

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.

April 26, 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 April 26, 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 OLS estimation	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	10
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	23
2.2.2 Choosing initial \mathbf{m}_1 subset	26
2.2.3 Speed-up of the algorithm	27
2.2.4 Putting all together	28
2.3 Feasible solution	30
2.4 OEA, MOEA, MMEA	34
2.4.1 Multiplicative formula	34

2.4.2	The OEA and its properties	38
2.4.3	Minimum-maximum exchange algorithm	40
2.4.4	Different method of computation	41
2.5	Combined algorithm	46
2.6	BAB algorithm	46
2.6.1	Improvements	47
2.7	BSA algorithm	47
3	Experiments	49
3.1	Data	49
3.2	Results	49
3.3	Outlier detection	49
	Conclusion	51
	Bibliography	53
	A Data sets	55
	B Contents of enclosed CD	57

List of Figures

2.1	todo caption	24
2.2	todo caption	25
2.3	Tree consiting of all index subsets for $n = 4$ and $h = 3$	47

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 it is usual to consider vectors 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

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

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

In practice, we usually deal with multiple dependent variables so we define multiple linear regression model.

Definition 2. The *multiple linear regression model* is

$$y_i = \mathbf{x}_i^T \mathbf{w} + \varepsilon_i, i = 1, \dots, n. \quad (1.2)$$

It is common to write the whole model in a matrix form

$$\mathbf{y} = \mathbf{X} \mathbf{w} + \boldsymbol{\varepsilon} \quad (1.3)$$

where

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

That means that we can think of rows of matrix \mathbf{X} as columns vectors x_i written into the row. This means that even when we think of all vectors as columns given the matrix \mathbf{X} we consider x_i as row vectors.

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

In this work we talk about the multiple linear regression model, but for simplicity, we call it merely the linear regression model.

Note 3. It is usual refer to given \mathbf{X} and y as to *data set* and to y_i with corresponding x_i as to *data sample* or *observation*. Sometimes it's also useful to refer to multiple observations as to subset of observations.

1.1.1 Prediction with the linear regression model

The linear regression model contains the vector of actual 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

$$\hat{y} = \hat{\mathbf{w}}^T \mathbf{x}. \quad (1.4)$$

True value of y is given by

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

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

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

$$\mathbb{E}(y) = \mathbb{E}(\mathbf{x}^T \mathbf{w}) + \mathbb{E}(\varepsilon) = \mathbb{E}(\mathbf{x}^T \mathbf{w}) \quad (1.6)$$

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 4. Let us consider that y represents price of the room and $x \in \mathbb{N}$ 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 $y = \mathbf{w}^T \mathbf{x} + \varepsilon$ 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 $x_0 = 1$ then the corresponding coefficient w_0 of \mathbf{w} is called *intercept*. We refer this model as a *model with an intercept*. Intercept then corresponds to expected value of y when all regressors are zero and prevent the problem from example 4. Given that y is a random variable due to the ε , intercept actually 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 unless otherwise specified. This means that we consider $\mathbf{x} = (x_1, x_2 \dots x_p)$ where constant $x_1 = 1$ represents intercept.

Note 5. Sometimes in model with 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.

Some other assumptions about ε are in practice invalid, and we describe this in section 1.3.

1.2 Ordinary least squares

We want to estimate \mathbf{w} so that the error of the model will be the least possible. Measurement of this error done by a *loss function* $L : \mathbb{R}^2 \rightarrow \mathbb{R}$, which in case of ordinary least squares (OLS) is quadratic loss function $L(y, \hat{y}) := (y - \hat{y})^2$. So the basic idea is to find $\hat{\mathbf{w}}$ so that it minimizes the sum of squared *residuals*

$$r_i(\mathbf{w}) = y_i - \hat{y}_i = y_i - \mathbf{x}_i^T \mathbf{w}, \quad i = 1, 2, \dots, n. \quad (1.7)$$

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

$$RSS(\mathbf{w}) = \sum_{i=1}^n r_i^2(\mathbf{w}) = (y_i - \mathbf{w}^T \mathbf{x}_i)^2. \quad (1.8)$$

Definition 6. The RSS as the function of \mathbf{w} is an *objective function* for OLS denoted as

$$\text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 = (\mathbf{Y} - \mathbf{X}\mathbf{w})^T (\mathbf{Y} - \mathbf{X}\mathbf{w}) \quad (1.9)$$

The point of the minimum of this function

$$\hat{\mathbf{w}}^{(OLS, n)} = \arg \min_{\mathbf{w} \in \mathbb{R}^p} \text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}(\mathbf{w}) \quad (1.10)$$

is a definition of ordinary least squares estimate of true regression coefficients \mathbf{w}^* .

1. LEAST TRIMMED SQUARES

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

$$\frac{\partial \text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}}{\partial w_j} = \sum_{i=1}^n 2(y_i - \mathbf{w}^T \mathbf{x}_i)(-x_{ij}), \quad j \in \{1, 2, \dots, p\}. \quad (1.11)$$

By this we obtain the gradient

$$\nabla \text{OF}^{(OLS, \mathbf{X}, \mathbf{y})} = - \sum_{i=1}^n 2(y_i - \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i. \quad (1.12)$$

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

$$- \sum_{i=1}^n 2(y_i - \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i = 0 \quad (1.13)$$

that can be rewritten in matrix form as

$$\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \mathbf{w} = 0. \quad (1.14)$$

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

$$\frac{\partial^2 \text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}}{\partial w_h \partial w_j} = \sum_{i=1}^n 2(-x_{ih})(-x_{ij}), \quad h \in \{1, 2, \dots, p\}. \quad (1.15)$$

We get

$$\mathbf{H}_{\text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}} = 2\mathbf{X}^T \mathbf{X}. \quad (1.16)$$

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

$$\mathbf{s}^T (2\mathbf{X}^T \mathbf{X}) \mathbf{s} = 2(\mathbf{X} \mathbf{s})^T (\mathbf{X} \mathbf{s}) = 2 \|\mathbf{X} \mathbf{s}\|^2 \quad (1.17)$$

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. Give that the solution of (1.14) gives us the global minimum.

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

$$\hat{\mathbf{w}}^{(OLS, n)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad (1.18)$$

where (OLS, n) denotes that $\hat{\mathbf{w}}$ is estimate of true regression coefficients \mathbf{w}^* 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 so the $\text{OF}^{(OLS, \mathbf{X}, \mathbf{y})}$ is strictly convex thus $\hat{\mathbf{w}}^{(OLS, n)}$ is then strict global minimum.

