, j_{swap}, S$ ;
```

Observation 20. The time complexity of the OSAA 4 is $\mathcal{O}(n^2 p^2)$

Proof. All $d_{i,i}$ and $d_{j,j}$ can be calculated before the for loops. To calculate $d_{i,i}$ resp. $d_{j,j}$ we need to multiply vector $\in \mathbb{R}^p$ with matrix $\in \mathbb{R}^{p \times p}$ and vector $\in \mathbb{R}^p$ that is $\mathcal{O}(p^2)$. For all $d_{i,i}$ this has to be done h times and for $d_{j,j}$ $n - h$ times. So all together it is $\mathcal{O}(np^2)$.

The two loops go through all pairs; thus it is $\mathcal{O}(n^2)$. $d_{i,j}$ can be calculated in the loop, and it can be done in $\mathcal{O}(p^2)$.

If we put everything together we get $\mathcal{O}(np^2 + n^2 p^2) \sim \mathcal{O}(n^2 p^2)$. \square

One run of this iteration process leads to some local optimum, i.e., set satisfying the strong necessary condition. In [8] they refer to this set as to *feasible set*. This is because it is not global optima the algorithm needs to be run multiple times say t times. A h subset with the smallest value of the objective function is chosen as the final solution.

Discussion about how to find the initial h -element subset resp. initial \mathbf{m} was already discussed when describing FAST-LTS algorithm. More importantly, as we already suggested h -element subset satisfying weak necessary

2. ALGORITHMS

condition do not need to satisfy strong necessary condition so passing such a h -element subset as input to this algorithm is another option and we discuss it in detail later. We now describe a feasible solution algorithm (FSA) with pseudocode, and we assume that we already have some function that generates for us h -element subsets, e.g., random one.

Algorithm 5: FSA

Input: $\mathbf{X} \in \mathbb{R}^{n \times p}, \mathbf{y} \in \mathbb{R}^{n \times 1}, i_{\max}, t$
Output: $\hat{\mathbf{w}}_{\text{final}}, \mathbf{m}_{\text{final}}$

```

1  $\hat{\mathbf{w}}_{\text{final}} \leftarrow \emptyset;$ 
2  $\mathbf{m}_{\text{final}} \leftarrow \emptyset;$ 
3  $R \leftarrow \emptyset;$ 
4 for  $k \leftarrow 0$  to  $t$  do
5    $\mathbf{m} \leftarrow \text{generate\_intial\_subset}();$            // e.g. random  $\mathbf{m} \in Q^{(n,h)}$ 
6    $l \leftarrow 0;$ 
7   while True do
8      $\mathbf{M} \leftarrow \text{diag}(\mathbf{m});$ 
9      $\mathbf{Z}_M = (\mathbf{X}^T \mathbf{M}^T \mathbf{X});$ 
10     $\mathbf{Z}_M^{-1} = \text{calculate inversion of } \mathbf{Z}_M;$ 
11     $\hat{\mathbf{w}} \leftarrow \text{OF}^{(OLS, \mathbf{M}\mathbf{X}, \mathbf{M}\mathbf{y})};$ 
12     $\mathbf{r}^{(m)} = \mathbf{Y} - \mathbf{X}\hat{\mathbf{w}};$ 
13     $S_{i,j}, i, j = \text{OSAA}(\mathbf{Z}_M^{-1}, \mathbf{r}^{(m)}, O_m, \mathbf{Z}_m, \mathbf{X});$ 
14    if  $S_{i,j} \geq 0$  or  $l \geq i_{\max}$  then
15       $\mathbf{m}_{\text{final}} \leftarrow \mathbf{m};$ 
16       $\hat{\mathbf{w}} \leftarrow \text{OF}^{(OLS, \mathbf{M}\mathbf{X}, \mathbf{M}\mathbf{y})};$ 
17       $R \leftarrow R \cup \{\mathbf{m}_{\text{final}}, \hat{\mathbf{w}}\};$ 
18      break;
19    end
20    else
21       $\mathbf{m} \leftarrow \mathbf{m}^{(i,j)};$ 
22    end
23  end
24 end
25  $\mathbf{m}_{\text{final}}, \hat{\mathbf{w}}_{\text{final}} \leftarrow \text{select best from } R \text{ based on smallest value of } \text{OF}_D^{\text{LTS}};$ 
26 return  $\hat{\mathbf{w}}_{\text{final}}, \mathbf{m}_{\text{final}};$ 

```

Observation 21. In the main loop beside running OSAA which time complexity is $\mathcal{O}(n^2 p^2)$ we need to recalculate $\hat{\mathbf{w}}$ using the matrix inversion which time complexity is $\mathcal{O}(p^2 n)$ (see Observation 11). Main loop of the FSA runs up to l iterations for each start t . So the time complexity of whole algorithm is $\mathcal{O}(lt(n^2 p^2 + p^2 n))$. Because l and t are usually quite low, we can see that the FSA time complexity is dominated by the OSAA.

In this section, we've described the FSA algorithm. In the next section, we

introduce a very similar algorithm, but which have higher numerical stability and also contains various optimizations for higher performance.

2.4 OEA, MOEA, MMEA

In the FSA algorithm, we assumed that after each cycle of OSAA 4 we need to recalculate inversion $(\mathbf{X}^T \mathbf{X})$ together with $\hat{\mathbf{w}}$. In this section, we introduce a different approach described in [5] so that we can update it instead of recalculating. Moreover, these ideas lead to bounding condition of FSA which not only improves the speed but also gives options to construct different algorithms.

In Section 2.3 we introduce additive formula (2.39) which corresponds to $\text{OF}_D^{\text{LTS}}(\mathbf{m}^{(i,j)}) - \text{OF}_D^{\text{LTS}}(\mathbf{m})$. Let us now try to obtain similar formula but with the focus on how the individual elements in our current algorithm change namely the inversion \mathbf{Z}^{-1} and $\hat{\mathbf{w}}$. We also split the idea of swap the i, j into the insertion of j and remove of i . Thus we describe how individual elements are affected after adding one row and also removing one row. Moreover, during this derivation, we obtain two different approaches to calculate it. Both are important, and we recapitulate them after.

2.4.1 Multiplicative formula

$\mathbf{Z} = \mathbf{X}^T \mathbf{X}$ Let us denote $\mathbf{A} = (\mathbf{X}, \mathbf{y})$, a matrix \mathbf{X} expanded by one column of corresponding dependent variables and and $\tilde{\mathbf{Z}} = \mathbf{A}^T \mathbf{A}$. Then

$$\tilde{\mathbf{Z}} = \begin{bmatrix} \mathbf{Z} & \mathbf{X}^T \mathbf{y} \\ \mathbf{y}^T \mathbf{X} & \mathbf{y}^T \mathbf{y} \end{bmatrix}. \quad (2.40)$$

Notice that \mathbf{Z} is symmetric square matrix $\in \mathbb{R}^{p \times p}$, moreover we suppose that \mathbf{Z} is regular. $\mathbf{X}^T \mathbf{y}$ is p dimensional column vector, $\mathbf{y}^T \mathbf{X}$ is p dimensional row vector and $\mathbf{y}^T \mathbf{y}$ is a scalar. OLS estimate $\hat{\mathbf{w}}$ is then given by

$$\hat{\mathbf{w}}^{(OLS,n)} = \mathbf{Z}^{-1} \mathbf{X}^T \mathbf{y} \quad (2.41)$$

and the RSS by

$$RSS = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{\mathbf{w}}. \quad (2.42)$$

Also, note that all necessary matrix multiplications are already in the matrix $\tilde{\mathbf{Z}}$.

Let us show that RSS can be expressed as ratio of determinants $\det(\tilde{\mathbf{Z}})$ and $\det(\mathbf{Z})$ (using determinant rule for block matrices) so that

$$\begin{aligned}
 RSS &= \frac{\det(\tilde{\mathbf{Z}})}{\det(\mathbf{Z})} \\
 &= \frac{\det \begin{pmatrix} \mathbf{Z} & \mathbf{X}^T \mathbf{y} \\ \mathbf{y}^T \mathbf{X} & \mathbf{y}^T \mathbf{y} \end{pmatrix}}{\det(\mathbf{M})} \\
 &= \frac{\det(\mathbf{Z}) \det(\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{Z}^{-1} \mathbf{y})}{\det(\mathbf{Z})} \\
 &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{\mathbf{w}}.
 \end{aligned} \tag{2.43}$$

If we assume that $RSS > 0$ then $\tilde{\mathbf{Z}}^{-1}$ can be expressed as

$$\tilde{\mathbf{Z}}^{-1} = \begin{bmatrix} \mathbf{Z}^{-1} + \frac{\hat{\mathbf{w}} \hat{\mathbf{w}}^T}{RSS} & -\frac{\hat{\mathbf{w}}^T}{RSS} \\ -\frac{\hat{\mathbf{w}}^T}{RSS} & \frac{1}{RSS} \end{bmatrix}. \tag{2.44}$$

For following equations it is important to notice that for any two row vectors $\mathbf{c}_i = (\mathbf{x}_i, y_i)$ and $\mathbf{c}_j = (\mathbf{x}_j, y_j)$, $\mathbf{c}_i, \mathbf{c}_j \in \mathbb{R}^{1 \times p+1}$ is

$$\mathbf{c}_i \tilde{\mathbf{Z}}^{-1} \mathbf{c}_i^T = \frac{(y_i - \mathbf{x}_i \hat{\mathbf{w}})^2}{RSS} + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T \tag{2.45}$$

and

$$\mathbf{c}_j \tilde{\mathbf{Z}}^{-1} \mathbf{c}_j^T = \frac{(y_j - \mathbf{x}_j \hat{\mathbf{w}})(y_i - \mathbf{x}_i \hat{\mathbf{w}})}{RSS} + \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T \tag{2.46}$$

Including the observation

Using above let us express how the determinant $\det(\mathbf{Z})$ and the inverse of the $\text{vec} \mathbf{Z}^{-1}$ changes when observation $\mathbf{c}_i = (\mathbf{x}_i, y_i)$ is added to the matrix \mathbf{A} . First let us notice that if we add this row to \mathbf{A} , then \mathbf{Z} changes as

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} & x_{i_1} \\ x_{21} & x_{22} & \dots & x_{2n} & x_{i_2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{p1} & x_{p2} & \dots & x_{pn} & x_{i_p} \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \\ x_{i_1} & x_{i_2} & \dots & x_{i_p} \end{bmatrix} = \mathbf{X}^T \mathbf{X} + \mathbf{x}_i^T \mathbf{x}_i = \mathbf{Z} + \mathbf{x}_i^T \mathbf{x}_i, \tag{2.47}$$

so determinant with appended row changes as

$$\det(\mathbf{Z} + \mathbf{x}_i^T \mathbf{x}_i) = \det(\mathbf{Z})(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T). \tag{2.48}$$

Finally the inversion \mathbf{Z}^{-1} can be obtained using Sherman-Morrison formula [22] so that

$$(\mathbf{Z} + \mathbf{x}_i^T \mathbf{x}_i)^{-1} = \mathbf{Z}^{-1} - \frac{\mathbf{Z}^{-1} \mathbf{x}_i^T \mathbf{x}_i \mathbf{Z}^{-1}}{1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T} \quad (2.49)$$

It is now convenient to denote

$$b = \frac{-1}{(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T)}, b \in \mathbb{R}, \quad (2.50)$$

and

$$\mathbf{u} = \mathbf{Z}^{-1} \mathbf{x}_i^T, \mathbf{u} \in \mathbb{R}^{p \times 1}. \quad (2.51)$$

Then (2.49) can be written as

$$(\mathbf{Z} + \mathbf{x}_i^T \mathbf{x}_i)^{-1} = \mathbf{Z}^{-1} + b \mathbf{u} \mathbf{u}^T. \quad (2.52)$$

Given that last missing piece to express updated $\hat{\mathbf{w}}$ is appended $\mathbf{X}^T \mathbf{y}$ by one row, which can be simply expressed by same idea as (2.47) so that updated $\hat{\mathbf{w}}$ which we denote as $\bar{\mathbf{w}}$ is

$$\bar{\mathbf{w}} = (\mathbf{Z}^{-1} + b \mathbf{u} \mathbf{u}^T)(\mathbf{X}^T \mathbf{y} + y_i \mathbf{x}_i^T). \quad (2.53)$$

This can be simplified so that we get

$$\bar{\mathbf{w}} = \hat{\mathbf{w}} - (y_i - \mathbf{x}_i \hat{\mathbf{w}}) b \mathbf{u}. \quad (2.54)$$

mention the mistake in original paper? $\mathbf{w} + (\mathbf{y} - \mathbf{x}\mathbf{w})b\mathbf{u}$??

