

Insert here your thesis' task.



**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 20, 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 20, 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

V několika větách shrňte obsah a přínos této práce v českém jazyce.

Klíčová slova LTS odhad, lineární regrese, optimalizace, nejmenších čtverců, usekané čtverce, metoda nejmenších čtverců, outliers

Abstract

The least trimmed squares (LTS) method is a robust version of the classical method of least squares used to find an estimate of coefficients in the linear regression model. Computing the LTS estimate is known to be NP-hard, and hence suboptimal probabilistic algorithms are used in practice.

Keywords LTS, linear regression, robust estimator, least trimmed squares, ordinary least squares, outliers, outliers detection

Contents

Introduction	1
Introduction	3
1 Least trimmed squares	5
1.1 Linear regression model	5
1.2 Ordinary least squares	7
1.3 Robust statistics	9
1.4 Least trimmed squares	11
2 Algorithms	15
2.1 Computing OLS	17
2.2 Computing LTS	18
2.3 FAST-LTS	18
2.4 Feasible solution	25
2.5 Combined algorithm	27
2.6 Improved FSA	28
2.7 Branch and bound aka. BAB	44
2.8 BSA algorithm	44
3 Experiments	45
3.1 Data	45
3.2 Results	45
3.3 Outlier detection	45
Conclusion	47
Bibliography	49
A Data sets	51

B Contents of enclosed CD

53

List of Figures

Introduction

Goal of this thesis is to research currently used algorithms for calculation Least trimmed squares. Least trimmed squares was first introduced in 1984 and since that time couple of researchers came up couple of solutions. Beside that we'll propose extension of current algorithms which was never used before. We'll show that this updated algorithm is fast enough and also gives same or better results than currently used algorithms. We'll also compare speed and performance of this algorithms on various data sets both from literature and also using our custom data generator which provides very flexible way of generating data with outliers. We'll extend all those algorithms and implement all of them both in C++ and python. We'll implement python library with C++ back end which will provide all of the currently used algorithms.

In first chapter we'll introduce linear regression and ordinary least squares method and its downfalls. We'll introduce robust statistic and methods of evaluating robust models. In second chapters we'll analyze all of the algorithms together with its time complexity etc

Conclusion

There is still lot of future works on least trimmed squares. Proof of couple of thoughts is still about to come.

Introduction

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 $\varepsilon \sim \mathcal{N}(0, \sigma^2)$.

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

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 .

1.1.1.1 Intercept

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

Example 3. Let us consider that y represents price of the room and $x \in \mathbb{N}$ 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 3. 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 4. 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*

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

Definition 5. 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}^* .

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

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

1. LEAST TRIMMED SQUARES

By this we obtain the gradient

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

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:

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

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 of 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 $\text{cov}(\varepsilon_i, \varepsilon_j) = 0$, $i, j = 1, 2, \dots, n$, $i \neq j$
- 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 6. As described in example 3 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.

1.2.2 Computing OLS

In this section, we describe a few of many methods that can be used to obtain $\hat{\mathbf{w}}^{(OLS,n)}$.

Matrix inversion

The $\hat{\mathbf{w}}^{(OLS,n)}$ can be computed directly by multiplying $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ and $\mathbf{B} = \mathbf{X}^T \mathbf{y}$ and then by finding inversion of \mathbf{A} and finally multiplying $\mathbf{A}^{-1} \mathbf{B}$.

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 7. 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

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

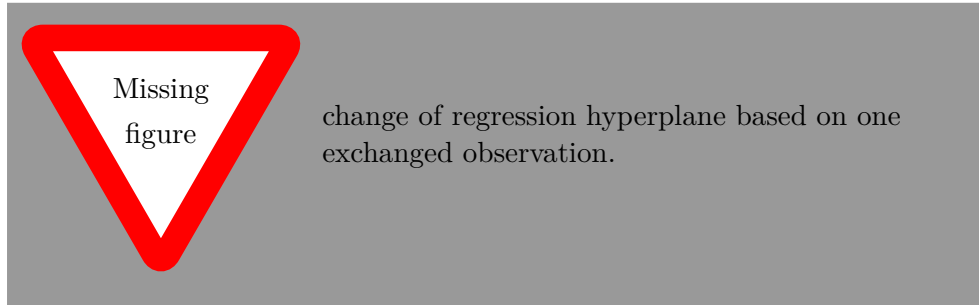
This definition is very general but let us specify it to our linear regression problem. It says that breakdown point is the function of the minimal number of observations needed to be changed for some others so then the estimator gives incorrect results. More robust estimators have higher breakdown point.

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 8. 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

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

where $r_{i:n}^2$ denotes i th smallest squared residuum out of n thus $r_{1:n}^2 \leq r_{2:n}^2 \dots r_{n:n}^2$. 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 8 is non-differentiable and non-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

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

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]

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

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

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

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

$$J(\mathbf{m}) = \left\| \mathbf{M}\mathbf{y} - \mathbf{M}\mathbf{X}(\mathbf{X}^T \mathbf{M}\mathbf{X})^{-1} \mathbf{X}^T \mathbf{M}\mathbf{y} \right\|^2. \quad (1.27)$$

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 $h = \lfloor (n/2) + \lfloor (p+1)/2 \rfloor$ [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 previous chapter we've covered necessary theory needed to implement algorithms that are in this chapter. Let's quickly recap most important fact that we know so far.

With robust linear regression problem we assume that our model with intercept

$$y_i = \mathbf{w}^T \mathbf{x}_i + \varepsilon_i \quad (2.1)$$

where $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ is *i.i.d.* random variable has some displacements in up-most half of the explanatory \mathbf{x}_i or dependent y_i variables. Thus only some subset $\tilde{\mathbf{X}} = \mathbf{M}\mathbf{X}$ with corresponding $\tilde{\mathbf{y}} = \mathbf{y}$ where $\mathbf{M} = \text{diag}(\mathbf{m})$ and $\mathbf{m} \in Q^{(n,h)}$ can be perceived so that

$$\tilde{y}_i \sim \mathcal{N}(\mathbf{w}^T \tilde{\mathbf{x}}_i, \sigma^2). \quad (2.2)$$

We'll denote this subset simply as *hsubset* of data in order to simplify following text. We've also learnt that finding solution of LTS we need to find correct h subset and then calculate estimate of regression coefficients \tilde{w} . So to find exact solution of LTS we need to go through all h subsets, find the correct one and calculate ordinary least squares fit to get the regression coefficients. There are not much of other options because objective function is non-differentiable and non-convex with lots of local minima.

We also know that exhaustive approach will fail due exponential size of $Q^{(n,h)}$. So what are the possibilities? First attempts were based on iterative removal data samples whose residuum had the highest value based on OLS fit on whole dataset. Such attempts were proven to be completely wrong because initial OLS fit is usually already heavily affected by outliers and we may end up removing data samples which represents original model.

Then there are algorithms based on purely on random approach.

On such algorithm is Random solution algorithm [6] which is basically randomly selects L h subsets and subsequently compute OLS fit on each of them and selecting fit with smallest *RSS* and consider it a approximate solution.

2. ALGORITHMS

Such approach is very simple, but in general probability of selecting at least one such subset from L subsets which don't contain outliers thus has a chance of producing good result goes to zero for increasing number of data samples n as we'll describe in detail in 18.

Another very similar algorithm called Resampling algorithm introduced in [7] have basically just a little difference and that it select vectors from $Q^{(n,p+1)}$ instead of $Q^{(n,h)}$. This minor tweak has not only higher chance to succeed because number of vectors in this set is significantly lower than in $Q^{(n,h)}$ (at least if h conservatively chosen thus $h = \lfloor n/2 \rfloor + \lfloor (p+1)/2 \rfloor$) but also because probability of selecting L subsets of size which is independent of n gives nonzero probability of selecting at least one subset such that it don't include outliers see 19 for more details.

Generating all possible h subsets is computationally hard and relying on selecting random subsets don't produce sufficiently good results. So what are our options? In [8] two criterions called *necessary conditions* are introduced. They talk about necessary properties which some h subset must satisfy so it could be set which leads to global optima of LTS. Let's introduce those two necessary conditions. For that it's convenient not to only label h subset of used observations but also complementary subset of not used observations. We'll refer to this complementary subset as g subset.