1.2.1 Properties of OLS estimation

Gauss-Markov theorem tells us that under specific conditions (assumptions about the regression model) is OLS estimation unbiased and efficient. In other words, Gauss-Markov theorem states that OLS is the best linear unbiased estimator (BLUE). Efficient means that any other linear unbiased estimator has the same or higher variance than OLS. Those conditions are:

- Expected value of errors $\mathbb{E}(\varepsilon_i) = 0$, $i = 1, 2, \dots, n$.
- Errors are independently distributed and uncorrelated thus $cov(\varepsilon_i, \varepsilon_j) = 0$, $i, j = 1, 2, \dots, n$, $i \neq j$
- Regressors \mathbf{x}_i , $i = 1, 2, \dots, n$ and corresponding errors ε_i , $i = 1, 2, \dots, n$ are uncorrelated.
- All errors have same finite variance. This is known as *homoscedasticity*.

Note 7. As described in example 4 the first assumption is usually fulfilled if we use the model with intercept. Rest of the conditions, on the other hand, are rarely met.

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 statistic methods expect some assumptions to work properly and fail if those assumptions are not met. A robust statistic is statistics that produce acceptable results even when data are from some unconventional distributions or if data contains errors which are not normally distributed. We should be highly motivated to use such methods because in practice we are often forced to work with such data on which classical statistics methods provide poor results.

Such assumptions about 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 ordinary least squares to give a good estimate of $\hat{\mathbf{w}}$. Unfortunately, in real conditions, these assumptions are often false so that ordinary least squares do not produce acceptable results. One of the most common reasons for false assumptions are abnormal observations called outliers.

Outliers are very common, and they are for instance erroneous measurements such as transmission errors or noise. Another common reason is that nowadays we are mostly given data which were automatically processed by computers. Sometimes we are also presented 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 most of the data we are dealing with are too big to check, and even more often we do not know how such should look.

Moreover, in higher dimensions, it is complicated to find outliers. Some methods try to find outliers, but such methods are only partially efficient. Let us note that robust models 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 have some terminology to describe certain types of outliers. We use terminology from [1]. Let us have observation (y_i, \mathbf{x}_i) . If observation is not outlying in any direction we call it *regular observation*. If it is outlying in \mathbf{x}_i direction 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 linear pattern we call it *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 y_i direction, we call it a *vertical outliers*.

1.3.2 Measuring robustness

There are a couple of tools to measure the robustness of statistics. The most popular one is called *breakdown point*. Then there are *empirical influence function* and *influence function and sensitivity curve*. For the sake of simplicity, we describe only breakdown point right now.

Definition 8. Let T be a statistics, $\mathbf{x} = (x_1, x_2, \dots, x_n)$ be an n -dimensional random sample and $T_n(\mathbf{x})$ value of this statistics with parameter \mathbf{x} . The breakdown point of T at sample \mathbf{x} can be defined using sample $\bar{\mathbf{x}}^{(k)}$ where we exchange k points from original sample \mathbf{x} with random values x_i . We get $T_n(\bar{\mathbf{x}}^{(k)})$. Then the *breakdown point* is

$$\text{bdpoint}(T, \mathbf{x}_n) = \frac{1}{n} \min S_{T, \mathbf{x}_n}, \quad (1.19)$$

where

$$S_{T, \mathbf{x}_n, D} = \left\{ k \in \{1, 2, \dots, n\} : \sup_{\mathbf{x}_{new}^{(k)}} \|T_n(\mathbf{x}) - T_n(\mathbf{x}_{new}^{(k)})\| = \infty \right\}. \quad (1.20)$$

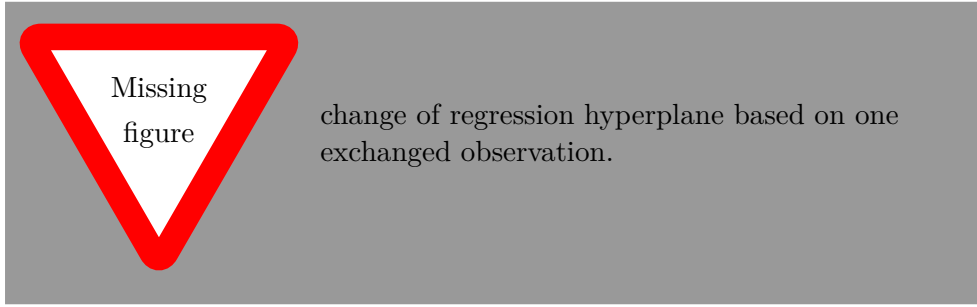
This definition is very general but let us specify it to our linear regression problem. It says that breakdown point is the function of the minimal number of observations needed to be changed for some others so then the estimator gives incorrect results. More robust estimators have higher breakdown point.

It is intuitive that reasonable breakdown point should not be higher than 0.5 [2] because if more than 50% of the data is exchanged, then the model of exchanged data overrides the model of the original data.

In the context of OLS estimator one outlier is enough to increase the value of T to any desired value [3] thus

$$\text{bdpoint}(OLS, \mathbf{x}_n) = \frac{1}{n}. \quad (1.21)$$

Figure [] gives us an example of how one outlier may change the hyperplane given by the OLS estimation of regression coefficients.



For an increasing number of data samples n this tends to zero. We can see that ordinary least squares estimator is not resistant to outliers at all. Due to this fact, multiple estimators similar to OLS have been proposed.

1.4 Least trimmed squares

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

Definition 9. 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 is then denoted as

$$\text{OF}^{(LTS,h,n)}(\mathbf{w}) = \sum_{i=1}^h r_{i:n}^2(\mathbf{w}) \quad (1.22)$$

where $r_{i:n}^2$ denotes i th smallest squared residuum out of n thus $r_{1:n}^2 \leq r_{2:n}^2 \leq \dots \leq r_{n:n}^2$. Even though that objective function of LTS seems similar to the OLS objective function, finding minimum is far more complex because the order of the least squared residuals is dependent on \mathbf{w} . Moreover $r_{i:n}^2(\mathbf{w})$ residuum is not uniquely determined if more squared residuals have same value as $r_{i:n}^2(\mathbf{w})$. This makes from finding LTS estimate non-convex optimization problem and finding global minimum is NP-hard [4].

1.4.1 Discrete objective function

The LTS objective function from definition 9 is not differentiable and not convex, so we are unable to use the same approach as with OLS objective function. Let us transform this objective function to a discrete version which can appropriately be used by algorithms to minimize it.