Least but not last we want to express updated RSS which we denote as \overline{RSS} . This can be done easily from (2.43) and (2.48) as

$$\overline{RSS} = RSS + \frac{(y_i - \mathbf{x}_i \hat{\mathbf{w}})^2}{(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T)}. \quad (2.55)$$

It is convenient to mark

$$\gamma^+(\mathbf{c}_i) = \frac{(y_i - \mathbf{x}_i \hat{\mathbf{w}})^2}{(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T)}, \quad (2.56)$$

so that

$$\overline{RSS} = RSS + \gamma^+(\mathbf{c}_i) \quad (2.57)$$

We can see that $\gamma^+(\mathbf{c}_i)$ measures how RSS increase after we append our dataset with observation \mathbf{c}_i , thus $\gamma^+(\mathbf{c}_i) \geq 0$

Excluding the observation

Because we want to express both increment and decrement change in our dataset, let us now focus how RSS , \mathbf{Z}^{-1} and $\hat{\mathbf{w}}$ changes after we exclude one observation.

Consider that we already included one observation \mathbf{c}_i in our dataset and mark $\bar{\mathbf{Z}} = \mathbf{Z} + \mathbf{x}_i \mathbf{x}_i^T$. If we exclude one observation $\mathbf{c}_j = (\mathbf{x}_j, y_j) \in \mathbb{R}^{p+1 \times 1}$ from already updated matrix \mathbf{A} then the determinant $\det(\bar{\mathbf{Z}})$ changes as

$$\det(\bar{\mathbf{Z}} - \mathbf{x}_j^T \mathbf{x}_j) = \det(\bar{\mathbf{Z}})(1 - \mathbf{x}_j \bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T) \quad (2.58)$$

and the inversion changes (again, according to Sherman-Morrison formula) as

$$(\bar{\mathbf{Z}} - \mathbf{x}_j^T \mathbf{x}_j)^{-1} = \bar{\mathbf{Z}}^{-1} + \frac{\bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T \mathbf{x}_j \bar{\mathbf{Z}}^{-1}}{1 - \mathbf{x}_j \bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T}. \quad (2.59)$$

Once again, it is convenient to denote

$$\bar{b} = \frac{-1}{(1 \mathbf{x}_j \bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T)}, \bar{b} \in \mathbb{R}, \quad (2.60)$$

and

$$\bar{\mathbf{u}} = \bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T, \bar{\mathbf{u}} \in \mathbb{R}^{p \times 1}, \quad (2.61)$$

so that we can write

$$(\bar{\mathbf{Z}} - \mathbf{x}_j^T \mathbf{x}_j)^{-1} = \bar{\mathbf{Z}}^{-1} - \bar{b} \bar{\mathbf{u}} \bar{\mathbf{u}}^T. \quad (2.62)$$

Using the same approach as before we can denote express down-dated estimate which we denote as $\bar{\hat{\mathbf{w}}}$

$$\bar{\hat{\mathbf{w}}} = +\hat{\mathbf{w}}(y_j - \mathbf{x}_j \hat{\mathbf{w}}) \bar{b} \bar{\mathbf{u}}. \quad (2.63)$$

Finally lets also express updated \overline{RSS} which we denote as $\overline{\overline{RSS}}$. This can be done easily from (2.43) and (2.48) and (2.62) as

$$\overline{\overline{RSS}} = \overline{RSS} - \frac{(y_j - \mathbf{x}_j \bar{\hat{\mathbf{w}}})^2}{(1 - \mathbf{x}_j \bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T)}. \quad (2.64)$$

Let us mark

$$\gamma^-(\mathbf{c}_j) = \frac{(y_j - \mathbf{x}_j \bar{\hat{\mathbf{w}}})^2}{(1 - \mathbf{x}_j \bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T)}, \quad (2.65)$$

so that

$$\overline{\overline{RSS}} = \overline{RSS} - \gamma^-(\mathbf{c}_j). \quad (2.66)$$

Swapping two observations

Let us now express the equation for including and excluding observation at once. First, notice that from (2.62) we can express

$$\mathbf{x}_j \bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T = \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T - \frac{(\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T)^2}{1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T}, \quad (2.67)$$

so that

$$\begin{aligned} \det(\mathbf{Z} + \mathbf{x}_i^T \mathbf{x}_i - \mathbf{x}_j^T \mathbf{x}_j) = \\ \det(\mathbf{Z})(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T + (\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T)^2 - \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T) \end{aligned} \quad (2.68)$$

Finally we can express the \overline{RSS} as

$$\overline{RSS} = RSS \rho(\mathbf{c}_i, \mathbf{c}_j), \quad (2.69)$$

where

$$\rho(\mathbf{c}_i, \mathbf{c}_j) = \frac{(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T + \frac{e_j^2}{RSS})(1 - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T - \frac{e_i^2}{RSS}) + (\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T + \frac{e_i e_j}{RSS})^2}{1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T + (\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T)^2 - \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T}, \quad (2.70)$$

where $e_i = y_i - \mathbf{x}_i^T \hat{\mathbf{w}}$ and $e_j = y_j - \mathbf{x}_j^T \hat{\mathbf{w}}$.

We can see that this formula is similar to the (2.39) but here the $\rho(\mathbf{c}_i, \mathbf{c}_j)$ represents multiplicative increment. Moreover $0 < \rho(\mathbf{c}_i, \mathbf{c}_j)$ and if $0 < \rho(\mathbf{c}_i, \mathbf{c}_j) < 1$ then the swap leads to improvement in terms of feasible solution.

We are now able to modify the FSA so that we do not need to recompute $\hat{\mathbf{w}}$ and inversion $(\mathbf{X}^T \mathbf{X})^{-1}$ but we can update it. Authors call this algorithm optimal exchange algorithm (OEA)

2.4.2 The OEA and its properties

We can apply theory from the previous section to the FSA. Sets O_m and Z_m (see Section 2.3) contains indexes of observations to exclude and include respectively. The matrix \mathbf{Z}_M can be used instead of the matrix \mathbf{Z} and \mathbf{X}_M with \mathbf{y}_M to calculate $\hat{\mathbf{w}}$.

In terms of the FSA 5 we first need to change the OSAA. There are only minor tweaks. First, we need to pass one more argument, and that is RSS . Second, We want to find minimal $\rho(\mathbf{c}_i, \mathbf{c}_j)$ calculated by (2.70) so that $0 < \rho(\mathbf{c}_i, \mathbf{c}_j) < 1$. Other parts of the algorithm remains the same. Time complexity thus remains the same.

This algorithm produces $\rho(\mathbf{c}_i)$, and indexes $i_{swap} \in O_m$ and $j_{swap} \in Z_m$ of observations we want to swap. Given that, we can update the RSS by (2.69), \mathbf{Z}_M^{-1} by (2.52) and (2.62) and finally update $\hat{\mathbf{w}}$ by (2.54) and (2.63).

Let us now talk about the time complexity of such a solution. As we said, the time complexity of modifies OSAA is $\mathcal{O}(n^2p^2)$. Thus there is no asymptotic improvement. On the other hand, there are some significant constants outside the OSAA so let us look at how they improve.

When an improvement is found, thus when $0 < \rho(\mathbf{c}_i, \mathbf{c}_j) < 1$, then we update RSS , which is constant. Next we update inversion by (2.52) which is $\mathcal{O}(4p^2)$ and (2.62) which is also $\mathcal{O}(4p^2)$. Finally we need to update $\hat{\mathbf{w}}$ by (2.54) and (2.63) which are both $\mathcal{O}(p^2 + p)$.

Observation 22. Time complexity of updating all the quantities is $\mathcal{O}(8p^2 + 2p^2 + p) \sim \mathcal{O}(p^2)$. That is quite an improvement if we compare it to time complexity $\mathcal{O}(p^2n)$ of updating those quantities in the FSA (see Observation 21).

Note that right now it actually does not matter if we use additive formula (2.39) with the stopping criterion $\Delta S_{i,j}^{(m)} \geq 0$ — thus unmodified OSAA — or multiplicative formula (2.70) with the stopping criterion $\rho(\mathbf{c}_i, \mathbf{c}_j) \geq 1$. Both results we can update RSS and other quantities can be used to calculate both formulas.

However, the advantage of the multiplicative formula (2.70) is that we can use the following bounding condition to improve the performance of the modified OSAA.