Theorem 9. Strong necessary condition. The criterion cannot be improved by exchanging any of the observations from g subset for any of the currently used observations in h subset. Thus $\mathbf{m} \in Q^{(n,h)}$ meets the criterion if $J(\mathbf{m}) \leq J(\mathbf{m}_{\text{swap}})$ where \mathbf{m}_{swap} is any vector from $Q^{(n,h)}$ such that it has same values except one swapped.

Proof. Trivial. LTS uses subset of h observations that minimize it's objective function. To be this true none of the swaps between observations from h subset and g subset must not improve (reduce) it's objective function. \square

Based on this idea algorithm can be created. We'll discuss it in detail in 2.4.

Second necessary condition named *weak necessary condition*

Theorem 10. $\mathbf{m} \in Q^{(n,h)}$ meets the criterion if for each observation from h subset has smaller (or equal) squared residual than any observation from g subset.

Again, based on this criterion an algorithm can be crated. Corollary of this criteria together with proof can be found in 2.3.

Very interesting consequence which we'll use later gives us following lemma.

Lemma 11. Strong necessary condition is not satisfied unless weak necessary condition in satisfied. Thus if strong condition is satisfied then weak is also.

Proof. We'll make proof by contradiction. Let's assume that we have $\mathbf{m} \in Q^{(n,h)}$ and $J(\mathbf{m})$ for which strong necessary condition is satisfied but weak

necessary condition is not. That means there exists \mathbf{x}_i with y_i from h subset and \mathbf{x}_j and y_j from g subset such that $r_j^2 < r_i^2$. Thus

$$J(\mathbf{m}) > J(\mathbf{m}) + r_j^2 - r_i^2 \quad (2.3)$$

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

$$J(\mathbf{m}) + r_j^2 - r_i^2 \geq J(\mathbf{m}_{\text{swap}}). \quad (2.4)$$

That's indeed trivial because $J(\mathbf{m}_{\text{swap}})$ is in fact just OLS that minimize objective function on given subset of observations. That's of course contradiction with our assumption which says that strong necessary condition is already satisfied. \square

When we'll discuss algorithms based on this conditions we'll show that algorithm based on weak necessary condition is much faster than algorithm based on strong necessary condition which will lead us to another algorithm where we'll use 11.

Now we've covered all necessary theoretical background and it's time to introduce currently popular algorithms of computing LTS estimate.

2.1 Computing OLS

In this section we'll describe what algorithms of computing OLS exists. We'll see that beside computing OLS directly from the objective function there are better ways. Let's now start with this straightforward approach.

Lemma 12. Time complexity of OLS on $\mathbf{X}^{n \times p}$ and $\mathbf{Y}^{n \times 1}$ is $O(p^2n)$.

Proof. Normal equation of OLS is $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$. Time complexity of matrix multiplication $\mathbf{A}^{m \times n}$ and $\mathbf{B}^{n \times p}$ is $\sim \mathcal{O}(mnp)$. Time complexity of matrix $\mathbf{C}^{m \times m}$ is $\sim \mathcal{O}(m^3)$ So we need to compute $\mathbf{A} = \mathbf{X}^T \mathbf{X} \sim \mathcal{O}(p^2n)$ and $\mathbf{B} = \mathbf{X}^T \mathbf{Y} \sim \mathcal{O}(pn)$ and $\mathbf{C} = \mathbf{A}^{-1} \sim \mathcal{O}(p^3)$ and finally $\mathbf{CB} \sim \mathcal{O}(p^2)$. That gives us $\mathcal{O}(p^2n + pn + p^3 + p^2)$. Because $\mathcal{O}(p^2n)$ and $\mathcal{O}(p^3)$ asymptotically dominates over $\mathcal{O}(p^2)$ and $\mathcal{O}(pn)$ we can write $\mathcal{O}(p^2n + p^3)$.

\square

17

CO zo toho je vic?
Neni casove naroc-
nejši vynasobeni
 $\mathbf{X}^T \mathbf{X}$ nez inver-
sion, kdyz bereme
v uvahu $n \gg p$
???

complete this sec-
tion with describ-
ing computation of
OLS using matrix
decomposition

2.2 Computing LTS

Note 13. When discussing following algorithms, we'll refer to given \mathbf{X} and y as to **data set** and to y_i with corresponding \mathbf{x}_i as to **data sample** or **observation**. Sometimes it's also useful to refer to multiple observations as to subset of observations. When we want to mark subset of observations $y_i, \mathbf{x}_i, i \in H, H \subset \{1, 2, \dots, n\}$ we can simply refer to it as to subset of observations H . Sometimes it's also useful to mark matrix \mathbf{X} with only some subset of observations which we'll do by \mathbf{X}_H .

2.3 FAST-LTS

In this section we will introduce FAST-LTS algorithm[9]. It is, as well as in other cases, iterative algorithm. We will discuss all main components of the algorithm starting with its core idea called concentration step which ' authors simply call C-step.

2.3.1 C-step

We will show that from existing LTS estimate $\hat{\mathbf{w}}_{old}$ we can construct new LTS estimate $\hat{\mathbf{w}}_{new}$ which objective function is less or equal to the old one. Based on this property we will be able to create sequence of LTS estimates which will lead to better results.

Theorem 14. Consider dataset consisting of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ explanatory variables where $\mathbf{x}_i \in \mathbb{R}^p, \forall \mathbf{x}_i = (x_1^i, x_2^i, \dots, x_p^i)$ where $x_1^i = 1$ and its corresponding y_1, y_2, \dots, y_n response variables. Let's also have $\hat{\mathbf{w}}_0 \in \mathbb{R}^p$ any p -dimensional vector and $H_0 = \{h_i; h_i \in \mathbb{Z}, 1 \leq h_i \leq n\}, |H_0| = h$. Let's now mark $RSS(\hat{\mathbf{w}}_0) = \sum_{i \in H_0} (r_0(i))^2$ where $r_0(i) = y_i - (w_1^0 x_1^i + w_2^0 x_2^i + \dots + w_p^0 x_p^i)$. Let's 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 $H_1 = \{\pi(1), \pi(2), \dots, \pi(h)\}$ set of h indexes corresponding to h smallest absolute residuals $r_0(i)$. Finally take $\hat{\mathbf{w}}_1^{OLS(H_1)}$ ordinary least squares fit on H_1 subset of observations and its corresponding $RSS(\hat{\mathbf{w}}_1) = \sum_{i \in H_1} (r_1(i))^2$ sum of least squares. Then

$$RSS(\hat{\mathbf{w}}_1) \leq RSS(\hat{\mathbf{w}}_0) \quad (2.5)$$

Proof. Because we take h observations with smallest absolute residuals r_0 , then for sure $\sum_{i \in H_1} (r_0(i))^2 \leq \sum_{i \in H_0} (r_0(i))^2 = RSS(\hat{\mathbf{w}}_0)$. When we take into account that Ordinary least squares fit OLS_{H_1} minimize objective function of H_1 subset of observations, then for sure $RSS(\hat{\mathbf{w}}_1) = \sum_{i \in H_1} (r_1(i))^2 \leq \sum_{i \in H_1} (r_0(i))^2$. Together we get

$$RSS(\hat{\mathbf{w}}_1) = \sum_{i \in H_1} (r_1(i))^2 \leq \sum_{i \in H_1} (r_0(i))^2 \leq \sum_{i \in H_0} (r_0(i))^2 = RSS(\hat{\mathbf{w}}_0)$$

□

Corollary 15. Based on previous theorem, using some $\hat{\mathbf{w}}^{OLS(H_{old})}$ on H_{old} subset of observations we can construct H_{new} subset with corresponding $\hat{\mathbf{w}}^{OLS(H_{new})}$ such that $RSS(\hat{\mathbf{w}}^{OLS(H_{new})}) \leq RSS(\hat{\mathbf{w}}^{OLS(H_{old})})$. With this we can apply above theorem again on $\hat{\mathbf{w}}^{OLS(H_{new})}$ with H_{new} . This will lead to the iterative sequence of $RSS(\hat{\mathbf{w}}_{old}) \leq RSS(\hat{\mathbf{w}}_{new}) \leq \dots$. One step of this process is described by following pseudocode. Note that for C-step we actually need only $\hat{\mathbf{w}}$ without need of passing H .