Let us assume for now that we know the $\hat{\mathbf{w}}^{(LTS,h,n)}$ vector of estimated regression coefficients. (LTS, h, n) denotes that coefficients are estimated using LTS on the h subset of n data samples. With this in mind, let π be the permutation of $\hat{n} = \{1, 2, \dots, n\}$ such that

$$r_{i:n} = r_{\pi(j)}, \quad j \in \hat{n}. \quad (1.23)$$

Put

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

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

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

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 $\mathbf{M}_{LTS} = \text{diag}(\mathbf{m}_{LTS}^T)$:

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

In other words, finding the minimum of the LTS objective function can be done by finding the OLS estimate (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 h subsets of \mathbf{X} and \mathbf{Y} , then [5]

$$\begin{aligned}
\min_{\mathbf{w} \in \mathbb{R}^p} \text{OF}^{(LTS, h, n)}(\mathbf{w}) &= \min_{\mathbf{w} \in \mathbb{R}^p} \sum_{i=1}^h r_{i:n}^2(\mathbf{w}) \\
&= \min_{\mathbf{w} \in \mathbb{R}^p, \mathbf{m} \in Q^{(n, h)}} \sum_{i=1}^n m_i r_i^2(\mathbf{w}) \\
&= \min_{\mathbf{m} \in Q^{(n, h)}} \left(\min_{\mathbf{w} \in \mathbb{R}^p} \text{OF}^{(OLS, \mathbf{M}\mathbf{X}, \mathbf{M}\mathbf{y})}(\mathbf{w}) \right) \\
&= \min_{\mathbf{m} \in Q^{(n, h)}} \left(\min_{\mathbf{w} \in \mathbb{R}^p} \|\mathbf{M}\mathbf{y} - \mathbf{M}\mathbf{X}\mathbf{w}\|^2 \right).
\end{aligned}$$

Substituting \mathbf{w} with (1.25) we get

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

We can easily see, that have the objective function with argument $\mathbf{m} \in Q^{(n, h)}$. This objective function we mark as

$$\text{OF}_D(\mathbf{m}) = \|\mathbf{M}\mathbf{y} - \mathbf{M}\mathbf{X}(\mathbf{X}^T \mathbf{M}\mathbf{X})^{-1} \mathbf{X}^T \mathbf{M}\mathbf{y}\|^2. \quad (1.27)$$

We can also clearly see that the minimizing this OF could be done straightforwardly by iterating over the $Q^{(n, h)}$ set. That means we transformed the objective function to the discrete space of $Q^{(n, h)}$.

Unfortunately, this set is vast, namely $\binom{n}{h}$, so this approach is infeasible on more massive 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 $h = \lfloor (n/2) + \lfloor (p+1)/2 \rfloor \rfloor$ [3] (where $\lfloor \cdot \rfloor$ denotes largest integer function). That 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 cover the theoretical background required to implement algorithms that are in this chapter. We introduced the discrete version of the LTS objective function which minimum is equivalent to the continuous one. Finding the minimum of this function requires to find the correct h subset and then calculate the estimate of regression coefficients \mathbf{w} . To achieve this, we need to examine all h subsets. The exhaustive approach fails due to the exponential size of $Q^{(n,h)}$. So what are the possibilities?

2.0.1 First attempts

Initial algorithms were based on iterative removal data samples whose residuum had the highest value based on OLS fit on the whole dataset. Such attempts are known to be flawed [6] because the initial OLS fit can be already profoundly affected by outliers and we may remove data samples which represent the actual model.

Other algorithms were based purely on a random approach. On such algorithm is Random solution algorithm [7] which randomly selects l of h subsets and subsequently compute OLS fit on each of them and chooses fit with a minimum value of the objective function. Such approach is straightforward, but the probability of selecting at least one subset from l subsets which don't contain outliers thus has a chance of producing good result tends to zero for an increasing number of data samples n as we'll describe in detail in section 2.2.2.

Another very similar algorithm called Resampling algorithm introduced in [8]. 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 l subsets of size $p+1$ thus it is independent of n gives the nonzero probability of selecting at least one subset such that it does not include 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 subsets is computationally exhaustive and relying on randomly generating “good” h -element subsets does not lead to reliable results. So what are our options? In [9] two criteria called *weak necessary conditions* and *strong necessary conditions* are introduced. They state necessary properties which some h subset must satisfy to be 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 subset of *non-trimmed* observations but also the complementary subset of not used observations. We’ll refer to this complementary subset as to *trimmed* subset.

Definition 10. 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 an 1 with one 0, we get $\text{OF}_D(\mathbf{m}) \leq \text{OF}_D(\mathbf{m}_{\text{swap}})$. Thus the value of 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.

The second necessary condition is named *weak necessary condition*

Definition 11. A h -element subset corresponding to $\mathbf{m} \in Q^{(n,h)}$ satisfies the *strong necessary condition* if $r_i^2(\hat{\mathbf{w}}^{(OLS, MX, My)})$ for all trimmed observation is greater or equal to the greatest non-trimmed squared residuum $r_j^2(\hat{\mathbf{w}}^{(OLS, MX, My)})$.

Again, based on this criteria an algorithm can be created. The corollary of this criteria together with proof can be found in section 2.2.

The interesting consequence which we’ll use later gives us the following lemma.

Lemma 12. The strong necessary condition is not satisfied unless the weak necessary condition is satisfied. Thus if a strong condition is satisfied then weak is also.

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

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

Thus

$$\text{OF}_D(\mathbf{m}) > \text{OF}_D(\mathbf{m}) + r_j^2 - r_i^2 \quad (2.1)$$

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

$$\text{OF}_D(\mathbf{m}) + r_j^2 - r_i^2 \geq \text{OF}_D(\mathbf{m}_{\text{swap}}). \quad (2.2)$$

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

Now we have covered all the necessary theoretical background, and it is time to introduce currently known algorithms for computing the LTS estimate.

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 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 matrix multiplication.

There are multiple algorithms for matrix multiplication, for example:

- Naive algorithm; its time complexity is $\mathcal{O}(n^3)$.
- Strassen algorithm; its time complexity is $\mathcal{O}(n^{2.8074})$ [10].
- Coppersmith-Winograd; its time complexity is $\mathcal{O}(n^{2.375477})$ [11]

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 case of the naive algorithm. For this reason, they are not used even when they might improve performance [12].