Bounding condition improvement

The $\rho(\mathbf{c}_i, \mathbf{c}_j)$ is expressed as a ratio. We can see that in the numerator we have

$$(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T + \frac{e_j^2}{RSS})(1 - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T - \frac{e_i^2}{RSS}) + (\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T + \frac{e_i e_j}{RSS})^2, \quad (2.71)$$

and because $\frac{e_i e_j}{RSS})^2 \geq 0$ then whole numerator is greater or equal to

$$(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T + \frac{e_j^2}{RSS})(1 - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T - \frac{e_i^2}{RSS}). \quad (2.72)$$

On the other hand, we can see that denominator is

$$1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T + (\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T)^2 - \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T, \quad (2.73)$$

and because $(\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T)^2$ and $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T$ are actually inner products of \mathbf{x}_i and \mathbf{x}_j (\mathbf{Z}^{-1} is positive definite) thus

$$(\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T)^2 \leq \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T \quad (2.74)$$

using the Cauchy-Schwarz inequality. This means that denominator is less or equal to

$$1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T. \quad (2.75)$$

Given that we can denote $\rho_b(\mathbf{c}_i, \mathbf{c}_j)$ as

$$\rho_b(\mathbf{c}_i, \mathbf{c}_j) = \frac{(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T + \frac{e_j^2}{RSS})(1 - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T - \frac{e_i^2}{RSS})}{1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T} \leq \rho(\mathbf{c}_i, \mathbf{c}_j). \quad (2.76)$$

The actual speed improvement is then given by that we do not need to compute $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T$ in each of $h(n-h)$ pairs swap comparison. As we know this quantity cannot be computed outside the loop; thus it is the reason for such a high time complexity. The modified OSAA can be further modified as follows.

First we set $\rho_{min} := 1$. Then for each pair we only compute $\rho_b(\mathbf{c}_i, \mathbf{c}_j)$ and if it is greater or equal to ρ_{min} we can continue to next pair without computing $\rho(\mathbf{c}_i, \mathbf{c}_j)$. It is very useful because all quantities necessary for calculating $\rho_b(\mathbf{c}_i, \mathbf{c}_j)$ can be precalculated outside of the loop.

If $\rho_b(\mathbf{c}_i, \mathbf{c}_j)$ is less than ρ_{min} we actually compute $\rho(\mathbf{c}_i, \mathbf{c}_j)$ and set $\rho_{min} := \rho(\mathbf{c}_i, \mathbf{c}_j)$.

That means in the double loop which time complexity is $\mathcal{O}(n^2)$ we do not always need to calculate $\rho(\mathbf{c}_i, \mathbf{c}_j)$ which time complexity is $\mathcal{O}(n^2)$. This does not improve the asymptotic time complexity, but (as can be seen in the table) it can improve the speed of the algorithm.

create experiment

Finally, let us note that in [5] authors call this algorithm with bounding condition as *modified optimal exchange algorithm* (MOEA).

2.4.3 Minimum-maximum exchange algorithm

The minimum-maximum exchange algorithm (MMEA) is very similar to the FSA respectively to its modified version the OEA. The main difference is the greediness of this algorithm. What we mean by that is that this algorithm does not find the optimal swap, but rather first find the \mathbf{c}_j which increase the RSS by the minimum value and include this observation. Next, it finds \mathbf{c}_j such that in decrease RSS by maximum and exclude this observation.

The minimum increase can be found by calculating $\gamma^+(\mathbf{c}_i)$ by (2.56) for each trimmed observation in Z_m .

Next, this observation is included, so we get $h+1$ untrimmed observations. We update \mathbf{Z}_M^{-1} to $\overline{\mathbf{Z}}_M^{-1}$ by (2.52) and $\hat{\mathbf{w}}$ to $\overline{\hat{\mathbf{w}}}$ by (2.54).

Then we can find maximum $\gamma^-(\mathbf{c}_j)$ by (2.65) among $O_m^{(i)}$ where $m^{(i)}$ denotes included observation \mathbf{c}_i .

Next, we can update $\overline{\mathbf{Z}}_M^{-1}$ to the $\overline{\overline{\mathbf{Z}}}_M^{-1}$ by (2.62) and $\overline{\hat{\mathbf{w}}}$ to $\overline{\overline{\hat{\mathbf{w}}}}$ by (2.63).

Finally, we can update RSS to $\overline{\overline{RSS}}$ by (2.57) and (2.66). We can repeat this process until $\gamma^-(\mathbf{c}_j) > \gamma^+(\mathbf{c}_i)$.

Observation 23. One step of the algorithm MMEA has time complexity $\mathcal{O}(p^2n)$. So the whole algorithm has time complexity $\mathcal{O}(tlp^2n)$ where t is number of the starts and l is maximum number of iterations of each start.

Proof. Because we do not iterate through all pairs but only over $n-h$ followed by $h+1$ subsets the loops takes only $\mathcal{O}(n)$ time. In the loops, we are computing $\gamma^-(\mathbf{c}_j)$ and $\gamma^+(\mathbf{c}_i)$ both takes $\mathcal{O}(p^2)$ time. Outside the loops, all the each of the quantities we are updating take $\mathcal{O}(p^2)$ time. That means the one loop of the whole algorithm last $\mathcal{O}(p^2n)$. As in the case of the FSA, if we introduce parameters t and l , then the time complexity of the algorithm is $\mathcal{O}(tlp^2n)$ \square

2.4.4 Different method of computation

In the last section we introduced a way of calculating the OEA so that we update $\hat{\mathbf{w}}$ and inversion $(\mathbf{X}^T \mathbf{X})^{-1}$. This, however, requires to compute inversion at the start of the algorithm (this is also the case for the FSA, MOEA and MMEA) As we know, calculating inversion is not practical due to low numerical stability. In practice, we usually use QR decomposition. In this section, we describe how we can modify OEA to use the QR decomposition (the same idea can also be applied to the FSA, MOEA and the MMEA).

Let us start by describing how we can update the QR factorization, which is a critical part of this modified computation. Assume that we have QR decomposition of \mathbf{X} , and we need to exchange i th observation from O_m with j th observation from Z_m . We can simulate this by inserting j th row and consequently extracting i th row from the QR decomposition.

QR insert

First, let us discuss how to update QR decomposition when a row \mathbf{x}_j is inserted. If we add a row to \mathbf{A} as the last row, our decomposition looks like

$$\bar{\mathbf{R}} = \bar{\mathbf{Q}} \mathbf{A}^{(+)} = \begin{bmatrix} \times & \cdots & \cdots & \cdots & \times \\ 0 & \times & \cdots & \cdots & \times \\ \vdots & 0 & \ddots & \cdots & \vdots \\ \vdots & \vdots & 0 & \ddots & \vdots \\ \vdots & \vdots & \vdots & 0 & \times \\ \vdots & \vdots & \vdots & \vdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \times & \cdots & \cdots & \cdots & \times \end{bmatrix} \quad (2.77)$$

Where $\bar{\mathbf{R}}$ denotes matrix \mathbf{R} which is now not upper triangular and needs to be updated and $\bar{\mathbf{Q}} \in \mathbb{R}^{p+1 \times p+1}$ denotes matrix created from \mathbf{Q} by adding one row and one column with zeros. To preserve zeroes on the diagonal we in addition

put $\overline{Q}_{n+1,n+1} = 1$. To update \overline{R} to upper triangular matrix $R^{(+)}$, we need to create additional orthogonal Givens matrices $Q_{e+1} \dots Q_{e+p}$. Then updated upper triangular matrix is $R^{(+)} = Q_{e+p} \dots Q_{e+2} Q_{e+1} R$.

Updated matrix \overline{Q} which we denote as $Q^{(+)}$ is updated in the same manner, thus $Q^{(+)} = Q Q_e^T Q_{e+1}^T \dots Q_{e+p}^T$. Note that if we do not want to have inserted row x_j as the last but as (say k th) row, then in terms of QR decomposition, we only need to move (last) x_j row to the k th position. In other words, we create a permutation matrix P so that

$$PQ^{(+)} = \begin{bmatrix} A(1 : k - 1, 1 : p) \\ x_j \\ A(1 : k + 1, 1 : p) \end{bmatrix} \quad (2.78)$$

where $A(a : b, 1 : p)$ denotes rows of matrix A from a to b . We can describe the whole process by the following pseudocode.

Algorithm 6: QR insert

Input: $Q \in \mathbb{R}^{n \times n}$, $R \in \mathbb{R}^{n \times p}$, $x_j \in \mathbb{R}^p$, k
Output: $Q^{(+)} \in \mathbb{R}^{n+1 \times n+1}$, $R^{(+)} \in \mathbb{R}^{n+1 \times p}$

```

1  $R^{(+)} \leftarrow R$  with appended last row by  $x_j$ ;
2  $Q^{(+)} \leftarrow Q$  with appended last row and last column by zeors;
3  $Q^{(+)}_{n+1,n+1} \leftarrow 1$ ; // i.e.  $Q^{(+)}$  now has 1 on the diagonal
4 for  $i \leftarrow n$  to  $p$  do
5    $Q(i, n+1) \leftarrow$  create Givens matrix  $Q(i, n+1) \in \mathbb{R}^{n+1 \times n+1}$ ;
6    $R^{(+)} \leftarrow Q(i, n+1) R^{(+)}$ ;
7    $Q^{(+)} \leftarrow Q^{(+)} Q^T(i, n+1)$ ;
8 end
9 for  $i \leftarrow n+1$  to  $k$  do
10   $Q^{(+)} \leftarrow Q^{(+)}$  where we swap the  $i$  row with the  $i-1$  row
11 end
12 return  $Q^{(+)}$ ,  $R^{(+)}$ ;
```

Observation 24. Computing $R^{(+)}$ is $\mathcal{O}(p^2)$ and computing $Q^{(+)}$ is $\mathcal{O}(np)$.

Proof. We have loop over p $Q(i, j)$ Givens matrices. We do not need to create those matrices, because by multiplying $R^{(+)}$ or $Q^{(+)}$ with the matrix $Q(i, j)$ only one row of is affected. That means we can simulate this matrix-matrix multiplication only by iterating over those matrices and multiplying elements with $\cos(\varphi)$ and $\sin(\varphi)$ adequately. This means in case of $R^{(+)}$ we are iterating p times over nonzero rows of $R^{(+)}$ thus p rows. That gives us time complexity of $\mathcal{O}(p^2)$. In the case of $Q^{(+)}$ we are iterating p times over columns of $Q^{(+)}$ and that gives us time complexity of $\mathcal{O}(np)$. \square

Note 25. This process can also be done in case we are using an economic version of matrices R and Q thus matrices R_1 and Q_1 (see (2.11)). In such a case $Q_1 \in \mathbb{R}^{p \times p}$ thus then updating Q_1 to $Q_1^{(+)}$ is only $\mathcal{O}(p^2)$.

Note 26. The matrix $\overline{\mathbf{R}}$ can be updated to $\mathbf{R}^{(+)}$ without the presence of matrix \mathbf{Q} . We use this observation in the algorithm described in Section 2.6.