Algorithm 1: 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: $\hat{\mathbf{w}}_{new}$, H_{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 $H_{new} \leftarrow$ select set of h smallest absolute residuals from R ;
- 6 $\hat{\mathbf{w}}_{new} \leftarrow OLS(H_{new})$;
- 7 **return** $\hat{\mathbf{w}}_{new}$, H_{new} ;

include and reference nice image showing one c-step

Observation 16. Time complexity of algorithm C-step 1 is the same as time complexity as OLS. 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 $\mathcal{O}(p)$. Rest is in constant time. So time of computation n residuals is $\mathcal{O}(np)$. Next we must select set of h smallest residuals which can be done in $\mathcal{O}(n)$ using modification of algorithm QuickSelect. Finally we must compute $\hat{\mathbf{w}}$ OLS estimate on 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 $\mathcal{O}(np + n)$. □

create better proof. And take into account both versions - directly vs. using decomposition

reference or define quick select

As we stated above, repeating algorithm C-step will lead to sequence of $\hat{\mathbf{w}}_1, \hat{\mathbf{w}}_2 \dots$ on subsets $H_1, H_2 \dots$ with corresponding residual sum of squares $RSS(\hat{\mathbf{w}}_1) \geq RSS(\hat{\mathbf{w}}_2) \geq \dots$. One could ask if this sequence will converge, so that $RSS(\hat{\mathbf{w}}_i) == RSS(\hat{\mathbf{w}}_{i+1})$. Answer to this question will be presented by the following theorem.

Theorem 17. Sequence of C-step will converge to $\hat{\mathbf{w}}_m$ after maximum of $m = \binom{n}{h}$ so that $RSS(\hat{\mathbf{w}}_m) == RSS(\hat{\mathbf{w}}_n)$, $\forall n \geq m$ where n is number of data samples and h is size of subset H_i .

Proof. Since $RSS(\hat{\mathbf{w}}_i)$ is non-negative and $RSS(\hat{\mathbf{w}}_i) \leq RSS(\hat{\mathbf{w}}_{i+1})$ the sequence will converge. $\hat{\mathbf{w}}_i$ is computed out of subset $H_i \subset \{1, 2, \dots, n\}$. When

2. ALGORITHMS

there is finite number of subsets of size h out of n samples, namely $\binom{n}{h}$, the sequence will converge at the latest after this number of steps. \square

Above theorem gives us clue to create algorithm described by following pseudocode.

Algorithm 2: Repeat-C-step

Input: dataset consisting of $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^{n \times 1}$, $\hat{\mathbf{w}}_{old} \in \mathbb{R}^{p \times 1}$, H_0
Output: $\hat{\mathbf{w}}_{final}$, H_{final}

```
1  $\hat{\mathbf{w}}_{new} \leftarrow \emptyset$ ;  
2  $H_{new} \leftarrow \emptyset$ ;  
3  $RSS_{new} \leftarrow \infty$ ;  
4 while True do  
5    $RSS_{old} \leftarrow RSS(\hat{\mathbf{w}}_{old})$ ;  
6    $\hat{\mathbf{w}}_{new}, H_{new} \leftarrow \mathbf{X}, \mathbf{y}, \hat{\mathbf{w}}_{old}$ ;  
7    $RSS_{new} \leftarrow RSS(\hat{\mathbf{w}}_{new})$ ;  
8   if  $RSS_{old} == RSS_{new}$  then  
9     break  
10  end  
11   $\hat{\mathbf{w}}_{old} \leftarrow \hat{\mathbf{w}}_{new}$   
12 end  
13 return  $\hat{\mathbf{w}}_{new}, H_{new}$ ;
```

It is important to note, that although maximum number of steps of this algorithm is $\binom{n}{h}$ in practice it is very low, most often under 20 steps. That is not enough for the algorithm *Repeat-C-step* to converge to global minimum, but it is necessary condition. That gives us an idea how to create the final algorithm. [9]

Choose a lot of initial subsets H_1 and on each of them apply algorithm Repeat-C-step. From all converged subsets with corresponding $\hat{\mathbf{w}}$ estimates choose that which has lowest $RSS(\hat{\mathbf{w}})$.

Before we can construct final algorithm we must decide how to choose initial subset H_1 and how many of them mean “a lot of”. First let’s focus on how to choose initial subset H_1 .

2.3.2 Choosing initial H_1 subset

It is important to note, that when we choose H_1 subset such that it contains outliers, then iteration of C-steps usually won’t converge to good results, so we should focus on methods with non zero probability of selecting H_1 such that it won’t contain outliers. There are a lot of possibilities how to create initial hH_1 subset. Lets start with most trivial one.

include sime nice
grap which show
this. or table ?

2.3.2.1 Random selection

Most basic way of creating H_1 subset is simply to choose random $H_1 \subset \{1, 2, \dots, n\}$. Following observation will show that it not the best way.

Observation 18. With increasing number of data samples, thus with increasing n , the probability of choosing among m random selections of H_{1_1}, \dots, H_{1_m} the probability of selecting at least one H_{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 $h = (n + p + 1)/2$ and m is number of selections random $|H| = h$ subsets. Then

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

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

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

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

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

Because $n \rightarrow \infty \Rightarrow (1 - \epsilon)^h \rightarrow 0 \Rightarrow 1 - (1 - \epsilon)^h \rightarrow 1 \Rightarrow (1 - (1 - \epsilon)^h)^m \rightarrow 1 \Rightarrow 1 - (1 - (1 - \epsilon)^h)^m \rightarrow 0$ \square

That means that we should consider other options of selecting H_1 subset. Actually if we would like to continue with selecting some random subsets, previous observation gives us clue, that we should choose it independent of n . Authors of algorithm came with such solution and it goes as follows.

2.3.2.2 P-subset selection

Let's choose subset $J \subset \{1, 2, \dots, n\}, |J| = p$. Next compute rank of matrix $\mathbf{X}_{J:}$. If $\text{rank}(\mathbf{X}_{J:}) < p$ add randomly selected rows to $\mathbf{X}_{J:}$ without repetition until $\text{rank}(\mathbf{X}_{J:}) = p$. Let's from now on suppose that $\text{rank}(\mathbf{X}_{J:}) = p$. Next let us mark $\hat{\mathbf{w}}_0 = OLS(J)$ and corresponding $(r_0(1)), (r_0(2)), \dots, (r_0(n))$ 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 put $H_1 = \{\pi(1), \pi(2), \dots, \pi(h)\}$ set of h indexes corresponding to h smallest absolute residuals $r_0(i)$.

Observation 19. With increasing number of data samples, thus with increasing n , the probability of choosing among m random selections of J_{1_1}, \dots, J_{1_m} the probability of selecting at least one J_{1_i} such that its corresponding data samples does not contains outliers, goes to

$$1 - (1 - (1 - \epsilon)^h)^m > 0$$

Proof. Similarly as in previous observation. \square

Note that there are other possibilities of choosing H_1 subset other than these presented in [9]. We'll properly discuss them in chapter

Last missing piece of the algorithm is determining number of m initial H_1 subsets, which will maximize probability to at least one of them will converge to good solution. Simply put, the more the better. So before we will answer this question properly, let's discuss some key observations about algorithm.

2.3.3 Speed-up of the algorithm

In this section we will describe important observations which will help us to formulate final algorithm. In two subsections we'll briefly describe how to optimize current algorithm.

2.3.3.1 Selective iteration

The most computationally demanding part of one C-step is computation of OLS on H_i subset and then 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 in the same time didn't lose performance of algorithm.

Due to that convergence of repeating C-step is very fast, it turns out, that we are able to distinguish between starts that will lead to good solutions and those who won't even after very little 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 $RSS(\hat{w}_3)$ or $RSS(\hat{w}_4)$ respectively.

So even though authors don't specify size of m explicitly, they propose that after a few C-steps we can choose (say 10) best solutions among all H_1 starts and continue C-steps till convergence only on those best solutions. This process is called Selective iteration.

We can choose m with respect to observation 19. In ideal case we would like to have probability of existence at least one initial H_1 subset close to 1. As we see m is exponentially dependent on p and at the same time in practice we don't know percentage of outliers in dataset. So it is difficult to mention exact value. Specific values of m in respect to data size is visible in table . So we can say that with $p < 10$ choosing $m = 500$ is usually safe starting point.