Moreover current implementations of the linear algebra namely LAPACK and more importantly BLAS level 3 routines implementing matrix-matrix multiplication which are widely used in such linear algebra software implement naive matrix-matrix multiplication [13].

For that reason we assume in this work that 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

Solving OLS using 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 inversion of \mathbf{A} .
3. Multiply $\mathbf{A}^{-1} \mathbf{B}$.

2. ALGORITHMS

Observation 13. 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 matrix inversion is $\mathcal{O}(p^2n)$.

Proof. First, we compute $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ 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 $\hat{\mathbf{w}}^{(OLS,n)} = \mathbf{C}\mathbf{B}$ which is $\mathcal{O}(p^2)$. Together 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 $n \geq p$ and usually even $n \gg p$ we get $\mathcal{O}(p^2n + p^3) \sim \mathcal{O}(p^2n)$. \square

Computing OLS estimate using matrix inversion is possible, however in most cases not used in practice because of the low numerical stability computing matrix inversion.

2.1.2 Computation using the Cholesky decomposition

Definition 14. Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be an symmetric positive definite matrix. Then it is possible to find lower triangular matrix \mathbf{L} so that

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

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

If we look at our problem of finding a solution to

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

we can easily rewrite it as

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

where $\mathbf{Z} = \mathbf{X}^T \mathbf{X}$ is symmetric positive definite matrix and $\mathbf{b} = \mathbf{X}^T \mathbf{y}$. Because \mathbf{Z} can be factorized as $\mathbf{L}\mathbf{L}^T$ 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} is solved by forward substitution and \mathbf{w} is then solved by backward substitution which 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 $\mathbf{Z} = \mathbf{L}\mathbf{L}^T$ 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 15. 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 $\mathcal{O}(\frac{1}{3}p^3)$ [14] In the third and fourth step, we are solving triangular systems both of them require $\mathcal{O}(\frac{1}{2}p^2)$, so together $\mathcal{O}(p^2)$.

Putting all steps together we get

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

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

Computation OLS estimate using the Cholesky factorization is more numerically stable than using matrix inversion and time complexity is asymptotically similar. This approach is however not usually used in practice as well.

2.1.3 Computation using the QR decomposition

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

Definition 16. Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be any square matrix. Then it is possible to find 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 orthogonal matrix and $\mathbf{R} \in \mathbb{R}^{n \times n}$ is upper triangular matrix. This decomposition is known as *QR decomposition*.

If matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ thus 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 orthogonal matrix, $\mathbf{R}_1 \in \mathbb{R}^{n \times n}$ is upper triangular matrix, $\mathbf{0} \in \mathbb{R}^{(m-n) \times n}$ is zero matrix, $\mathbf{Q}_1 \in \mathbb{R}^{n \times n}$ is matrix with orthogonal columns and $\mathbf{Q}_2 \in \mathbb{R}^{(m-n) \times n}$ is also matrix with orthogonal columns.

2. ALGORITHMS

That means that \mathbf{X} can be factorized as

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

Then

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

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

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

Because $\mathbf{X} \in \mathbb{R}^{n \times p}$ is not square matrix then as we shown $\mathbf{R} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}$ so that

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

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

Solution to

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

is then given by

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

and because we assume \mathbf{X} have full column rank, thus is invertible, we can simplify it 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 triangular matrix, solution of \mathbf{w} is trivial using backward substitution. \mathbf{w} then represents our OLS estimate $\hat{\mathbf{w}}^{(OLS, n)}$.

Note 17. We can see that Cholesky factorization $\mathbf{X}^T \mathbf{X} = \mathbf{L}\mathbf{L}^T$ is closely connected to the QR decomposition $\mathbf{X} = \mathbf{Q}_1 \mathbf{R}_1$ so that

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

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

1. Calculate 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 $\mathbf{R}_1 \mathbf{w} = \mathbf{b}_1$ for \mathbf{w} .

QR factorization can be calculated in multiple ways. The most basic method is applying the Gram-Schmidt process to columns of matrix \mathbf{X} . This approach is not numerically stable so in practice is not used as much as two other methods. Those are *Householder transformations* and *Givens rotations*.

The time complexity of both algorithms is similar. QR decomposition of matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ is

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

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

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

so although Householder transformations are about 50% faster, both have asymptotically same 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. This property of Givens rotations we use in section 2.4.4.

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

Observation 18. The time complexity of solving OLS estimate using 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 QR decomposition. Householder transformations can be done in $\mathcal{O}(2np^2 - \frac{2}{3}p^3)$. Second step consist of matrix multiplication $\mathbf{Q}_1^T \mathbf{y}$ which is $\mathcal{O}(np)$. In the last step we are solving upper triangular systems 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 in case we use 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 OLS estimate exists, and that is 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 in terms of including and excluding row from factorized matrix \mathbf{X} which we use in the algorithm described in section 2.4.4.

[find citation](#)

2. ALGORITHMS

We compute the QR factorization of $\mathbf{A} \in \mathbb{R}^{m \times n}$ which has full column rank so that we apply orthogonal transformation by a matrix \mathbf{Q}^T as [15]

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

where \mathbf{Q} is product of orthogonal matrices. One such example of orthogonal matrix can be:

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

This matrix is indeed orthogonal since

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

Moreover, if we multiply this orthogonal matrix with some column vector $\mathbf{x} \in \mathbb{R}^{2 \times 1}$ thus $\mathbf{Q}\mathbf{x}$; as a result, we get a vector of the same length but which is rotated clockwise by φ radians. 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 column vector $\mathbf{x} \in \mathbb{R}^{m \times 1}$ 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 column vectors of \mathbf{A} , so we obtain zeros under diagonal. One such method - *Givens rotation* uses *Givens matrices* that zero one element under diagonal at a time.