QR delete

When we extract the row \mathbf{x}_i from matrix \mathbf{A} we can use following trick [18]. First, we move such row as the first row of the matrix \mathbf{A} so we create permutation matrix \mathbf{P} so that

$$\mathbf{PA} = \begin{bmatrix} \mathbf{x}_i \\ \mathbf{A}(1:i-1, 1:p) \\ \mathbf{A}(1:i+1, 1:p) \end{bmatrix} = \begin{bmatrix} \mathbf{x}_i \\ \mathbf{A}^{(-)} \end{bmatrix} = \mathbf{PQR} \quad (2.79)$$

where $\mathbf{A}(a:b, 1:p)$ means rows of matrix \mathbf{A} from a to b . We can see that we only need to introduce zeros in the first row \mathbf{q}_1 (except q_{11}) of the matrix \mathbf{Q} . We can do this by $n-1$ matrices $\mathbf{Q}(i, j) \in \mathbb{R}^{n \times n}$ so that

$$\mathbf{Q}(1, 2) \dots \mathbf{Q}(n-1, n) \mathbf{q}_1^T = \begin{bmatrix} 1 \\ 0 \\ \vdots \end{bmatrix} \quad (2.80)$$

To propagate the change into \mathbf{R} we then consequently need to update \mathbf{R} so that

$$\mathbf{Q}(1, 2) \dots \mathbf{Q}(n-1, n) \mathbf{R} = \begin{bmatrix} \mathbf{v} \\ \mathbf{R}^{(-)} \end{bmatrix}. \quad (2.81)$$

The result is then

$$\mathbf{PA} = (\mathbf{PQQ}^T(n-1, n) \dots \mathbf{Q}^T(1, 2))(\mathbf{Q}(1, 2) \dots \mathbf{Q}(n-1, n) \mathbf{R}) = \begin{bmatrix} 1 & 0 \\ 0 & \mathbf{Q}^{(-)} \end{bmatrix} \begin{bmatrix} \times \\ \mathbf{R}^{(-)} \end{bmatrix}, \quad (2.82)$$

so that

$$\mathbf{A}^{(-)} = \mathbf{Q}^{(-)} \mathbf{R}^{(-)} \quad (2.83)$$

Pseudocode is very similar to the algorithm QR Insert 6 so it is not necessary to provide it here.

Observation 27. Time complexity of QR delete is $\mathcal{O}(n^2)$.

Proof. We need to create $n-1$ Givens matrices, and with each of this we need to multiply $\mathbf{Q}^{(-)}$ and $\mathbf{R}^{(-)}$. As we stated in Observation 24 this can be done in $\mathcal{O}(n)$ for each matrix. So together we get $\mathcal{O}(n^2)$. \square

Note 28. In case of QR delete, it is not possible to use the economic version of matrices.

Calculation of OEA using QR decomposition

Given all the required theory above, let us describe the computation. Let us start with (2.70). Here we need inversion to calculate $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T$, $\mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T$ and $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T$. But this can also be done without inversion, only using the decomposition. We can write this equation

$$\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T = \mathbf{v}^T \mathbf{v} \iff \mathbf{x}_i (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{x}_i^T = \mathbf{v}^T \mathbf{v} \quad (2.84)$$

where \mathbf{v} can be obtained by solving lower triangular system

$$\mathbf{R}^T \mathbf{v} = \mathbf{x}_i^T. \quad (2.85)$$

The same can be done with $\mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T$. Least but not last we need to solve $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T$. We can write this

$$\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T = \mathbf{x}_j \mathbf{u} \iff \mathbf{x}_i (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{x}_j^T = \mathbf{x}_j^T \mathbf{u}, \quad (2.86)$$

where column vector \mathbf{u} is defined by (2.51). We can see that

$$\mathbf{u} = \mathbf{Z}^{-1} \mathbf{x}_i^T \iff (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{x}_i^T = \mathbf{u} \quad (2.87)$$

so that \mathbf{u} can be obtained by solving the upper triangular system

$$\mathbf{R} \mathbf{u} = \mathbf{v}, \quad (2.88)$$

where \mathbf{v} is solution of (2.85). Other quantities of (2.70) does not require \mathbf{Z}^{-1} .

Observation 29. We can see that time complexity of all these calculations is the same as in case of calculating with inversion thus $\mathcal{O}(p^2)$.

Proof. Indeed, in the case of the inversion, we are multiplying quantities such as $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T$ which can be done in $\mathcal{O}(p^2)$. In the case of decomposition, we solve this problem by solving the triangular system of p equations. This can also be done in $\mathcal{O}(p^2)$. \square

When optimal exchange $\mathbf{c}_i, \mathbf{c}_j$ is found then we need to update RSS which can be done same way by (2.69). Updating $\hat{\mathbf{w}}$ to $\bar{\hat{\mathbf{w}}}$ by (2.54) requires \mathbf{u} which in this case we calculate by (2.88) and b which requires $\mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T$. This can be done in the same manner as (2.84). Analogous operations can be used to update $\bar{\hat{\mathbf{w}}}$ to $\bar{\bar{\hat{\mathbf{w}}}}$ by means of (2.63). Only thing we need to realize that we can express (2.67) as $\mathbf{x}_j \bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T = \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T + b(\mathbf{u} \mathbf{x}_j^2)$.

On the other hand, we cannot use equation (2.52) and (2.62) for updating the inversion, because we do not have one. So we need to update our QR decomposition somehow. We describe algorithms for updating the QR decomposition in Section 2.4.4. First, we can use the QR Insert algorithm to add row \mathbf{x}_j to the decomposition and consequently QR delete to extract \mathbf{x}_j from the decomposition.

Observation 30. This approach is slower than updating inversion directly. On the other hand, this solution is numerically stable. Updating inversion can be done in $\mathcal{O}(p^2)$ (see Observation 22) while QR insert has $\mathcal{O}(np)$ (see Observation 24) and time complexity and QR delete even $\mathcal{O}(n^2)$ (see Observation 27).

Note 31. Because time complexity of the QR delete is $\mathcal{O}(n^2)$, it is up to consideration if instead of recycling QR decomposition is not worth it to recalculate it from scratch which takes $\mathcal{O}(p^2n)$ (see (2.20)).

Finally let us talk about matrix $\tilde{\mathbf{Z}}$ (2.40) which we used for derivation of our equations. If we use Observation 2.19 then we realize that if we make QR factorization of $\mathbf{A} = (\mathbf{X}, \mathbf{y})$ then

$$\tilde{\mathbf{Z}} = \begin{bmatrix} \mathbf{Z} & \mathbf{X}^T \mathbf{y} \\ \mathbf{y}^T \mathbf{X} & \mathbf{y}^T \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{R}^T & 0 \\ \phi^T & r \end{bmatrix} \begin{bmatrix} \mathbf{R} & \phi \\ 0 & r \end{bmatrix} \quad (2.89)$$

where

$$\tilde{\mathbf{R}} = \begin{bmatrix} \mathbf{R} & \phi \\ 0 & r \end{bmatrix} = \tilde{\mathbf{Q}}^T \mathbf{A} \quad (2.90)$$

is matrix from QR factorization of \mathbf{A} . Next, we can realize that $\mathbf{R} \in \mathbb{R}^{p \times p}$ is matrix \mathbf{R} which we can obtain by QR factorization of \mathbf{X} . Moreover $\phi \in \mathbb{R}^{p \times 1}$ is column vector which is actually equal to

$$\phi = \mathbf{Q}^T \mathbf{y} \quad (2.91)$$

where \mathbf{Q} is matrix \mathbf{Q} from QR factorization of \mathbf{R} . Due to this fact $\hat{\mathbf{w}}$ is solution of upper triangular system

$$\mathbf{R} \hat{\mathbf{w}} = \phi \quad (2.92)$$

Finally $r \in \mathbb{R}$ is scalar such that

$$r^2 = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{\mathbf{w}} = RSS \quad (2.93)$$

Remark 32. We can see that all the quantities we use in the algorithm can be obtained from matrix $\tilde{\mathbf{R}}$.

Now we have described both versions of calculation the algorithm OEA. As we already stated, it easy to use this approach on algorithms FSA, MOEA, and MMEA. Using inversion \mathbf{Z}^{-1} is slightly faster, but numerically less stable, second, using QR decomposition is more numerically stable. Asymptotically both approaches provide the same performance. Finally, we have also shown that using matrix $\tilde{\mathbf{R}}$ is useful because it contains all necessary quantities for the algorithms.

2.5 Combined algorithm

In this section, we shortly describe how we can utilize Lemma 10; thus that strong necessary condition is not satisfied unless the weak necessary condition is satisfied.

Algorithm FAST-LTS outputs h element subset satisfying the weak necessary condition. One step of this algorithm has time complexity $\mathcal{O}(p^2n)$ (see Observation 16). The iteration of this step is quite low and moreover usually limited by parameter.

On the other hand, algorithms for finding strong necessary condition FSA, OEA and MOEA have time complexity of one step $\mathcal{O}(p^2n^2)$ (we now don't take into account eager version MMEA, because we have no proof that it finds h element subset satisfying strong necessary condition).

We can use this in our favor so that we can find h element subset satisfying weak necessary condition by the FAST-LTS algorithm and consequently use this h element subset as an input to some of the algorithms which find h element subset satisfying the strong necessary condition.

Let us consider the combination of the FAST-LTS and the MOEA. We have two options on how to perform this: z

1. Let the FAST-LTS converge and use this h element subset as input to the MOEA which we let converge.
2. Let the FAST-LTS converge to some h element subset. Perform one step of the MOEA and use the result as the input to the FAST-LTS. We can then iterate between steps of these two algorithms till convergence.

On the large data sets, where even one step of the MOEA is too exhausting, we can use the eager MMEA which time complexity of one step is lower than in the case of MOEA.

In Chapter 3 we show the experimental results of various combinations of these algorithms.

2.6 BAB algorithm

In this section, we describe the algorithm from [5]. Very similar algorithm can also be found in [23].

Unlike previous algorithms, this one is exact. As we discussed in Section 1.4.1, exact algorithm calculate OLS fit on all of the $\binom{n}{h}$ h element subsets. For larger data sets this approach is computationally prohibitive.

This version of the exact algorithm is based on the branch and bound design paradigm; thus it tries to avoid exhaustive computation on all h element subsets. First, let us describe how the combination tree is built.

Let us denote subset of indexes $J_k = (j_1, j_2, \dots, j_k) \subset 1, 2, \dots, n$. We can then mark \mathbf{X}_{J_k} and \mathbf{Y}_{J_k} matrices created from \mathbf{X} and \mathbf{Y} so that all rows that

2. ALGORITHMS

are not indexed by J_k are removed. We can see that the number of subsets is given by $\binom{n}{k}$.

Let us now consider tree such that at level k is $\binom{n}{k}$ nodes representing all J_k subsets. The depth of the tree is h , so this tree has $\binom{n}{h}$ leaves representing all h index subsets. Example of such a tree we can see at Figure 2.3.