2.3.3.2 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 H_i subset of size h which is dependent on n . And then calculate n absolute residuals.

Authors came up with solution they call Nested extension. We will describe it briefly now.

- If n is greater than limit l , we'll 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'll set number of starts $m_{P_i} = \frac{m}{l}$.
- Next in every P_i we'll create m_{P_i} number of initial $H_{P_{i_1}}$ subsets and iterate C-steps for two iterations.
- Then we'll choose 10 best results from each subsets and merge them together. We'll 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'll again iterate 2 C-steps and then choose 10 best results.
- Finally we'll use these best 10 subsets and use them to iterate C-steps till convergence.
- As a result we'll choose best of those 10 converged results.

2.3.3.3 Putting all together

We've described all major parts of the algorithm FAST-LTS. One last thing we need to mention is that even though C-steps iteration usually converge under 20 steps it is appropriate to introduce two parameters *max_iteration* and *threshold* which will limit number of C-steps iterations in some rare cases when convergence is too slow. Parameter *max_iteration* denotes maximum number of iterations in final C-step iteration till convergence. Parameter *threshold* denotes stopping criterion such that $|RSS(\hat{\mathbf{w}}_i) - RSS(\hat{\mathbf{w}}_{i+1})| \leq threshold$ instead of $RSS_i == RSS_{i+1}$. When we put all together, we'll get

2. ALGORITHMS

FAST-LTS algorithm which is described by following pseudocode.

Algorithm 3: FAST-LTS

Input: $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{y} \in \mathbb{R}^{n \times 1}$, $m, l, s, \max_iteration, threshold$
Output: $\hat{\mathbf{w}}_{final}, H_{final}$

```

1  $\hat{\mathbf{w}}_{final} \leftarrow \emptyset;$ 
2  $H_{final} \leftarrow \emptyset;$ 
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  $H_i$  in  $F_{initial}$  do
11         $H_i \leftarrow \text{Iterate } C \text{ step few times}(H_i);$ 
12         $F_{selected} \leftarrow F_{selected} \cup \{H_i\};$ 
13      end
14    end
15     $F_{merged} \leftarrow F_{merged} \cup \text{Select 10 best subsets from } F_{selected};$ 
16  end
17  for  $H_i$  in  $F_{merged}$  do
18     $H_i \leftarrow \text{Iterate } C \text{ step few times}(H_i);$ 
19     $F_{best} \leftarrow F_{best} \cup \{H_i\};$ 
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  $H_i$  in  $F_{best}$  do
29    $H_i, \hat{\mathbf{w}}_i \leftarrow$ 
30      $\text{Iterate } C \text{ step till convergence}(H_i, \max\_iteration, threshold);$ 
31    $F_{final} \leftarrow F_{final} \cup \{H_i\};$ 
32    $W_{final} \leftarrow W_{final} \cup \{\hat{\mathbf{w}}_i\};$ 
33 end
34  $\hat{\mathbf{w}}_{final}, H_{final} \leftarrow \text{select what with best RSS}(F_{final}, W_{final});$ 
35 return  $\hat{\mathbf{w}}_{final}, H_{final};$ 

```

2.4 Feasible solution

In this section we'll introduce feasible solution algorithm from [10]. It is based on strong necessary condition we've described at 9. The basic idea can be described as follows.

Let's consider that we have some $\mathbf{m} \in Q^{(n,h)}$. It'll be convenient when we'll mark in the following text $O_m = \{i \in \{1, 2, \dots, n\}; w_i = 1\}$ 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 correspond to vector \mathbf{m}_{swap} which we marked at 9.

With this in mind we can mark let's mark

$$\Delta S_{i,j}^{(m)} = J(\mathbf{m}^{(i,j)}) - J(\mathbf{m}) \quad (2.6)$$

thus change of the LTS objective function by swapping one observation from h subset with another from g subset. To calculate this we can obviously first calculate $J(\mathbf{m})$ and $J(\mathbf{m}^{(i,j)})$ and finally subtract both results. Although it is a option, it is computationally hard. So the question is if there is easier way of calculating $\Delta S_{i,j}^{(m)}$ and the answer is positive.

Let's mark $M = \text{diag}(\mathbf{m})$ and $M^{(i,j)} = \text{diag}(\mathbf{m}^{(i,j)})$ and also $\mathbf{H} = (\mathbf{X}^T \mathbf{M}^T \mathbf{X})$ For now let's assume that we already calculated \mathbf{H}^{-1} and also $\hat{\mathbf{w}} = \mathbf{H}^{-1} \mathbf{X}^T \mathbf{M} \mathbf{y}$ we now want to calculate $\Delta S_{i,j}^{(m)}$ Let's mark vector of residuals $\mathbf{r}^{(m)} = \mathbf{Y} - \mathbf{X} \hat{\mathbf{w}}$ and also $d_{r,s} = \mathbf{x}_r \mathbf{H}^{-1} \mathbf{x}_s$ then by equation introduced in [11] we get

$$\Delta S_{i,j}^{(m)} = \frac{(r_j^{(m)})^2(1 - d_{i,i}) - (r_i^{(m)})^2(1 + d_{j,j}) + 2r_i^{(m)}r_j^{(m)}d_{j,j}}{(1 - d_{i,i})(1 + d_{j,j}) + d_{i,j}^2}. \quad (2.7)$$

Let's now describe core of the algorithm. It's a similar to the FAST-LTS algorithm in the sense of iterative refinement of h subset. So let's assume that we have some vector $\mathbf{m} \in Q^{(n,h)}$. Now we'll 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 positive
3. multiple $S_{i,j}^{(m)}$ are positive

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

2. ALGORITHMS

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

1. use the first swap that leads to the improvement (so don't even try to find different swap)
2. from all possible swaps choose one that has highest improvement value $J(\mathbf{m}^{(i,j)})$
3. use the first swap that has improvement higher than some given threshold

O(n²p) or O(n³p)
remains

In the terms of complexity all three options give the same results, because to find feasible solution you need to evaluate all pair swaps. In practice third options is winner because it'll lead to least amount of iterations. On the other hand as we said this won't improve complexity of the algorithm, so from now on let's assume that we'll use case number two. So if there positive $S_{i,j}^{(m)}$ we'll make the swap and repeat the process again. Algorithm ends when there is no possible improvement i.e. when all $S_{i,j}^{(m)}$ are non-negative. Number of iterations needed till algorithm stops is usually quite low, but for practical usage it's still convenient to use some parameter *max_iteration* after which algorithm will stop without finding h subset satisfying strong necessary condition.

put sem part of
the pseudocode
describing core
iteration?

and its time com-
pleixty

One run of this iteration process will lead to some local optima i.e. set satisfying strong necessary condition. In [10] they refer to this set as to *feasible set*. This because it is not global optima the algorithm needs to be run multiple times say N times. As a final solution is taken such h subset with the smallest residual sum of squares.

experiment with
this and refer here

We didn't yet mention how to create initial h subset respectively initial \mathbf{m} . We already had this discussion when describing FAST-LTS algorithm. The [10] describes only random starting h subsets, but using p subsets instead may lead to improvement. More importantly as we already suggested h subset satisfying weak necessary condition don't need to satisfy strong necessary condition so passing such a h subset as input to this algorithm is another option and we'll discuss it in detail later. For that reasons we'll now describe feasible algorithm with pseudocode and we'll assume that we already have

some function that generates for us h subsets e.g. random one.