The example of (2.24) is not random. Let us look at this orthogonal matrix 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 introduce this zeroing effect, we need to rotate this vector so that it is parallel to $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$. So we only need to find $\cos(\varphi)a + \sin(\varphi)b = r$. As we know rotation with \mathbf{Q} preserve L^2 norm. That means if we want to $z = 0$ then r must be equal to L^2 norm of the vector $\begin{bmatrix} a \\ b \end{bmatrix}$; thus $r = \sqrt{a^2 + b^2}$. Then the solution is trivial. We do not even need to calculate φ because if we put

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

and

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

then

$$r = \frac{a^2}{\sqrt{a^2 + b^2}} + \frac{b^2}{\sqrt{a^2 + b^2}} = \sqrt{a^2 + b^2} \quad (2.30)$$

$$z = \frac{-ab}{\sqrt{a^2 + b^2}} + \frac{ab}{\sqrt{a^2 + b^2}} = 0. \quad (2.31)$$

Practically used algorithm of computing $\cos \varphi$ and $\sin \varphi$ is slightly different because we want to prevent 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 this same idea to higher dimensions. Let us denote matrix $Q(i, j) \in \mathbb{R}^{m \times m}$ defined as

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

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 means of 2.28 and 2.29 for $a := x_i$ and $b := x_j$. Finally if we multiply $\mathbf{Q}(i, j)\mathbf{x} = \mathbf{p}$ we can see that \mathbf{p} is same as vector \mathbf{x} except p_i and p_j so that:

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

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. So if we have such matrix we need to create matrix $\mathbf{Q}(i, j)$ 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} then $a_{13} \dots a_{1n}$. Then we start with second column a_{23} and so on. That means we need to create in total $e = \frac{p^2 - p}{2} + np - p^2$ matrices $\mathbf{Q}(i, j)$. We can denote this sequence of matrices as $\mathbf{Q}_1, \mathbf{Q}_2, \dots, \mathbf{Q}_e$. The QR decomposition then looks like

$$\mathbf{Q}_e \dots \mathbf{Q}_2 \mathbf{Q}_1 \mathbf{A} = \mathbf{R} \quad (2.32)$$

where $\mathbf{Q}_e \dots \mathbf{Q}_2 \mathbf{Q}_1$ is actually \mathbf{Q}^T . So \mathbf{Q} is obtained by

$$\mathbf{Q} = \mathbf{Q}_1^T \mathbf{Q}_2^T \dots \mathbf{Q}_e^T. \quad (2.33)$$

We can see that for each row of the matrix, we need one additional matrix \mathbf{Q}_i , and with such matrix, we are multiplying only two rows. Moreover, the rows are getting shorter after each finished column. So the total number of operations can be expressed as

$$\sum_{i=1}^n \sum_{j=i+1}^p 6(n-i+1) \approx 6np^2 - 3np^2 - 3p^3 + 2p^3 = 3np^2 - p^3. \quad (2.34)$$

2.2 FAST-LTS

In this section, we introduce the FAST-LTS algorithm[16]. It is, as well as other algorithms we introduce, 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 which leading to better results.

Theorem 19. Consider $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^{n \times 1}$. Let us also have $\mathbf{w}_0 \in \mathbb{R}^p$ and $\mathbf{m}_0 \in Q^{(n,h)}$. Let us put $L_0 = \sum_{i=1}^n m_i^{(0)} r_i^2(\mathbf{w}_0)$ where $r_i^2(\mathbf{w}_0) = y_i - (w_1^0 x_1^i + w_2^0 x_2^i + \dots + w_p^0 x_p^i)$. Let us take $\hat{n} = \{1, 2, \dots, n\}$ and mark $\pi : \hat{n} \rightarrow \hat{n}$ permutation of \hat{n} such that $|r_0(\pi(1))| \leq |r_0(\pi(2))| \leq \dots \leq |r_0(\pi(n))|$ 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 subset with smallest absolute residuals $r_i^2(\mathbf{w}_0)$.

Finally, take $OF^{(OLS, \mathbf{M}_1 \mathbf{X}, \mathbf{M}_1 \mathbf{y})}(\mathbf{w})$ least squares fit on m_1 subset of observations and its corresponding $L_1 = \sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_1)$. Then

$$L_1 \leq L_0 \quad (2.35)$$

Proof. Because we take h observations with smallest absolute residuals $r_i^2(\mathbf{w}_0)$, then for sure $\sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_0) \leq \sum_{i=1}^n m_i^{(0)} r_i^2(\mathbf{w}_0) = L_0$. When we take into account that the OLS estimate minimizes the objective function of \mathbf{m}_1 subset of observations, then for sure $L_1 = \sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_1) \leq \sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_0)$.

Together we get $L_1 = \sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_1) \leq \sum_{i=1}^n m_i^{(1)} r_i^2(\mathbf{w}_0) \leq \sum_{i=1}^n m_i^{(0)} r_i^2(\mathbf{w}_0) = L_0$ \square

Corollary 20. Based on the previous theorem, using some $\hat{\mathbf{w}}_{old}^{OLS(h,n)}$ of \mathbf{m}_{old} subset of observations we can construct \mathbf{m}_{new} subset with corresponding $\hat{\mathbf{w}}_{new}^{OLS(H_{new})}$ 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 *C-step* of this process is described by the following pseudocode. Note that for C-step we need only $\hat{\mathbf{w}}$ without the need of passing \mathbf{m} .