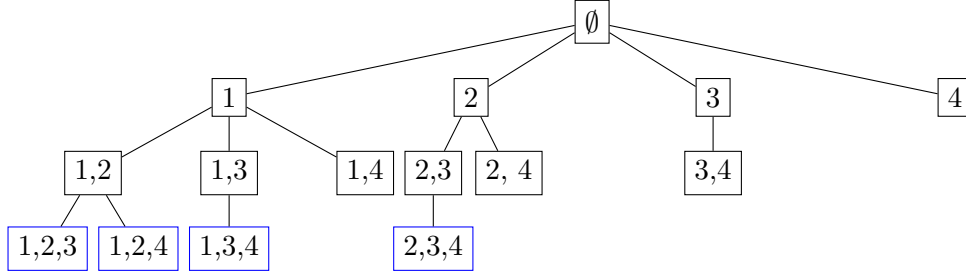


Figure 2.3: Tree consisting of all element subsets for $n = 4$ and $h = 3$. Leaves with blue border color represents h -element subsets.

In the case of the exhaustive approach, we can then traverse the tree starting in the root using the right to left (RTL) preorder traversal, and in each node, we can calculate the OLS fit on \mathbf{X}_{J_k} and \mathbf{Y}_{J_k} . If we are in the level $k = h$, we calculate the fit on the h element subset. This approach is exhaustive, so let us now describe how we can improve efficiency.

First, it is not necessary to calculate fit on k th level of the tree in a case $k < p$ because this would mean that matrix \mathbf{X}_{J_k} is not regular. The main improvement can be because we can skip parts of the tree (perform pruning) based on the following criteria.

If at any step (at any node) of the tree traversal we calculate the $\text{OF}_D^{\text{LTS}}(m_k)$ where $m_{k_i} = 1$ when $k_i \in J_k$ and $m_{k_i} = 0$ and we find out that this value is higher than minimal value found for some $\text{OF}_D^{\text{LTS}}(m_h)$ thus value of OLS fit on some h element subset, then we can discard all subsets which contain J_k . That means we can trim all descendants of the node which represent J_k subset.

We can describe the algorithm as follows:

1. Set $RSS^* = \inf. m^*$
2. For each node in RTL preorder traversal calculate $\text{OF}_D^{\text{LTS}}(m_k)$
3. If $\text{OF}_D^{\text{LTS}}(m_k) > RSS^*$, skip all siblings of this node in the traversal.
4. If $\text{OF}_D^{\text{LTS}}(m_k) < RSS^*$ and $k = h$, then set $RSS^* = \text{OF}_D^{\text{LTS}}(m_k)$ and m^* so that $m_{k_i} = 1$ when $k_i \in J_k$ and $m_{k_i} = 0$.
5. After the traversal return m^* and $\hat{\mathbf{w}}$

During the traversal of the tree, we only need to know the path back to the root. Thus the whole tree does not have to be in the memory.

Remark 33. If we calculated $\text{OF}_D^{\text{LTS}}(m_k)$ than we do not need to calculate $\text{OF}_D^{\text{LTS}}(m_{k+1})$ from scratch. In Section 2.4 we described the way how to update it when a row is added using the matrix inversion any we can apply it here. In Section 2.4.4 we described the way of doing this using the matrix decomposition which is another option.

Moreover, we described that if rows are not extracted and only inserted, then matrix \mathbf{Q} do not have to be present and only matrix \mathbf{R} can be updated. We can use this in our favor because we can keep all the decompositions on the path from the root in the memory. If we use matrix $\tilde{\mathbf{Z}}$ and $\tilde{\mathbf{R}}$ we can (see Remark 32) calculate all the quantities required for the update only using the matrix $\tilde{\mathbf{R}}$. This means that both versions of updating are possible and both can be done in $\mathcal{O}(p^2)$. Thus we recommend using the matrix decomposition version of the updating because it can provide higher numerical stability at the same speed.

2.6.1 Improvements

Because the tree is pruned based on the smallest value RSS^* , it would be convenient if we were able to obtain small value RSS^* early in the traversal. We can obtain this value by calculating an approximative LTS estimate. This can be done using any probabilistic algorithm described in the previous sections of this chapter. An ideal candidate may be some combined algorithm described in Section 2.5. When we find such small value, it is then also convenient to permute the observations based on its absolute residuals based on this LTS estimate in the decreasing order.

Another improvement can be made using the following sorting rule *sorting rule*: If we are at the level k (we assume that $k > p$) in the node J_k and this node has s siblings indexed as $s_1 \dots s_s$ then we can change the order of those siblings first by calculating partial increment (2.56) for each sibling and consequently ordered them in the descending order from left to right. This means we visit siblings with a lowest partial increment first. Trimming can then be done in the same manner. Moreover, if we calculate RSS in one of those siblings we and realize that $RSS > RSS^*$, then we can on the top of trimming all siblings of this sibling trim also all brothers left to this sibling. That means we can trim parent node J_k because all of his siblings have already been explored or can be trimmed. We can see such a procedure at Figure 2.4.

2.7 BSA algorithm

In this section, we'll introduce another exact algorithm. It has a little different approach than previous algorithms. First of all, build a theoretical basis for

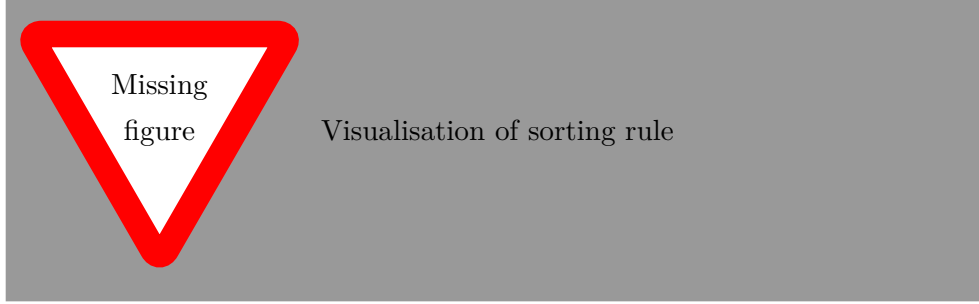


Figure 2.4: Step in the algorithms when due to the sorted siblings node j_k can be trimmed

this algorithm.

2.7.0.1 Domain of OF-LTS

We introduced two version of OF-LTS. First one is $\text{OF}^{(LTS,h,n)}(\mathbf{w})$, $\mathbf{w} \in \mathbb{R}^p$ and the second is discrete version $\text{OF}_D^{\text{LTS}}(\mathbf{m})$, $\mathbf{m} \in Q^{(n,h)}$. We also know that $\min_{\mathbf{w} \in \mathbb{R}^p} \text{OF}^{(LTS,h,n)}(\mathbf{w}) = \min_{\mathbf{m} \in Q^{(n,h)}} \text{OF}_D^{\text{LTS}}(\mathbf{m})$.

(see (1.27)) Let us now talk about the non-discrete version and introduce some new feature.

Definition 34. Let $Z \subset \mathbb{R}^p \times Q^{(n,h)}$ be an relation defined as

$$(\mathbf{w}, \mathbf{m}) \in Z \Leftrightarrow \sum_{i=1}^h r_{i:n}^2(\mathbf{w}) = \sum_{i=1}^n m_i r_i^2(\mathbf{w}) \quad (2.94)$$

Z is not a mapping. To show this we can take simple example so that $r_h^2 = r_{h+1}^2$. Then we have for the given \mathbf{w} two different vectors \mathbf{m} that are in relation with it.

For that reason, let us define $\mathcal{U} \subset \mathbb{R}^p$ as a maximal set where Z is a mapping. Next we define $\mathcal{H} = \mathbb{R}^p \setminus \mathcal{U}$ as a complement of \mathcal{U} .

Let's now describe some properties based on which \mathbf{w} is in \mathcal{U} or \mathcal{H} set.

Theorem 35. $\mathbf{w} \in \mathcal{U}$, $\mathbf{w} \in \mathbb{R}^p$ if and only if $r_{(h:n)}^2(\mathbf{w}) < r_{(h+1:n)}^2(\mathbf{w})$

Proof. As shown in example above. If $r_i^2(\mathbf{w}) = r_{(h:n)}^2(\mathbf{w}) = r_{(h+1:n)}^2(\mathbf{w}) = r_j^2(\mathbf{w})$, $i, j \in \{1, 2, \dots, n\}$ then $(\mathbf{w}, \mathbf{m}_1) \in Z$ and also $(\mathbf{w}, \mathbf{m}_2) \in Z$ where \mathbf{m}_1 has ones at indexes of h smallest residuals and \mathbf{m}_2 has ones at the same indexes except swap i th one with j th zero. \square

Corollary 36. Complement of vectors \mathbf{w} from Theorem 35 are vectors such that

$$r_{(h:n)}^2(\mathbf{w}) = r_{(h+1:n)}^2(\mathbf{w}) \quad (2.95)$$

thus

$$\mathcal{H} = \{\mathbf{w} \in \mathbb{R}^p | r_{(h:n)}^2(\mathbf{w}) = r_{(h+1:n)}^2(\mathbf{w})\}. \quad (2.96)$$

That means that for each $\mathbf{w} \in \mathcal{H}$ there are two different (\mathbf{x}_i, y_i) and (\mathbf{x}_j, y_j) so that

$$(y_i - \mathbf{x}_i \mathbf{w})^2 = r_i^2(\mathbf{w}) = r_{(h:n)}^2(\mathbf{w}) = r_{(h+1:n)}^2(\mathbf{w}) = r_j^2(\mathbf{w}) = (y_j - \mathbf{x}_j \mathbf{w})^2. \quad (2.97)$$

We can see that

$$(y_i - \mathbf{x}_i \mathbf{w})^2 = (y_j - \mathbf{x}_j \mathbf{w})^2 \iff y_i \pm y_j + (\mathbf{x}_i \pm \mathbf{x}_j) \mathbf{w} = 0 \quad (2.98)$$

Assumption 37. For $\mathbf{X} \in \mathbb{R}^{n \times p}$ let us assume that for all $i, j \in \{1, 2, \dots, n\}$ if $i \neq j$ then $\mathbf{x}_i \neq \pm \mathbf{x}_j$ and $\|\mathbf{x}_i\| \neq 0$.

If Assumption 37 is fulfilled then $y_i \pm y_j + (\mathbf{x}_i \pm \mathbf{x}_j) \mathbf{w} = 0$ represents a hyperplane.

It's easy to show that set \mathcal{U} is an open and Lebesgue measure of \mathcal{H} is 0 [24]. Moreover \mathcal{H} splits \mathbb{R}^p into finite number of m open disjoint subsets of \mathcal{U} .

Definition 38. Let's define set of $m \in \mathbb{N}$ sets $\mathcal{U}^{(set)} = \{U_i\}_{i=1}^m$ so that all U_i are open and connected sets, U_i, U_j are disjoint thus $U_i \cap U_j = \emptyset$, $\cup_{i=1}^m U_i = \mathcal{U}$ and $\cup_{i=1}^m \partial U_i = \mathcal{H}$.