Algorithm 4: Feasible solution

Input: $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{y} \in \mathbb{R}^{n \times 1}$, $max_iteration$, N
Output: $\hat{\mathbf{w}}_{final}$, H_{final}

```

1  $\hat{\mathbf{w}}_{final} \leftarrow \emptyset$ ;
2  $H_{final} \leftarrow \emptyset$ ;
3  $Results \leftarrow \emptyset$ ;
4 for  $k \leftarrow 0$  to  $N$  do
5    $\mathbf{m} \leftarrow generate\_intial\_subset()$  ;           // e.g. random  $\mathbf{m} \in Q^{(n,h)}$ 
6   while  $True$  do
7      $best_i \leftarrow 0$  ;
8      $best_j \leftarrow 0$  ;
9      $best_{delta} \leftarrow 0$  ;
10    for  $i \in O_m$  do
11      for  $j \in Z_m$  do
12         $delta \leftarrow calculate\Delta S_{i,j}^{(m)}$ ;
13        if  $delta > best_{delta}$  then
14           $best_{delta} \leftarrow delta$ ;
15           $best_i \leftarrow i$  ;
16           $best_j \leftarrow j$  ;
17        end
18      end
19    end
20    if  $best_{delta} > 0$  then
21       $\mathbf{m} \leftarrow \mathbf{m}^{(i,j)}$ ;
22    end
23    else
24       $H \leftarrow h$  subset corresponding to  $\mathbf{m}$ ;
25       $\mathbf{M} \leftarrow diag(\mathbf{m})$ ;
26       $\hat{\mathbf{w}} \leftarrow OF^{(OLS, \mathbf{M}\mathbf{X}, \mathbf{M}\mathbf{y})}$ ;
27       $Results \leftarrow Results \cup \{H, \hat{\mathbf{w}}\}$ ;
28      break;
29    end
30  end
31 end
32  $H_{final}, \hat{\mathbf{w}}_{final} \leftarrow$  select best from  $Results$  based on smallest  $RSS$ ;
33 return  $\hat{\mathbf{w}}_{final}$ ,  $H_{final}$ ;

```

time complexity

2.5 Combined algorithm

Here we'll describe what we've indicated above and that is combination of

todo this section
and move it some-
where else qwerty

both previous algorithms. Two options. z

1. Let fast-lts converge and then run FSA and let it converge
2. Let fast-lts converge make one step of FSA and try fast-LTS again and iterate between this

Second option would be faster ? Let's think about it.. etc etc..

2.6 Improved FSA

So far as we've described FSA algorithm we assumed that after each cycle of the algorithm (after each best swap) we need to recalculate inversion ($\mathbf{X}^T \mathbf{X}$) together with $\hat{\mathbf{w}}$. In this section we'll introduce different approaches described in [3] which will improve it so that we'll be able to update it instead of recalculate. Moreover these ideas will lead us to bounding condition of FSA which will improve the speed and we'll also be able to construct another algorithms based on this idea.

We've introduced additive formula 2.7. Let's now try to obtain similar formula but let's try to focus on how individual elements in our current algorithm will change not only on the swap of two rows but how those elements will be affected after adding one row and also removing one row. Moreover during this derivation we'll be able to obtain two different approaches of calculating one thing. Both are important and we'll recapitulate them after.

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

$$\tilde{\mathbf{Z}} = \begin{bmatrix} \mathbf{Z} & \mathbf{X}^T \mathbf{y} \\ \mathbf{y}^T \mathbf{X} & \mathbf{y}^T \mathbf{y} \end{bmatrix}. \quad (2.8)$$

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 column vector $\in \mathbb{R}^{1 \times p}$, $\mathbf{y}^T \mathbf{X}$ is vector $\in \mathbb{R}^{p \times 1}$ and $\mathbf{y}^T \mathbf{y}$ is scalar. OLS estimate $\hat{\mathbf{w}}$ is then

$$\hat{\mathbf{w}}^{(OLS,n)} = \mathbf{Z}^{-1} \mathbf{X}^T \mathbf{y} \quad (2.9)$$

and residual sum of squares

$$RSS = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{\mathbf{w}} \quad (2.10)$$

Also note that all necessary multiplications are already in the matrix $\tilde{\mathbf{Z}}$. Let's now realize 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

$$\begin{aligned}
 RSS &= \frac{\det(\tilde{\mathbf{Z}})}{\det(\mathbf{Z})} \\
 &= \frac{\det \begin{pmatrix} \mathbf{Z} & \mathbf{X}^T \mathbf{y} \\ \mathbf{y}^T \mathbf{X} & \mathbf{y}^T \mathbf{y} \end{pmatrix}}{\det(\mathbf{M})} \\
 &= \frac{\det(\mathbf{Z}) \det(\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{Z}^{-1} \mathbf{y})}{\det(\mathbf{Z})} \\
 &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{\mathbf{w}}.
 \end{aligned} \tag{2.11}$$

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

$$\tilde{\mathbf{Z}}^{-1} = \begin{bmatrix} \mathbf{Z}^{-1} + \frac{\hat{\mathbf{w}} \hat{\mathbf{w}}^T}{RSS} & -\frac{\hat{\mathbf{w}}^T}{RSS} \\ -\frac{\hat{\mathbf{w}}}{RSS} & \frac{1}{RSS} \end{bmatrix}. \tag{2.12}$$

For following equations it's important to notice that for any two observations $\mathbf{z}_i = (x_i, y_i)$ and $\mathbf{z}_j = (x_j, y_j)$ $\mathbf{z}_i, \mathbf{z}_j \in \mathbb{R}^{p+1 \times 1}$ is

$$\mathbf{z}_i \tilde{\mathbf{Z}}^{-1} \mathbf{z}_i^T = \frac{(y_i - \mathbf{x}_i \hat{\mathbf{w}})^2}{RSS} + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T \tag{2.13}$$

and

$$\mathbf{z}_j \tilde{\mathbf{Z}}^{-1} \mathbf{z}_j^T = \frac{(y_j - \mathbf{x}_j \hat{\mathbf{w}})(y_i - \mathbf{x}_i \hat{\mathbf{w}})}{RSS} + \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T \tag{2.14}$$

Using above let's express how the determinant $\det(\mathbf{Z})$ and inverse of $\text{vec} \mathbf{Z}^{-1}$ will change when we'll add one observation $\mathbf{z}_i = (x_i, y_i) \in \mathbb{R}^{p+1 \times 1}$ to the matrix \mathbf{A} . First lets notice that if we add this row to \mathbf{A} , then \mathbf{Z} will change

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} & x_{1i} \\ x_{21} & x_{22} & \dots & x_{2n} & x_{2i} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{p1} & x_{p2} & \dots & x_{pn} & x_{pi} \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \\ x_{i1} & x_{i2} & \dots & x_{ip} \end{bmatrix} = \mathbf{X}^T \mathbf{X} + \mathbf{x}_i^T \mathbf{x}_i = \mathbf{Z} + \mathbf{x}_i^T \mathbf{x}_i, \tag{2.15}$$

so determinant with appended row will be

$$\det(\mathbf{Z} + \mathbf{x}_i^T \mathbf{x}_i) = \det(\mathbf{Z})(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T) \tag{2.16}$$

and inversion can be obtained using Sherman-Morrison formula [12] so that

$$(\mathbf{Z} + \mathbf{x}_i^T \mathbf{x}_i)^{-1} = \mathbf{Z}^{-1} - \frac{\mathbf{Z}^{-1} \mathbf{x}_i^T \mathbf{x}_i \mathbf{Z}^{-1}}{1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T} \tag{2.17}$$

2. ALGORITHMS

For following equations it'll be convenient for us to mark

$$b = \frac{-1}{(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T)}, b \in \mathbb{R}, \quad (2.18)$$

and

$$\mathbf{u} = \mathbf{Z}^{-1} \mathbf{x}_i^T, \mathbf{u} \in \mathbb{R}^{p \times 1}, \quad (2.19)$$

so that 2.17 can be written as

$$(\mathbf{Z} + \mathbf{x}_i^T \mathbf{x}_i)^{-1} = \mathbf{Z}^{-1} + b \mathbf{u} \mathbf{u}^T. \quad (2.20)$$

Given that last piece that we're missing 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 2.15 so that updated $\hat{\mathbf{w}}$ which we'll denote as $\bar{\hat{\mathbf{w}}}$ is

$$\bar{\hat{\mathbf{w}}} = (\mathbf{Z}^{-1} + b \mathbf{u} \mathbf{u}^T)(\mathbf{X}^T \mathbf{y} + y_i \mathbf{x}_i^T). \quad (2.21)$$

This can be simplified so that we get

$$\bar{\hat{\mathbf{w}}} = \hat{\mathbf{w}} - (y_i - \mathbf{x}_i \hat{\mathbf{w}}) b \mathbf{u}. \quad (2.22)$$