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: $\hat{\mathbf{w}}_{new}$, 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  $m_{new} \leftarrow$  select set of  $h$  smallest absolute residuals from  $R$ ;
6  $\hat{\mathbf{w}}_{new} \leftarrow$  OLS fit on  $m_{new}$  subset; return  $\hat{\mathbf{w}}_{new}$ ,  $m_{new}$ ;

```

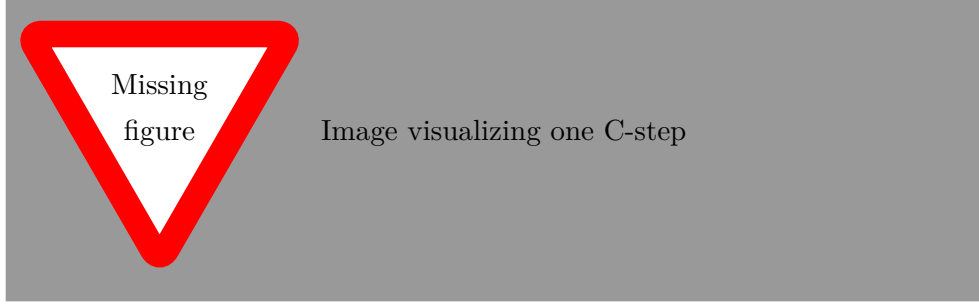


Figure 2.1: todo caption

One step of algorithm C-step is visualized on figure 2.1.

Observation 21. The time complexity of algorithm C-step 2 is asymptotically similar to the time complexity as OLS fit. Thus $O(p^2n)$.

Proof. In 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 $\mathcal{O}(p)$. So time of computation 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 algorithm QuickSelect [17]. Finally, we must compute \hat{w} OLS estimate on a h subset of data. Because h is linearly dependent on n , we can say that it is $\mathcal{O}(p^2n + p^3)$ which is asymptotically dominant against previous steps which are $\mathcal{O}(np + n)$. Because we assume $n \geq p$ then $\mathcal{O}(p^2n + p^3) \sim \mathcal{O}(p^2n)$ \square

As we stated above, repeating C-step leads to sequence of $\hat{\mathbf{w}}_1, \hat{\mathbf{w}}_2 \dots$ on subsets $\mathbf{m}_1, \mathbf{m}_2 \dots$ with corresponding $L_1 \geq L_2 \geq \dots$. 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 22. Sequence of C-step converges to $\hat{\mathbf{w}}_k$ after maximum of $k = \binom{n}{h}$ so that

$$L_k = L_i, \forall i \geq k. \quad (2.36)$$

Proof. Since $\text{OF}_D(\hat{\mathbf{w}}_i)$ is non-negative and $\text{OF}_D(\hat{\mathbf{w}}_i) \leq \text{OF}_D(\hat{\mathbf{w}}_{i+1})$ the sequence converges. $\hat{\mathbf{w}}_i$ is computed out of subset $\mathbf{m}_i \in Q^{(n,h)}$. Since $Q^{(n,h)}$ is finite, namely its size is $\binom{n}{h}$, the sequence converges 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    $L_{old} \leftarrow \text{OF}_D(\hat{\mathbf{w}}_{old})$ ;
6    $\hat{\mathbf{w}}_{new}, H_{new} \leftarrow \text{C-step}(\mathbf{X}, \mathbf{y}, \hat{\mathbf{w}}_{old})$ ;
7    $L_{new} \leftarrow \text{OF}_D(\hat{\mathbf{w}}_{new})$ ;
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 shallow, most often under 20 steps as can be seen on figure 2.2.

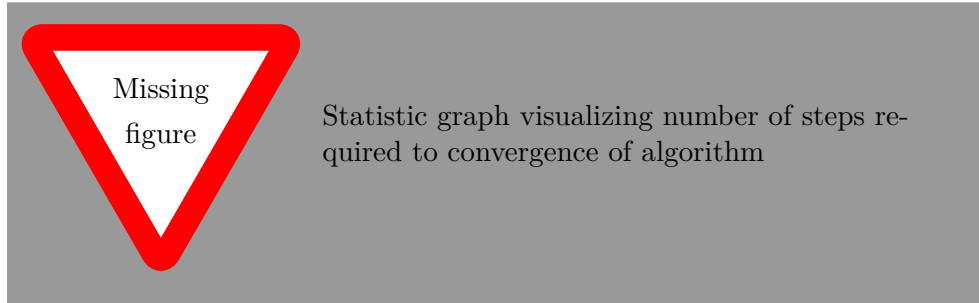


Figure 2.2: todo caption

That is not enough for the algorithm Repeat-C-step to converge to the global minimum. That gives us an idea of how to create the final algorithm. [16]

Choose a lot of initial subsets \mathbf{m}_1 and on each of them apply algorithm Repeat-C-step. From all converged subsets with corresponding $\hat{\mathbf{w}}$ estimates choose the one with the lowest value of $\text{OF}_D(\hat{\mathbf{w}})$.

Before we can construct final the algorithm, we must decide how to choose initial subset \mathbf{m}_1 and how many of them mean “a lot”. First, let us focus on how to choose an initial subset \mathbf{m}_1 .

2.2.2 Choosing 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 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 \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 23. With an increasing number of data samples, thus with increasing n , the probability of choosing among k random selections of $\mathbf{m}_{1_1}, \dots, \mathbf{m}_{1_k}$ the probability of selecting at least one \mathbf{m}_{1_i} such that its corresponding data samples does not contains outliers, goes to 0.

Proof. Consider dataset of n containing $\epsilon > 0$ relative amount of outliers. Let h be chosen conservatively so that $h = \lceil (n/2) \rceil + \lceil (p+1)/2 \rceil$ and k is the number of selections random h subsets. Then

$$P(\text{one random data sample not outliers}) = (1 - \epsilon)$$

$$P(\text{one subset without outliers}) = (1 - \epsilon)^h$$

$$P(\text{one subset with at least one outlier}) = 1 - (1 - \epsilon)^h$$

$$P(m \text{ subsets with at least one outlier in each}) = (1 - (1 - \epsilon)^h)^k$$

$$P(m \text{ subsets with at least one subset without outliers}) = 1 - (1 - (1 - \epsilon)^h)^k$$

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

That means that we should consider other options for selecting \mathbf{m}_1 subset. If we would like to continue with selecting some random subsets, previous observation gives us a clue, that we should choose it independent of n . Authors of the algorithm came with such a solution, and it goes as follows.

P-subset selection

Let us choose vector $\mathbf{c} \in Q^{(n,p)}$. Next we compute rank of matrix $\mathbf{X}_C = \mathbf{C}\mathbf{X}$, where $C = \text{diag}(c)$. If $\text{rank}(\mathbf{X}_C) < p$ add randomly selected rows to \mathbf{X}_C without repetition until $\text{rank}(\mathbf{X}_C) = p$. Let us from now on suppose that $\text{rank}(\mathbf{X}_C) = p$. Next let us mark $\hat{\mathbf{w}}_0 = \text{OF}_D(c)$ and corresponding $r_1(\hat{\mathbf{w}}_0), r_2(\hat{\mathbf{w}}_0), \dots, r_n(\hat{\mathbf{w}}_0)$ residuals. Now mark $\hat{n} = \{1, 2, \dots, n\}$ and let $\pi : \hat{n} \rightarrow \hat{n}$ be permutation of \hat{n} such that $|r(\pi(1))| \leq |r(\pi(2))| \leq \dots \leq |r(\pi(n))|$.

Finally, 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.

Observation 24. With the increasing number of data samples, thus with increasing n , the probability of choosing among k random selections of $\mathbf{c}_{1_1}, \dots, \mathbf{c}_{1_k}$ the probability of selecting at least one \mathbf{c}_{1_i} such that its corresponding data samples do not contains outliers, goes to

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

Proof. Similarly as in previous observation. \square

Last missing piece of the algorithm is determining the number of k initial \mathbf{m}_1 subsets, which maximize the probability to at least one of them converges to correct solution. 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 OLS fit on \mathbf{m}_i subset and then the calculation of n absolute residuals. How 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.

Due to that convergence of repeating C-step is very fast, it turns out, that we can distinguish between starts that leads to good solutions and those which does not, even after a small amount of the 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 $\text{OF}_D(\hat{\mathbf{w}}_3)$ or $\text{oflts}(\hat{\mathbf{w}}_4)$ respectively.

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 iteration of the C-steps till convergence only on those best solutions. Authors refer to this process as to *selective iteration*.

Nested extension

C-step computation is usually very fast for small n . Problem starts with very high n say $n > 10^3$ because we need to compute OLS on \mathbf{m}_i subset of size h which is dependent on n . And then calculate n absolute residuals.

Authors came up with a solution they call *nested extension*. We describe it briefly now.

2. ALGORITHMS

- If n is greater than limit l , we create subset of data samples L , $|L| = l$ and divide this subset into s disjunctive 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 number of starts $m_{P_i} = \frac{m}{l}$.
- Next in every P_i we create m_{P_i} number of initial $H_{P_{i_1}}$ subsets and iterate C-steps for two iterations.
- Then we choose 10 best results from each subsets and merge them together. We get family of sets F_{merged} containing 10 best $H_{P_{i_3}}$ subsets from each P_i .
- On each subset from F_{merged} family of subsets we again iterate 2 C-steps and then choose 10 best results.
- Finally we use these best 10 subsets and use them to iterate C-steps till convergence.
- As a result we choose best of those 10 converged results.

2.2.4 Putting all together

We 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 t 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 final C-step iteration till convergence. Parameter t denotes threshold stopping criterion such that $|\text{OF}_D(\hat{\mathbf{w}}_i) - \text{OF}_D(\hat{\mathbf{w}}_{i+1})| \leq t$ instead of $\text{OF}_D(\hat{\mathbf{w}}_i) = \text{OF}_D(\hat{\mathbf{w}}_{i+1})$. When we put all together, we get *FAST-LTS* algorithm which is described

by the following pseudocode.

Algorithm 4: FAST-LTS

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