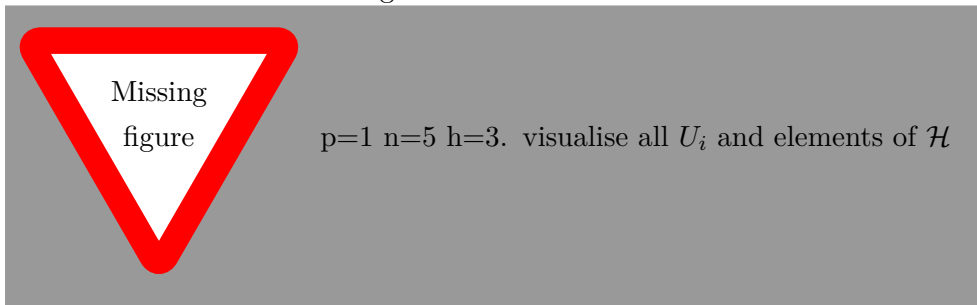
We define neighbor sets $U_i, U_j, i \neq j$ so that $U_i \cap U_j \neq \emptyset$.

We also define $M^{(min)}$ set of m vectors $\mathbf{m} \in Q^{(n,h)}$ so that

$$M^{(min)} = \{\mathbf{m}_1, \dots, \mathbf{m}_m | \mathbf{m}_i = Z(\mathbf{w}) \mathbf{w} \in U_i\}.$$

where $Z(\mathbf{w})$ represent unique vector \mathbf{m}_i . That means $Z(\mathbf{w}) = Z(\mathbf{w}')$ for all $\mathbf{w}, \mathbf{w}' \in U_i$.

For a better understanding definition above see



We say that matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ has h -full if for all $\mathbf{m} \in Q^{(n,h)}$ matrix $\mathbf{M}\mathbf{X}$ has rank p .

Theorem 39. If the matrix \mathbf{X} has h -full rank then for every local minima at point satisfying weak-necessary condition \mathbf{w}_0 of the OF-LTS function, then there exists a vector $\mathbf{m} \in Q^{(n,h)}$ such that $(\mathbf{w}_0, \mathbf{m}) \in Z$.

Proof can be found in [24, Theorem 7].

2.7.1 One-dimensional version of the algorithm

Based on the Theorem 39 we can see that

$$\min_{\mathbf{m} \in Q^{(n,h)}} \text{OF}^{(OLS, \mathbf{M}\mathbf{X}, \mathbf{M}\mathbf{Y})}(\mathbf{w}) = \min_{\mathbf{m} \in M^{(min)}} \text{OF}^{(OLS, \mathbf{m}^i \mathbf{X}, \mathbf{m}^i \mathbf{Y})}(\mathbf{w}) \quad (2.99)$$

where $\mathbf{M} = \text{diag}(\mathbf{m})$, $\mathbf{m}^i = \text{diag}(m)$. Set $M^{(min)}$ is very useful because we can iterate only through this set and not through whole $Q^{(n,h)}$.

Then minimizing objective function $\text{OF}_D^{\text{LTS}}(\mathbf{m})$ (1.31) can be reformulated as

$$\min_{\mathbf{m} \in Q^{(n,h)}} \text{OF}_D^{\text{LTS}}(\mathbf{m}) = \min_{\mathbf{m} \in M^{(min)}} \text{OF}_D^{\text{LTS}}(m) \quad (2.100)$$

That means if we can find all $m \in M^{(min)}$ then minimizing OF_D^{LTS} would be easy. So how we can find the elements of $M^{(min)}$ set?

For now, let us assume that we know elements of \mathcal{H} .

Because for each $m \in M^{(min)}$ exists at least one $\mathbf{w} \in \mathcal{H}$ such that $(\mathbf{w}, \mathbf{m}) \in Z$ then for given $\mathbf{w} \in \mathcal{H}$ we can find all \mathbf{m} by algorithm described by following pseudocode:

Algorithm 7: Find all \mathbf{m}

Input: $\mathbf{Q} \in \mathbb{R}^{n \times n}$, $\mathbf{R} \in \mathbb{R}^{n \times p}$, $\mathbf{x}_j \in \mathbb{R}^p$, k

Output: $\mathbf{Q}^{(+)} \in \mathbb{R}^{n+1 \times n+1}$, $\mathbf{R}^{(+)} \in \mathbb{R}^{n+1 \times p}$

1 **return** $\mathbf{Q}^{(+)}$, $\mathbf{R}^{(+)}$;

Finally we need to find all elements of \mathcal{H} . As we know for all elements $\mathbf{w} \in \mathcal{H}$ $r_{(h:n)}^2(\mathbf{w}) = r_{(h+1:n)}^2(\mathbf{w})$. So if we find all \mathbf{w} that solve this, then we found \mathcal{H} . The problem is that residuals are sorted. To overcome this issue by finding some candidates which would possible be equal $(h:n)$ th and $(h+1:n)$ th squared residuals. In [24] we are provided necessary condition for such a candidates.

Definition 40. If $\mathbf{w} \in \mathcal{H}$ then there exists $i, j, i \neq j$ so that $r_i^2(\mathbf{w}) = r_j^2(\mathbf{w})$. We denote set H containing \mathbf{w} satisfying this necessary condition, thus

$$H = \{\mathbf{w} \in \mathbb{R}^p | r_i^2(\mathbf{w}) = r_j^2(\mathbf{w}), i \neq j\}. \quad (2.101)$$

Note that $\binom{|H| \leq 2}{np+1}$. So in the case of the one dimensional case $\binom{|H| \leq 2}{n2=n(n-1)}$.

Finding all elements of H is then equal for for solving $\binom{|H| \leq 2}{np+1}$ equations. Moreover, because equations are quadratic, we can have up to two solutions for each equation.

So idea of the whole algorithm is to find all elements of H (by solving quadratic equations), consequently finding which of them are part of \mathcal{H} (by ordering squared residuals and checking if $r_{(h:n)}^2(\mathbf{w}) = r_{(h+1:n)}^2(\mathbf{w})$ and finally for each $\mathbf{w} \in \mathcal{H}$ find \mathbf{m} that are in relation. All of those \mathbf{m} vectors form $M^{(min)}$ set. Finally we can use $\mathbf{m} \in M^{(min)}$ for minimizing $\text{OF}_D^{\text{LTS}}(\mathbf{m})$ by the terms of (2.100). The whole algorithm, which authors call Border Scanning Algorithm (BSA) is described by the following pseudocode:

Algorithm 8: BSA**Input:** $Q \in \mathbb{R}^{n \times n}, R \in \mathbb{R}^{n \times p}, x_j \in \mathbb{R}^p, k$ ~~**Output:** $Q^{(+)} \in \mathbb{R}^{n+1 \times n+1}, R^{(+)} \in \mathbb{R}^{n+1 \times p}$~~ **1 return** $Q^{(+)}, R^{(+)}$;

write the pseudo-code

2.7.2 Multidimensional BSA

If we would like to scale the algorithm into the higher dimension, we have to face one major complication. The algorithm finds all elements of H and \mathcal{H} . This is possible in one dimension, but because \mathcal{H} forms a hyperplane, then for higher dimensions is infinite.

Authors then propose to find finite set \mathcal{H}_p such that

$$\forall \mathbf{m} \in M^{(min)} \exists \mathbf{w} \in \mathcal{H}_p : (\mathbf{w}, \mathbf{m}) \in Z \quad (2.102)$$

and also finite set H_p , so that $\mathcal{H}_p \subset H_p \subset H$. We can see that for $p = 1$ we need only one equation so the solution has dimension 0. If we would like to have some solution with dimension equal to 0 for $p > 1$ we need system of p equations. Moreover this system of equations have to be independent. So we can consider system of p equations of type $r_i^2(\mathbf{w}) = r_j^2(\mathbf{w})$

$$\begin{aligned} r_{i_1}^2(\mathbf{w}) &= r_{i_2}^2(\mathbf{w}) \\ r_{i_2}^2(\mathbf{w}) &= r_{i_3}^2(\mathbf{w}) \\ &\vdots \\ r_{i_p}^2(\mathbf{w}) &= r_{i_{p+1}}^2(\mathbf{w}) \end{aligned} \quad (2.103)$$

where $i_1, i_2 \dots i_{p+1}$ is $(p+1)$ -element subset of $1, 2, \dots n$.

System of those quadratic equations can be rewritten as Because this is a system of quadratic equations, we can have up to 2^p solutions.

Moreover total number of $(p+1)$ -element subsets out of n is $\binom{n}{p+1}$. This means that the total number of solutions can be up to $\binom{n}{p+1} 2^p$. We denote set of all these solutions H_p .

So if set H_p satisfies (2.102) then it can be used for finding \mathcal{H}_p subset, thus we would be able to use this set for the multidimensional version of the algorithm.

Let $\circ \in +, -$ be either operation $+$ or $-$. Then system of quadratic equations (2.103) is equivalent to 2^p linear systems

$$\begin{aligned} (\mathbf{x}_{i_1} \circ_1 \mathbf{x}_{i_2})^T \mathbf{w} &= y_{i_1} \circ_1 y_{i_2} \\ &\vdots \\ (\mathbf{x}_{i_p} \circ_p \mathbf{x}_{i_{p+1}})^T \mathbf{w} &= y_{i_p} \circ_p y_{i_{p+1}}. \end{aligned} \quad (2.104)$$

2. ALGORITHMS

Moreover (2.104) is equivalent to the set of equations

$$\begin{aligned} (\mathbf{x}_{i_1} \circ_1 \mathbf{x}_{i_2})^T \mathbf{w} &= y_{i_1} \circ_1 y_{i_2} \\ &\vdots \\ (\mathbf{x}_{i_1} \circ_p \mathbf{x}_{i_{p+1}})^T \mathbf{w} &= y_{i_1} \circ_p y_{i_{p+1}}. \end{aligned} \quad (2.105)$$

Elements of set \mathcal{H}_p then satisfies $r_{i_1^2}(\mathbf{w}) = r_{(h:n)}(\mathbf{w}) = r_{(h+1:n)}(\mathbf{w})$.

Set H_p was proven to be suitable for finding \mathcal{H}_p in [24] thus that set \mathcal{H}_p satisfies (2.102) if following assumptions are satisfied (we already assume that Assumption 37 holds). First is that for all \mathbf{w} $r_{(h:n)}(\mathbf{w}) > 0$, second, for all $(\circ_1, \circ_2, \dots, \circ_{n-1}) \in \times_{i=1}^n +, -$ the matrix system of $n-1$ equations

$$\begin{aligned} (\mathbf{x}_{i_1} \circ_1 \mathbf{x}_{i_2})^T \mathbf{w} &= y_{i_1} \circ_1 y_{i_2} \\ &\vdots \\ (\mathbf{x}_{i_1} \circ_{n-1} \mathbf{x}_{i_n})^T \mathbf{w} &= y_{i_1} \circ_p y_{i_n}. \end{aligned} \quad (2.106)$$

The second assumption prevents the set \mathcal{H}_p to be empty thus the case where all systems (2.105) are dependent.