Last but not least we want to express updated RSS which we denote as \overline{RSS} . This can be done easily from 2.11 and 2.16 as

$$\overline{RSS} = RSS + \frac{(y_i - \mathbf{x}_i \hat{\mathbf{w}})^2}{(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T)} \quad (2.23)$$

and again, it's convenient to mark

$$\gamma^+(z_i) = \frac{(y_i - \mathbf{x}_i \hat{\mathbf{w}})^2}{(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T)}, \quad (2.24)$$

so that

$$\overline{RSS} = RSS + \gamma^+(z_i) \quad (2.25)$$

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

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

Consider that we've already included one observation z_i in our dataset and mark $\bar{\mathbf{Z}} = \mathbf{Z} + \mathbf{x}_i \mathbf{x}_i^T$. If we exclude one observation $z_j = (\mathbf{x}_j, y_j) \in \mathbb{R}^{p+1 \times 1}$ from already updated matrix \mathbf{A} then determinant $\det(\bar{\mathbf{Z}})$ will change as

$$\det(\bar{\mathbf{Z}} - \mathbf{x}_j^T \mathbf{x}_j) = \det(\bar{\mathbf{Z}})(1 - \mathbf{x}_j \bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T) \quad (2.26)$$

and inversion will change (again according to Sherman-Morrison formula) as

$$(\bar{\mathbf{Z}} - \mathbf{x}_j^T \mathbf{x}_j)^{-1} = \bar{\mathbf{Z}}^{-1} + \frac{\bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T \mathbf{x}_j \bar{\mathbf{Z}}^{-1}}{1 - \mathbf{x}_j \bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T}. \quad (2.27)$$

mistake in original paper where is $w + (y - \mathbf{x} \mathbf{w}) b \mathbf{u}$

Once again, it's convenient to denote

$$\bar{b} = \frac{-1}{(1\mathbf{x}_j\bar{\mathbf{Z}}^{-1}\mathbf{x}_j^T)}, \bar{v} \in \mathbb{R}, \quad (2.28)$$

and

$$\bar{\mathbf{u}} = \bar{\mathbf{Z}}^{-1}\mathbf{x}_j^T, \bar{\mathbf{u}} \in \mathbb{R}^{p \times 1}, \quad (2.29)$$

so that we can write

$$(\bar{\mathbf{Z}} + \mathbf{x}_j^T\mathbf{x}_j)^{-1} = \mathbf{Z}^{-1} - \bar{b}\bar{\mathbf{u}}\bar{\mathbf{u}}^T. \quad (2.30)$$

Using same approach as before we can denote express downdated estimate which we denote as $\bar{\bar{\mathbf{w}}}$

$$\bar{\bar{\mathbf{w}}} = +\bar{\mathbf{w}}(y_j - \mathbf{x}_j\bar{\mathbf{w}})\bar{b}\bar{\mathbf{u}}. \quad (2.31)$$

Finally lets also express updated \bar{RSS} which we denote as $\bar{\bar{RSS}}$. This can be done easily from 2.11 and 2.16 and 2.27 as

$$\bar{\bar{RSS}} = \bar{RSS} - \frac{(y_j - \mathbf{x}_j\bar{\mathbf{w}})^2}{(1 - \mathbf{x}_j\bar{\mathbf{Z}}^{-1}\mathbf{x}_j^T)} \quad (2.32)$$

and let's also mark

$$\gamma^-(z_j) = \frac{(y_j - \mathbf{x}_j\bar{\mathbf{w}})^2}{(1 - \mathbf{x}_j\bar{\mathbf{Z}}^{-1}\mathbf{x}_j^T)}, \quad (2.33)$$

so that

$$\bar{\bar{RSS}} = \bar{RSS} - \gamma^-(z_j) \quad (2.34)$$

Lets now express equation for including and removing observation at once. First let's notice that from 2.27 we can express

$$\mathbf{x}_j\bar{\mathbf{Z}}^{-1}\mathbf{x}_j^T = \mathbf{x}_j\mathbf{Z}^{-1}\mathbf{x}_j^T - \frac{(\mathbf{x}_i\mathbf{Z}^{-1}\mathbf{x}_j^T)^2}{1 + \mathbf{x}_i\mathbf{Z}^{-1}\mathbf{x}_i^T}, \quad (2.35)$$

so that

$$\begin{aligned} \det(\mathbf{Z} + \mathbf{x}_i^T\mathbf{x}_i - \mathbf{x}_j^T\mathbf{x}_j) = \\ \det(\mathbf{Z})(1 + \mathbf{x}_i\mathbf{Z}^{-1}\mathbf{x}_i^T - \mathbf{x}_j\mathbf{Z}^{-1}\mathbf{x}_j^T + (\mathbf{x}_i\mathbf{Z}^{-1}\mathbf{x}_j^T)^2 - \mathbf{x}_i\mathbf{Z}^{-1}\mathbf{x}_i^T\mathbf{x}_j\mathbf{Z}^{-1}\mathbf{x}_j^T) \end{aligned} \quad (2.36)$$

Finally we can express $\bar{\bar{RSS}}$ as

$$\bar{\bar{RSS}} = RSS\rho(z_i, z_j), \quad (2.37)$$

2. ALGORITHMS

where

$$\rho(z_i, z_j) = \frac{(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T + \frac{e_j^2}{RSS})(1 - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T - \frac{e_i^2}{RSS}) + (\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T + \frac{e_i e_j}{RSS})^2}{1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T + (\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T)^2 - \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T}, \quad (2.38)$$

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

We can see that this formula is similar to the one 2.7 but here the $\rho(z_i, z_j)$ represents multiplicative increment. Moreover $0 < \rho(z_i, z_j)$ and if $0 < \rho(z_i, z_j) < 1$ then the swap leads to improvement in terms of feasible solution.

With all this we are now able to change FSA so that we don't need to recompute $\hat{\mathbf{w}}$ and inversion $(\mathbf{X}^T \mathbf{X})^{-1}$ but we can update it. That means after we find z_i, z_j that improve the current RSS the most (that means we find smallest $\rho(z_i, z_j)$ by 2.38). We can then update RSS by 2.37, update $(\mathbf{X}^T \mathbf{X})^{-1}$ by 2.20 and 2.27 and finally update $\hat{\mathbf{w}}$ by 2.22 and 2.31.

pseudokod

Let's now talk about time complexity of such solution. We can clearly see that core iteration takes the most time. So let's first describe time complexity of finding optimal swap and updating result with such pair. We need to go through all pairs between z_i, z_j where z_i is from h -subset and z_j is from g -subset. That is

$$h(n-h) \approx \left(\frac{n^2 - p^2}{4}\right)^2 \sim \mathcal{O}(n^2 - p^2) \quad (2.39)$$

For each of this pair we need to calculate 2.7. Note that e_i and e_j can be calculated before this loop. The same goes for $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T$ and $\mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T$. So inside the loop there is only one vector-matrix-vector multiplication that is $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T$ which is $\sqrt{\epsilon}$. Rest are scalars. The whole loop is then

$$\mathcal{O}((n^2 - p^2)p^2) \quad (2.40)$$

Before we'll move to the updating we must not to forget for quantities which we said we can calculate in advance. That is

$\mathcal{O}(2np^2)$ for all $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T$ and $\mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T$ and $\mathcal{O}(2np)$ for all e_i and e_j .

Finally we need to update RSS which is trivial. Next inversion by 2.20 which is $\mathcal{O}(4p^2)$ and 2.27 which is also $\mathcal{O}(4p^2)$. Finally we need to update $\hat{\mathbf{w}}$ by 2.22 and 2.31 which are both $\mathcal{O}(p^2 + p)$. Putting everything together we get

$$\mathcal{O}((n^2 - p^2)p^2 + 2np^2 + 2np + 8p^2 + 2p^2 + p) \sim \mathcal{O}(n^2 p^2 - p^4) \quad (2.41)$$

nejak zminit ko-
likrat ten loop
vetsinou projede
atd..

Note that right now it doesn't matter if we use 2.7 with stopping criterion $\Delta S_{i,j}^{(m)} \geq 0$ or 2.38 and use $\rho(z_i, z_j) \geq 1$ stopping criterion. With both results we can update RSS . But the advantage of 2.38 is that we can use following bounding condition to improve performance.