```

1  $\hat{w}_{final} \leftarrow \emptyset$ ;
2  $m_{final} \leftarrow \emptyset$ ;
3  $F_{best} \leftarrow \emptyset$ ;
4 if  $n \geq l$  then
5    $F_{merged} \leftarrow \emptyset$ ;
6   for  $i \leftarrow 0$  to  $s$  do
7      $F_{selected} \leftarrow \emptyset$ ;
8     for  $j \leftarrow 0$  to  $\frac{l}{s}$  do
9        $F_{initial} \leftarrow \text{Selective iteration}(\frac{m}{l})$ ;
10      for  $m_i$  in  $F_{initial}$  do
11         $m_i \leftarrow \text{Iterate } C \text{ step few times}(m_i)$ ;
12         $F_{selected} \leftarrow F_{selected} \cup \{m_i\}$ ;
13      end
14    end
15     $F_{merged} \leftarrow F_{merged} \cup \text{Select 10 best subsets from } F_{selected}$ ;
16  end
17  for  $m_i$  in  $F_{merged}$  do
18     $m_i \leftarrow \text{Iterate } C \text{ step few times}(m_i)$ ;
19     $F_{best} \leftarrow F_{best} \cup \{m_i\}$ ;
20  end
21   $F_{best} \leftarrow \text{Select 10 best subsets from } F_{best}$ ;
22 else
23    $F_{initial} \leftarrow \text{Selective iteration}(m)$ ;
24    $F_{best} \leftarrow \text{Select 10 best subsets from } F_{initial}$ ;
25 end
26  $F_{final} \leftarrow \emptyset$ ;
27  $W_{final} \leftarrow \emptyset$ ;
28 for  $m_i$  in  $F_{best}$  do
29    $m_i, \hat{w}_i \leftarrow \text{Iterate } C \text{ step till convergence}(m_i, i_{max}, t)$ ;
30    $F_{final} \leftarrow F_{final} \cup \{m_i\}$ ;
31    $W_{final} \leftarrow W_{final} \cup \{\hat{w}_i\}$ ;
32 end
33  $\hat{w}_{final}, m_{final} \leftarrow$ 
    $\text{select } \hat{w}_i \text{ and } m_i \text{ with smallest } \text{OF}_D(\hat{w}_i) \text{ from } F_{final}$ ;
34 return  $\hat{w}_{final}, m_{final}$ ;

```

2.3 Feasible solution

In this section we introduce feasible solution algorithm from [6]. It is based on the strong necessary condition described at definition 10. The basic idea can be described as follows.

Let us consider that we have some $\mathbf{m} \in Q^{(n,h)}$ and 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 10.

With this in mind, let us mark

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

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 $\text{OF}_D(\mathbf{m})$ and the $\text{OF}_D(\mathbf{m}^{(i,j)})$ 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 $\mathbf{r}^{(m)} = \mathbf{Y} - \mathbf{X} \hat{\mathbf{w}}$ and also $d_{r,s} = \mathbf{x}_r \mathbf{Z}_M^{-1} \mathbf{x}_s$ then by equation introduced in [18] we get

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

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 subset. So let us assume that we have some vector $\mathbf{m} \in Q^{(n,h)}$. Now we compute $\Delta S_{i,j}^{(m)}$ for all $i \in O$ and $j \in Z$. This may lead to several different outcomes:

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

In the first case, all $\text{OF}_D(\mathbf{m}^{(i,j)})$ are higher or the same as the $\text{OF}_D(\mathbf{m})$ 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 $\text{OF}_D(\mathbf{m}^{(i,j)})$

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 $S_{i,j}^{(m)}$, 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 $S_{i,j}^{(m)}$ 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 subset satisfying the strong necessary condition. One step of this algorithm we call optimal

2. ALGORITHMS

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

Algorithm 5: 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 25. The time complexity of the OSAA 5 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 [6] they refer to this set as *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 subset resp. initial \mathbf{m} was already discussed when describing FAST-LTS algorithm. More importantly, as we already suggested h subset satisfying weak necessary condition do not need to satisfy strong necessary condition so passing such a h 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 subsets, e.g., random one.

Algorithm 6: 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     $\hat{\mathbf{w}} \leftarrow \text{OF}^{(\text{OLS}, \mathbf{M}\mathbf{X}, \mathbf{M}\mathbf{y})}$ ;
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       $\hat{\mathbf{w}} \leftarrow \text{OF}^{(\text{OLS}, \mathbf{M}\mathbf{X}, \mathbf{M}\mathbf{y})}$ ;
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  $\mathbf{m}_{\text{final}}, \hat{\mathbf{w}}_{\text{final}} \leftarrow \text{select best from } R \text{ based on smallest value of } \text{OF}_D$ ;
26 return  $\hat{\mathbf{w}}_{\text{final}}, \mathbf{m}_{\text{final}}$ ;

```

Observation 26. 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 13). 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 5 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 [3] 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 $\text{OF}_D(\mathbf{m}^{(i,j)}) - \text{OF}_D(\mathbf{m})$. 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 $\tilde{\mathbf{Z}} = \mathbf{A}^T \mathbf{A}$. 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

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

2. ALGORITHMS

Finally the inversion \mathbf{Z}^{-1} can be obtained using Sherman-Morrison formula [19] 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 $\bar{\hat{\mathbf{w}}}$ is

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

This can be simplified so that we get

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

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 $\gamma^+(\mathbf{c}_i) \geq 0$

mention the mistake in original paper? $w + (y - xw)bu$??