This assumption can cause problem because when we assume model with the intercept, thus for all rows \mathbf{x}_i of matrix \mathbf{X} is $x_{i1} = 1$ then in the case $\circ_1, \circ_2, \dots, \circ_{n-1} = -$, first column of matrix \mathbf{X} only contains zeroes. We can modify this assumption so that instead of using $(\circ_1, \circ_2, \dots, \circ_{n-1}) \in \times_{i=1}^n +, -$ we use only $(\circ_1, \circ_2, \dots, \circ_{n-1}) \in \times_{i=1}^n +, - \setminus (-, \dots, -)$ [24].

Given that we now know that set H_p and \mathcal{H}_p are suitable for our algorithm and we also know how to calculate them, so let us describe the multidimensional version of the algorithm.

write the pseudo-code

Let us now describe the time complexity of this algorithm.

write the time complexity

Experiments

In this chapter we describe experiments and their results for our implementation of all algorithms described in the previous chapter. For testing performance of algorithms we have implemented data set generator which provide wide variety of configuration options of generating the data.

3.1 Data set generator

When we want to generate n observations without outliers that satisfies linear regression model we can do it as follows:

Algorithm 41 (Generate clean data).

1. Generate regression coefficients $\mathbf{w} = (w_1, \dots, w_p)$ and set σ^2 .
2. Generate explanatory variable \mathbf{x}_i .
3. Generate random noise $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$.
4. Simulate dependent variable $y_i = \mathbf{w}^T \mathbf{x}_i + \varepsilon_i$.
5. Repeat steps 2–4 n times.

As a result we obtain \mathbf{X} and \mathbf{y} . Regression coefficients can be set to arbitrary values, but if we do not want to have dummy explanatory variables, then they are supposed to be non-zero. All \mathbf{x}_i should be generated independently and it is very common that we generate all \mathbf{x}_i from normal distribution.

Another thing that needs to be considered is the intercept. In this work we assumed that our data already include intercept, so in that case w_1 is equal to intercept and all x_{i1} should be equal to 1. Note that same result can be obtained by generating data without intercept and with $\varepsilon_i \sim \mathcal{N}(\mu, \sigma^2)$. We can then extend matrix \mathbf{X} so that we add first column that contains only 1s. This approach is very common and software for estimating regression coefficients usually allows to set parameter which determines if intercept should be

3. EXPERIMENTS

used; if so, the column of 1s is added. For that reason we generate our data sets using this approach. That means we generate $\varepsilon_i \sim \mathcal{N}(\mu, \sigma^2)$ and column of 1s is included only in case we set parameter for using intercept when fitting the data set.

3.1.1 Generating outliers

As we have already described in Section 1.3.1 we can describe different types of the outliers: vertical outliers and two types leverage points — good leverage points and bad leverage points. Those types of outliers are visualized on Figure 3.1. We can see that good leverage points are not deemed as an outliers here, even if they are distant observations, because they follow the linear pattern.

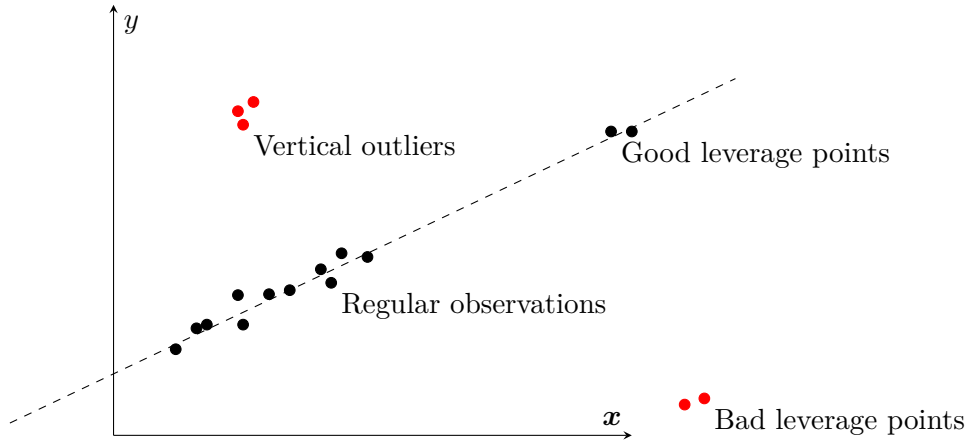


Figure 3.1: Different types of outliers.

Moreover, the data set can contain multiple observations that satisfies linear regression model but with different regression coefficients. That means such data set can contain data from multiple different models.

To generate the vertical outliers we only need to modify step 3 of Algorithm 41. We have multiple options:

- Generate ε_i from $\mathcal{N}(\mu, \sigma^2)$ but use different parameter μ and σ^2 than which were used for generating regular observations.
- Generate ε_i from some heavy tailed or asymmetrical distribution like Log-normal or exponential distribution.
- Combine both above so that we randomly choose distribution and randomly generate parameters for such distribution.

The last option is the most versatile so we use this option.

Because we generate \mathbf{x}_i from the normal distribution, we can generate leverage points just by changing parameter μ of this distribution. If we consequently generate ε_i from the same distribution with same parameters as for the regular observations, we obtain good leverage points. On the other hand if we generate ε_i as described above, we obtain bad leverage points.

If we want to generate outliers that correspond to the different model we can just set the regression coefficients \mathbf{w} differently. It is also possible to use different parameters of normal distribution for generating \mathbf{x}_i and parameters for generating ε_i . Theoretically we are able to introduce outliers even into this model, but when this model is “outlier” by itself relative to the original model, it is not needed. By this approach we would be able to generate the observations from arbitrary number of differ models, but for the sake of the simplicity we introduce only one different model.

3.2 Data sets

We have implemented ideas from previous section to the data set generator with following parameters:

- n and p for setting number of the generated observations and dimension of the explanatory variables
- *outlier ratio* for setting proportion of the outliers in the data set. This include vertical outliers, bad leverage points and also outliers from the second model.
- *leverage ratio* proportion of the explanatory variables that are generated as leverage points.
- $\mu_{\mathbf{x}}, \sigma_{\mathbf{x}}^2$ parameters of the normal distribution for generating non outlying \mathbf{x}_i
- $\mu_{\mathbf{x}_o}, \sigma_{\mathbf{x}_o}^2$ parameters of the normal distribution for generating outlying \mathbf{x}_i (leverage points)
- $\mu_{\varepsilon}, \sigma_{\varepsilon}^2$ parameters of the normal distribution for generating non outlying errors ε_i
- $\mu_{\varepsilon_o}, \sigma_{\varepsilon_o}^2$ parameters of the distribution for generating outlying errors ε_i .
- *distrib _{ε_o}* distribution from which outlying errors are generated — options are normal distribution, log-normal distribution and exponential distribution (when exponential distribution is chosen, then only $\sigma_{\varepsilon_o}^2$ parameter is used)
- *2m ratio* proportion from the outliers which are generated from the second model

3. EXPERIMENTS

- $\mu_{x_{M2}}, \sigma_{x_{M2}}^2$ and $\mu_{\varepsilon_{M2}}, \sigma_{\varepsilon_{M2}}^2$ parameters for normal distributions for generating \mathbf{x}_i and ε_i respectively from the second model

We used this generator to generate three data sets $D1$, $D2$ and $D3$ which differ by the types of the outliers they contain:

- $D1$ contains outliers which are not from the second model: vertical outliers, bad leverage points and good leverage points ($2m \text{ ratio} = 0$)
- $D2$ contains only outliers from the second model ($2m \text{ ratio} = 1$)
- $D3$ contain outliers of all described types. ($2m \text{ ratio} = 0.4$)

All three data sets are set to contain 20% leverage points (thus *leverage ratio* = 0.2) and non outlying \mathbf{x}_i are generated from $\mathcal{N}(0, 10)$ (thus $\mu_x = 0, \sigma_x^2 = 10$). Other parameters are independently randomly generated from uniform distribution so that:

- $\mu_{x_o} \sim \mathcal{U}(20, 60), \sigma_{x_o}^2 \sim \mathcal{U}(10, 20)$
- $\mu_{\varepsilon} \sim \mathcal{U}(0, 10), \sigma_{\varepsilon}^2 \sim \mathcal{U}(1, 5)$
- $\mu_{\varepsilon_o} \sim \mathcal{U}(-50, 50), \sigma_{\varepsilon_o}^2 \sim \mathcal{U}(50, 200)$
- $\mu_{x_{M2}} \sim \mathcal{U}(-30, 30), \sigma_{x_{M2}}^2 \sim \mathcal{U}(10, 20)$
- $\mu_{\varepsilon_{M2}} \sim \mathcal{U}(-10, 10), \sigma_{\varepsilon_{M2}}^2 \sim \mathcal{U}(1, 5)$
- $\text{distrib}_{\varepsilon_o}$ is uniformly randomly set to normal, log-normal or exponential distribution

Finally parameters n, p and *outlier ratio* are set depending on the experiment.

3.3 Implementation of the algorithms

We have implemented all described algorithms, moreover because algorithms for computing the feasible solution could be implemented both by calculating inversion and by calculating QR decomposition, we have implemented both version of those algorithms. Here is the list of all the implemented algorithm with their acronyms we use for labeling them:

FAST-LTS from Section 2.2 with all described improvements

FSA-INV from Section 2.3

FSA-QR FSA using theory from Section 2.4

MOEA-INV from Section 2.4 it is the improved version of OEA

MOEA-QR MOEA using theory from Section 2.4.4

MMEA-INV from Section 2.4.3

MMEA-QR MMEA using theory from Section 2.4.4

BAB from Section 2.6

BSA is the implementation of improved BAB (BSABAB) from Section 2.7

FAST-LTS-MMEA-INV combination of algorithms from Section 2.5

FAST-LTS-MOEA-INV another combination of algorithm from Section 2.5

MOEA-QR-BAB BAB with sorting speedup as described in Section 2.6

MOEA-QR-BSA BSABAB using our idea from Section 2.7

P-BSA probabilistic BSA using our idea from Section 2.7

These algorithms were first implemented in Python using the NumPy package [25]. The performance was low, so we have implemented all of the algorithms in C++ using the Eigen library [26] for matrix manipulation. On the other hand it is very popular today to use Python for data processing and manipulation; for that reason we have used pybind11 library [27], that exposes C++ types for Python and vice versa and written Python wrappers around the C++ implementation.

Moreover pybind11 allows to bind Eigen types directly to the NumPy types (because both libraries are LAPACK compatible), so that it is possible to share pointers to the matrices between Eigen and NumPy so that the data does not have to be copied when transferring between Python and C++.