Bounding condition improvement

$\rho(z_i, z_j)$ is expressed as ratio. We can see that in numerator we have

$$(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T + \frac{e_j^2}{RSS})(1 - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T - \frac{e_i^2}{RSS}) + (\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T + \frac{e_i e_j}{RSS})^2 \quad (2.42)$$

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

$$(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T + \frac{e_j^2}{RSS})(1 - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T - \frac{e_i^2}{RSS}). \quad (2.43)$$

On the other hand we can see that denominator is

$$1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T + (\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T)^2 - \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T \quad (2.44)$$

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) so

$$(\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T)^2 \leq \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T \quad (2.45)$$

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

$$1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T. \quad (2.46)$$

Given that we can denote $\rho_b(z_i, z_j)$ so that

$$\rho_b(z_i, z_j) = \frac{(1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T + \frac{e_j^2}{RSS})(1 - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T - \frac{e_i^2}{RSS})}{1 + \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T - \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T} \leq \rho(z_i, z_j) \quad (2.47)$$

So the actual speed improvement is given by we don't need to compute $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T$ in each of $h(n-h)$ pairs swap comparison. Every time that algorithm starts to compare all pairs we set $\rho_{min} := 1$. Then every time we only compute $\rho_b(z_i, z_j)$ and if it is greater or equal to ρ_{min} we can continue to next pair without computing $\rho(z_i, z_j)$. On the other hand if $\rho_b(z_i, z_j)$ is less than ρ_{min} we compute $\rho(z_i, z_j)$ and set $\rho_{min} := \rho(z_i, z_j)$. This of course won't improve \setminus^ϵ time required to compare all z_i, z_j pairs. Finally let's note that [3] call FSA which compute $\rho(z_i, z_j)$ 2.38 as *Optimal exchange algorithm (OEA)* and OEA using bounding condition 2.47 as *Modified optimal exchange algorithm (MOEA)*.

pseudokod - tak
moc nechci :-D

2.6.0.1 Different method of computing Improved FSA

We've described way how to modify FSA so that we can update *hatw* and inversion $(\mathbf{X}^T \mathbf{X})^{-1}$. This requires to compute inversion at the start of the algorithm (this is also the case for the FSA). In practice, however, this is not the way how OLS estimate $\hat{\mathbf{w}}$ is computed. That is primarily due to the fact that computing inversion of large matrix is computationally exhaustive and first of all not numerically stable. In practice we usually use Cholesky factorization. In this subsection we'll describe it and also show that we can perform all above computations using such decomposition. Moreover we'll show close connection between Cholesky factorization and QR decomposition.

Let's start with Cholesky factorization.

Definition 20. 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

$$\mathbf{A} = \mathbf{L}\mathbf{L}^T \quad (2.48)$$

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

If we look at our problem of finding solution to

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (2.49)$$

We can easily rewrite it as

$$\mathbf{Z} \mathbf{w} = \mathbf{b} \quad (2.50)$$

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

$$\mathbf{L} \mathbf{d} = \mathbf{b} \quad (2.51)$$

$$\mathbf{L}^T \mathbf{w} = \mathbf{d}. \quad (2.52)$$

where \mathbf{d} is solved by forward substitution and \mathbf{w} is then solved by backward substitution which is our $\hat{\mathbf{w}}$ estimate.

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

1. Calculate $\mathbf{X}^T \mathbf{X}$ and $\mathbf{X}^T \mathbf{y}$
2. Calculate 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.

zmenit to oznaceni
Improved FSA na
neco jineho. mys-
lim ze to pouzivam
nekonzistentne

presunout cholesky
popis do 1. kapi-
toly

4. Solve upper triangular system $\mathbf{L}^T \mathbf{w} = \mathbf{d}$ for \mathbf{w}

Observation 21. Time complexity of solving OLS using Cholesky factorization is $\mathcal{O}(np^2)$

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 steps is calculating Cholesky factorization. This can be done in $\mathcal{O}(\frac{1}{3}p^3)$ [13]

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

$$\mathcal{O}(np^2 + np + \frac{1}{3}p^3 + p^2) \quad (2.53)$$

and because we assume $n \geq p$ and most usually $n \gg p$ then multiplication $\mathbf{X}^T \mathbf{X}$ asymptotically dominates so whole time complexity is $\mathcal{O}(np^2)$ \square

Note 22. In case we have $\mathbf{X}^T \mathbf{X}$ and $\mathbf{X}^T \mathbf{y}$ computed in advance, then it is asymptotically dominated by Cholesky factorization

$$\mathcal{O}(\frac{1}{3}p^3) \sim \mathcal{O}(p^3) \quad (2.54)$$

Let's now look on similar method of computing OLS which don't require multiplying $\mathbf{X}^T \mathbf{X}$. This can be done using QR decomposition.

Definition 23. 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

$$\mathbf{A} = \mathbf{Q}\mathbf{R} \quad (2.55)$$

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. We call this decomposition as *QR decomposition*.

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

$$\mathbf{A} = \mathbf{Q}\mathbf{R} = \mathbf{Q} \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{Q}_2 \end{bmatrix} \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix} = \mathbf{Q}_1 \mathbf{R}_1 \quad (2.56)$$

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, Where $\mathbf{Q}_1 \in \mathbb{R}^{n \times n}$ is matrix with orthogonal columns and $\mathbf{Q}_2 \in \mathbb{R}^{(m-n) \times n}$ is also matrix with orthogonal columns.

So that means that \mathbf{X} can be factorized so that

$$\mathbf{X} = \mathbf{Q}\mathbf{R}. \quad (2.57)$$

Then

$$\mathbf{X}^T \mathbf{X} = (\mathbf{QR})^T \mathbf{QR} = \mathbf{R}^T \mathbf{Q}^T \mathbf{QR} \quad (2.58)$$

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

$$\mathbf{X}^T \mathbf{X} = \mathbf{R}^T \mathbf{R}. \quad (2.59)$$

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

$$\mathbf{X}^T \mathbf{X} = \mathbf{R}_1^T \mathbf{R}_1. \quad (2.60)$$

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

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (2.61)$$

is then given by

$$\mathbf{R}_1^T \mathbf{R}_1 \mathbf{w} = \mathbf{R}_1^T \mathbf{Q}_1^T \mathbf{y} \quad (2.62)$$

and because we assume \mathbf{X} have full column rank, thus is invertible, we can simplify it to

$$\mathbf{R}_1 \mathbf{w} = \mathbf{b}_1 \quad (2.63)$$

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.

Note 24. 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

$$\mathbf{R}_1^T = \mathbf{L} \quad (2.64)$$

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

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

QR factorization can be calculated multiple ways. Most basic method is applying 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. Time complexity of both algorithms is similar. QR decomposition of matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ is

$$\mathcal{O}(2np^2 - \frac{2}{3}p^3). \quad (2.65)$$

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

$$\mathcal{O}(3np^2 - p^3) \quad (2.66)$$

so Householder transformations are about 50% faster but Givens rotations are more numerically stable and are suitable for sparse matrices. That will be suitable for our problem. That's why we'll describe Givens rotation in detail in 2.6.0.1. We'll also make a proof of 2.66 time complexity there.

For now let's look at time complexity of solving OLS using Householder transformations.

Observation 25. Time complexity of solving OLS using QR decomposition using Householder transformations is $\mathcal{O}(2np^2 - \frac{2}{3}p^3)$

Proof. In the first step we calculate QR decomposition. Householder transformations can be done in $\mathcal{O}(2np^2 - \frac{2}{3}p^3)$ which is slightly faster than Givens rotation $\mathcal{O}(3np^2 - p^3)$. 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

$$\mathcal{O}(2np^2 - \frac{2}{3}p^3 + np + \frac{1}{2}p^2) \quad (2.67)$$

We can see that again $\mathcal{O}(2np^2)$ asymptotically dominates. \square

Givens Rotation

In this section we'll describe Givens rotations in detail. Moreover we'll also show how to update QR decomposition in terms of including and excluding row from factorized matrix \mathbf{X} which will help us in our algorithm.