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

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

Let us mark

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

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 6 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 27. 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 26).

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

Finally, let us note that [3] calls 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 $\widehat{\mathbf{w}}$ 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 $\widehat{\mathbf{w}}$ to $\overline{\overline{\mathbf{w}}}$ 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 28. 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

2. ALGORITHMS

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 7: 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 29. 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 30. 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 31. The matrix $\bar{\mathbf{R}}$ can be updated to $\mathbf{R}^{(+)}$ without the presence of matrix \mathbf{Q} . We use this observation in the algorithm described in 2.6.

QR delete

When we extract the row \mathbf{x}_i from matrix \mathbf{A} we can use following trick [15]. First, we move such row as the first row of the matrix \mathbf{A} so we create permutation matrix \mathbf{P} so that

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

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

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

so that

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

Pseudocode is very similar to the QR Insert algorithm 7 so it is not necessary to provide it here.

Observation 32. 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 12 this can be done in $\mathcal{O}(n)$ for each matrix. So together we get $\mathcal{O}(n^2)$. \square

Note 33. 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 34. 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 $\tilde{\mathbf{w}}$ 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 $\bar{\mathbf{w}}$ to $\overline{\tilde{\mathbf{w}}}$ 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 the section . 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 35. This approach is slower than updating inversion directly. On the other hand, this solution is numerically stable. Updating inversion can be done in $\sqrt{\epsilon}$ (see observation 27) while QR insert has $\mathcal{O}(np)$ (see observation 29) and time complexity and QR delete even $\mathcal{O}(n^2)$ (see observation 32).

Note 36. 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 2.19 observation 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

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

Remark 37. 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 12; 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 21). 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 [3]. Very similar algorithm can also be found in [20].

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 calculation OLS 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

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.

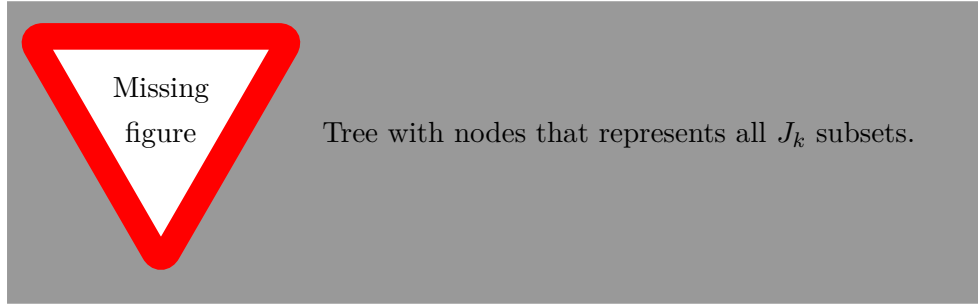


Figure 2.3: Tree consisting of all index subsets for $n = 4$ and $h = 3$

Le

2.6.1 Improvements

Sort data, so that the first h subset is smallest possible. We can use Any previous probabilistic algorithm. Moreover, a suitable solution may be some combined algorithm described in section 2.5.

Sorting rule. Compute partial increments. Go from the highest.

2.7 BSA algorithm

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

Experiments

3.1 Data

3.2 Results

3.3 Outlier detection

In this chapter we describe implementation of all algorithms as well as data sets which we use for experiments. Then a bunch of experiments will be presented. Right now we are finishing implementation of algorithms so we do not have any experimental results yet. This will change this or next week... Chapter will be about 4 pages long.

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] 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.
- [2] 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.
- [3] Agulló, J. New algorithms for computing the least trimmed squares regression estimator. *Computational Statistics & Data Analysis*, volume 36, no. 4, 2001: pp. 425–439.
- [4] Bernholt, T. Robust estimators are hard to compute. Technical report, Technical Report/Universität Dortmund, SFB 475 Komplexitätsreduktion, 2006.
- [5] Klouda, K. *Studium senzitivity odhadu metodou nejmenších usekaných čvterců*. Bachelor’s thesis, Technical University of Prague, Faculty of Nuclear Sciences and Physical Engineering, Department of Mathematics, 2006.
- [6] Hawkins, D. M. The feasible solution algorithm for least trimmed squares regression. *Computational statistics & data analysis*, volume 17, no. 2, 1994: pp. 185–196.
- [7] Bai, E.-W. A random least-trimmed-squares identification algorithm. *Automatica*, volume 39, no. 9, 2003: pp. 1651–1659.
- [8] Rousseeuw, P. J.; Leroy, A. M. *Robust regression and outlier detection*. John wiley & sons, 1987.

- [9] 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.
- [10] Strassen, V. Gaussian elimination is not optimal. *Numerische mathematik*, volume 13, no. 4, 1969: pp. 354–356.
- [11] Coppersmith, D.; Winograd, S. Matrix multiplication via arithmetic progressions. *Journal of symbolic computation*, volume 9, no. 3, 1990: pp. 251–280.
- [12] Ballard, G.; Benson, A. R.; et al. Improving the numerical stability of fast matrix multiplication. *arXiv preprint arXiv:1507.00687*, 2015.
- [13] 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).
- [14] Krishnamoorthy, A.; Menon, D. Matrix inversion using Cholesky decomposition. In *2013 Signal Processing: Algorithms, Architectures, Arrangements, and Applications (SPA)*, IEEE, 2013, pp. 70–72.
- [15] Hammarling, S.; Lucas, C. Updating the QR factorization and the least squares problem. 2008.
- [16] 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.
- [17] Hoare, C. A. Algorithm 65: find. *Communications of the ACM*, volume 4, no. 7, 1961: pp. 321–322.
- [18] 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.
- [19] Bartlett, M. S. An inverse matrix adjustment arising in discriminant analysis. *The Annals of Mathematical Statistics*, volume 22, no. 1, 1951: pp. 107–111.
- [20] 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.

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