That means that the data generator, all tests and experiments are implemented in Python; the C++ code is called only within the Python wrappers. It is also appropriate to mention that interface of the algorithms was created with regards to the popular scikit-learn package [28]; the interface is almost identical, so all of the classes implementing the algorithms can be used in the same manner as classes from scikit-learn linear regression module.

3.4 Results

In two previous sections we have described our experimental setup. Here we report results of our experiments.

In the first experiment we run multiple simulations where we compare speed and accuracy of the algorithm finding subsets satisfying strong necessary condition e.g. h -element feasible subsets. For each combination parameters n , p and *outlier ratio* we generate we generate data sets $D1$, $D2$ and $D3$ 100 times (each time new data sets are generated) and run all algorithms on those

3. EXPERIMENTS

data sets. Value of h is conservatively chosen so that $h = \lceil (n/2) + \lceil (p+1)/2 \rceil$. For all runs we use the intercept, so value p represents dimension of \mathbf{x} including intercept. All algorithm are set to run at most for 50 steps and each of them is starting only from 1 randomly selected h -element subset. The results are given in table 3.1 were average CPU time from 100 runs for each algorithm. We also measure cosine similarity of a given solution compared to the OLS solution on the subset which does not contain outliers. For $n > 500$ cells for algorithm FSA-I and FSA-QR are empty. That is because times of the run of these algorithms were too slow and it would take weeks to finish all simulations.

n	p	algorithm	FAST-LTS		FSA-I		FSA-QR		MMEA-I		MMEA-QR		MOEA-I		MOEA-QR	
		out	time	cos	time	cos	time	cos	time	cos	time	cos	time	cos	time	cos
20	2	0.10	0.004614	0.988359	0.000202	0.982532	0.006610	0.990883	0.000071	0.982120	0.000147	0.989808	0.000364	0.817895	0.000187	0.987816
		0.30	0.004291	0.985132	0.000213	0.986326	0.006309	0.963036	0.000068	0.814509	0.000152	0.993869	0.000370	0.656299	0.000202	0.972483
		0.45	0.004183	0.964252	0.000209	0.847824	0.004229	0.946639	0.000067	0.634588	0.000145	0.916261	0.000276	0.609704	0.000192	0.862572
100	3	0.10	0.010263	0.997496	0.019629	0.997817	0.178909	0.998082	0.000839	0.981119	0.002535	0.997572	0.002137	0.911982	0.006485	0.997943
		0.30	0.009968	0.998368	0.020972	0.998185	0.185812	0.998594	0.000854	0.924159	0.002704	0.996901	0.001721	0.877451	0.007001	0.998725
		0.45	0.009569	0.998725	0.020623	0.998749	0.150607	0.998839	0.000860	0.936188	0.002638	0.998588	0.001674	0.890495	0.006928	0.997922
	5	0.10	0.012194	0.996578	0.021666	0.996789	0.265649	0.997602	0.000884	0.990368	0.002995	0.997628	0.002414	0.885812	0.007262	0.997378
		0.30	0.012523	0.996829	0.024384	0.997198	0.286329	0.997886	0.000912	0.923997	0.003313	0.997691	0.002860	0.852555	0.007919	0.997339
		0.45	0.012282	0.997724	0.024928	0.996658	0.220120	0.994148	0.000933	0.891392	0.003232	0.997575	0.001798	0.865642	0.007690	0.993144
500	2	0.10	0.130725	0.999446	1.708373	0.999614	6.854093	0.999600	0.076088	0.999035	0.172415	0.999533	0.151864	0.998190	0.496092	0.999639
		0.30	0.134994	0.999591	1.661876	0.999801	6.661488	0.999791	0.075539	0.997468	0.173445	0.999742	0.157311	0.997426	0.488501	0.999702
		0.45	0.119165	0.999817	1.786019	0.999110	7.326588	0.999472	0.074628	0.998243	0.165817	0.999305	0.141019	0.998003	0.503233	0.999351
	5	0.10	0.100103	0.999442	2.013579	0.999605	8.534774	0.999680	0.076842	0.996679	0.161578	0.999745	0.111078	0.996261	0.480725	0.999641
		0.30	0.094455	0.999282	2.062271	0.999604	8.700572	0.999516	0.071619	0.994945	0.162314	0.999543	0.104943	0.994869	0.485604	0.999472
		0.45	0.091336	0.999652	2.077072	0.999163	8.865941	0.999343	0.067733	0.993651	0.152552	0.999013	0.103857	0.994371	0.484279	0.999033

Table 3.1: Table to test captions and labels

Conclusion

To date, multiple exact and probabilistic algorithms for calculating LTS estimate has been proposed. Those algorithms have been proposed across the last few decades. Most recently proposed algorithms are just a few years old. That means that research in the filed of LTS algorithms is still prevalent.

Although to date exact finite algorithm has polynomial time, we showed that currently used probabilistic algorithm provide sufficiently fast solutions which, even though are not exact, are good enough. Though it is proven that the exact solution cannot be obtained faster than in polynomial time, we showed that currently used algorithms could be combined to obtain better results.

Algorithms for calculating the LTS estimate is an exciting topic for further research. One of the ways is the research of combining exact algorithms with probabilistic ones. As our experimental results prompt, we would possibly propose probabilistic algorithms which provide even better solutions at the same time.

Another way would be to use algorithms for computing LTS estimate on different robust statistic methods because some of them are very similar to LTS problem. Just put there are still a lot of future works on least trimmed squares, and we are planning to research this in the future.

Bibliography

- [1] McCullagh, P. *Generalized linear models*. Routledge, 2018.
- [2] Rousseeuw, P.; C. van Zomeren, B. Unmasking Multivariate Outliers and Leverage Points. *Journal of The American Statistical Association - J AMER STATIST ASSN*, volume 85, 06 1990: pp. 633–639, doi:10.1080/01621459.1990.10474920.
- [3] Hampel, F. R.; Ronchetti, E. M.; et al. *Robust statistics*. Wiley Online Library, 1986.
- [4] Massart, D. L.; Kaufman, L.; et al. Least median of squares: a robust method for outlier and model error detection in regression and calibration. *Analytica Chimica Acta*, volume 187, 1986: pp. 171–179.
- [5] Agulló, J. New algorithms for computing the least trimmed squares regression estimator. *Computational Statistics & Data Analysis*, volume 36, no. 4, 2001: pp. 425–439.
- [6] Bernholt, T. Robust estimators are hard to compute. Technical report, Technical Report/Universität Dortmund, SFB 475 Komplexitätsreduktion, 2006.
- [7] Klouda, K. *Studium senzitivity odhadu metodou nejmenších usekaných čtvrců*. Bachelor’s thesis, Technical University of Prague, Faculty of Nuclear Sciences and Physical Engineering, Department of Mathematics, 2006.
- [8] Hawkins, D. M. The feasible solution algorithm for least trimmed squares regression. *Computational statistics & data analysis*, volume 17, no. 2, 1994: pp. 185–196.
- [9] Bai, E.-W. A random least-trimmed-squares identification algorithm. *Automatica*, volume 39, no. 9, 2003: pp. 1651–1659.

- [10] Rousseeuw, P. J.; Leroy, A. M. *Robust regression and outlier detection*. John Wiley & Sons, 1987.
- [11] Hawkins, D. M.; Olive, D. J. Improved feasible solution algorithms for high breakdown estimation. *Computational statistics & data analysis*, volume 30, no. 1, 1999: pp. 1–11.
- [12] Strassen, V. Gaussian elimination is not optimal. *Numerische mathematik*, volume 13, no. 4, 1969: pp. 354–356.
- [13] Coppersmith, D.; Winograd, S. Matrix multiplication via arithmetic progressions. *Journal of symbolic computation*, volume 9, no. 3, 1990: pp. 251–280.
- [14] Ballard, G.; Benson, A. R.; et al. Improving the numerical stability of fast matrix multiplication. *arXiv preprint arXiv:1507.00687*, 2015.
- [15] Anderson, E.; Bai, Z.; et al. *LAPACK Users' Guide*. Philadelphia, PA: Society for Industrial and Applied Mathematics, third edition, 1999, ISBN 0-89871-447-8 (paperback).
- [16] Krishnamoorthy, A.; Menon, D. Matrix inversion using Cholesky decomposition. In *2013 Signal Processing: Algorithms, Architectures, Arrangements, and Applications (SPA)*, IEEE, 2013, pp. 70–72.
- [17] Businger, P.; Golub, G. H. Linear least squares solutions by Householder transformations. *Numerische Mathematik*, volume 7, no. 3, 1965: pp. 269–276.
- [18] Hammarling, S.; Lucas, C. Updating the QR factorization and the least squares problem. 2008.
- [19] Rousseeuw, P. J.; Driessen, K. V. An Algorithm for Positive-Breakdown Regression Based on Concentration Steps. In *Data Analysis: Scientific Modeling and Practical Application*, edited by M. S. W. Gaul, O. Opitz, Springer-Verlag Berlin Heidelberg, 2000, pp. 335–346.
- [20] Hoare, C. A. Algorithm 65: find. *Communications of the ACM*, volume 4, no. 7, 1961: pp. 321–322.
- [21] Atkinson, A. C.; Weisberg, S. Simulated annealing for the detection of multiple outliers using least squares and least median of squares fitting. *Institute for Mathematics and Its Applications*, volume 33, 1991: p. 7.
- [22] Bartlett, M. S. An inverse matrix adjustment arising in discriminant analysis. *The Annals of Mathematical Statistics*, volume 22, no. 1, 1951: pp. 107–111.

- [23] Hofmann, M.; Kontoghiorghes, E. J. Matrix strategies for computing the least trimmed squares estimation of the general linear and sur models. *Computational Statistics & Data Analysis*, volume 54, no. 12, 2010: pp. 3392–3403.
- [24] Klouda, K. An exact polynomial time algorithm for computing the least trimmed squares estimate. *Computational Statistics & Data Analysis*, volume 84, 2015: pp. 27–40.
- [25] Oliphant, T. NumPy: A guide to NumPy. USA: Trelgol Publishing, 2006–, [Online; accessed `today`]. Available from: <http://www.numpy.org/>
- [26] Guennebaud, G.; Jacob, B.; et al. Eigen v3. <http://eigen.tuxfamily.org>, 2010.
- [27] Jakob, W.; Rhinelander, J.; et al. pybind11 – Seamless operability between C++11 and Python. 2017, <https://github.com/pybind/pybind11>.
- [28] Pedregosa, F.; Varoquaux, G.; et al. Scikit-learn: Machine Learning in Python . *Journal of Machine Learning Research*, volume 12, 2011: pp. 2825–2830.

Data sets

GUI Graphical user interface

XML Extensible markup language

Contents of enclosed CD

	readme.txt	the file with CD contents description
	exe	the directory with executables
	src	the directory of source codes
	wbdcm	implementation sources
	thesis	the directory of \LaTeX source codes of the thesis
	text	the thesis text directory
	thesis.pdf	the thesis text in PDF format
	thesis.ps	the thesis text in PS format