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 matrix \mathbf{Q}^T so that [14]

$$\mathbf{Q}^T \mathbf{A} = \mathbf{R} \quad (2.68)$$

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

$$\mathbf{Q} = \begin{bmatrix} \cos(\varphi) & \sin(\varphi) \\ -\sin(\varphi) & \cos(\varphi) \end{bmatrix} \quad (2.69)$$

This matrix is indeed orthogonal since

$$\begin{aligned} \mathbf{Q}\mathbf{Q}^T &= \mathbf{Q}^T\mathbf{Q} \\ &= \begin{bmatrix} \cos(\varphi)^2 + \sin(\varphi)^2 & \sin(\varphi)\cos(\varphi) - \cos(\varphi)\sin(\varphi) \\ \cos(\varphi)\sin(\varphi) - \sin(\varphi)\cos(\varphi) & \cos(\varphi)^2 + \sin(\varphi)^2 \end{bmatrix} = \mathbf{I} \end{aligned} \quad (2.70)$$

2. ALGORITHMS

Moreover if we multiply this orthogonal with some column vector $\mathbf{x} \in \mathbb{R}^{2 \times 1}$ thus $\mathbf{Q}\mathbf{x}$ as a result we'll get vector of same length but is rotated clockwise by φ radians. When we say that vector will have same length we mean that L-2 norm of such vector will be the same. This is important property of all orthogonal matrices. We can simply verify correctness of this claim. Let's have any orthogonal matrix $\mathbf{Q} \in \mathbb{R}^{m \times m}$ and any column vector $\mathbf{x} \in \mathbb{R}^{m \times 1}$ then

$$\|\mathbf{Q}\mathbf{x}\|^2 = (\mathbf{Q}\mathbf{x})^T(\mathbf{Q}\mathbf{x}) = \mathbf{x}^T \mathbf{Q}^T \mathbf{Q} \mathbf{x} = \mathbf{x}^T \mathbf{I} \mathbf{x} = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|^2 \quad (2.71)$$

So we would like to create series of orthogonal matrices which will gradually rotate column vectors of \mathbf{A} so there will be zeros under diagonal. One such method - *Givens rotation* uses *Givens matrices* that will zero one element under diagonal at a time. The example of 2.69 was not random. Let's look at this orthogonal matrix one more time and let's multiply this matrix with some vector.

$$\begin{bmatrix} \cos(\varphi) & \sin(\varphi) \\ -\sin(\varphi) & \cos(\varphi) \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} r \\ z \end{bmatrix} \quad (2.72)$$

To introduce this zeroing effect we need to rotate this vector so that it will be parallel to $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$. Now it's only up to find $\cos(\varphi)a + \sin(\varphi)b = r$. As we know rotation with \mathbf{Q} preserve L-2 norm. That means is we want to $z = 0$ then r must be equal to L-2 norm of vector $\begin{bmatrix} a \\ b \end{bmatrix}$ thus $r = \sqrt{a^2 + b^2}$. Then the solution seems trivial. We don't even need to calculate φ because if we put

$$\cos(\varphi) = \frac{a}{\sqrt{a^2 + b^2}} \quad (2.73)$$

and

$$\sin(\varphi) = \frac{b}{\sqrt{a^2 + b^2}} \quad (2.74)$$

then

$$r = \frac{a^2}{\sqrt{a^2 + b^2}} + \frac{b^2}{\sqrt{a^2 + b^2}} = \sqrt{a^2 + b^2} \quad (2.75)$$

$$z = \frac{-ab}{\sqrt{a^2 + b^2}} + \frac{ab}{\sqrt{a^2 + b^2}} = 0. \quad (2.76)$$

Practically used algorithm of computing $\cos \varphi$ and $\sin \varphi$ is slightly different because we want to prevent overflow. Real algorithm is described by following pseudocode:

Algorithm 5: Rotate

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

```

1  $\sin \leftarrow 0$ ;
2  $\cos \leftarrow 1$ ;
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    $\cotg \leftarrow \frac{a}{b}$ ;
8    $\sin \leftarrow \frac{1}{\sqrt{1+(\cotg)^2}}$ ;
9    $\cos \leftarrow \sin \cotg$ ;
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's denote matrix $Q(i, j) \in \mathbb{R}^{m \times m}$ defined as

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

where $c = \cos(\varphi)$ and $s = \sin(\varphi)$ for some φ . That means this matrix is orthogonal. Now if we have some column vector $\mathbf{x} \in \mathbb{R}^{m \times 1}$ and calculate c and s by means of 2.73 and 2.74 for $a := x_i$ and $b := x_j$. Finally if we multiply $Q(i, j)\mathbf{x} = \mathbf{p}$ we can see that \mathbf{p} is same as vector \mathbf{x} except p_i and p_j so that:

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

where $r = \sqrt{x_i^2 + x_j^2}$. If we have matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$ instead of only one column \mathbf{x} it works the same way. 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 exactly $e = \frac{p^2 - p}{2} + np - p^2$ 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

$$\mathbf{Q}_e \dots \mathbf{Q}_2 \mathbf{Q}_1 \mathbf{A} = \mathbf{R} \quad (2.77)$$

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

$$\mathbf{Q} = \mathbf{Q}_1^T \mathbf{Q}_2^T \dots \mathbf{Q}_e^T. \quad (2.78)$$

We can see that for each row of matrix we need one additional matrix \mathbf{Q}_i and with such matrix we are actually multiplying only two rows. Moreover the rows are getting shorter after each finished column. So time we see that time complexity is equal to

$$\sum_{i=1}^n \sum_{j=i+1}^p 6(n - i + 1) \approx 6np^2 - 3np^2 - 3p^3 + 2p^3 = 3np^2 - p^3 \quad (2.79)$$

todo pseudokod?
moc se mi nechce
:-D

move all this
about rotations
somewhere else

So we are now in situation when we have QR decomposition of \mathbf{X} and we need to exchange i th row for j th row. We can simulate this by first adding j th row and consequently removing i th row.

Updating QR decomposition

First let's discuss how to update QR decomposition when row is added. This is indeed very simple. If we add row to \mathbf{A} our decomposition looks like

$$\mathbf{R}^+ = \mathbf{Q}_e \dots \mathbf{Q}_2 \mathbf{Q}_1 \mathbf{A}_{new} = \begin{bmatrix} \times & \dots & \dots & \dots & \times \\ 0 & \times & \dots & \dots & \times \\ \vdots & 0 & \ddots & \dots & \vdots \\ \vdots & \vdots & 0 & \ddots & \vdots \\ \vdots & \vdots & \vdots & 0 & \times \\ \vdots & \vdots & \vdots & \vdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \times & \dots & \dots & \dots & \times \end{bmatrix} \quad (2.80)$$

So all we need to do is create matrices $\mathbf{Q}_{e+1} \dots \mathbf{Q}_{e+p}$ to zero newly added row. So the \mathbf{R} then will be equal to $\mathbf{R}_{new} = \mathbf{R} \mathbf{Q}_{e+1} \mathbf{Q}_e \dots \mathbf{Q}_2 \mathbf{Q}_1 \mathbf{A}$. \mathbf{Q} will be updated in the same way thus $\mathbf{Q}_{new} = \mathbf{Q} \mathbf{Q}_e^T \mathbf{Q}_{e+1}^T \dots \mathbf{Q}_{e+p}^T$.

blbe znaceni ale
nevim jak lepe

Algorithm 6: QR insert

Input: $\mathbf{Q}, \mathbf{R}, x_i$

Output: $\mathbf{Q}^+, \mathbf{R}^+$

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

todo pseudokod?
moc se mi nechce
:-D

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

Note 26. For adding rows matrix \mathbf{Q} is actually not needed so that we can update \mathbf{R} to \mathbf{R}^+ . Later we'll show that this is actually useful.

When we are deleting the row x_i from matrix \mathbf{A} we can use following trick [14]. First we move such row as a first row of matrix \mathbf{A} so we create permutation matrix \mathbf{P} so that

$$\mathbf{P} \mathbf{A} = \begin{bmatrix} x_i \\ \mathbf{A}(1:i-1, 1:p) \\ \mathbf{A}(1:i+1, 1:p) \end{bmatrix} = \begin{bmatrix} x_i \\ \mathbf{A}^- \end{bmatrix} = \mathbf{P} \mathbf{Q} \mathbf{R} \quad (2.81)$$

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

$$\mathbf{Q}(n-1, n) \dots \mathbf{Q}(1, 2) \mathbf{q}_1^T = \begin{bmatrix} 1 \\ 0 \\ \vdots \end{bmatrix} \quad (2.82)$$

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

$$\mathbf{Q}(n-1, n) \dots \mathbf{Q}(1, 2) \mathbf{q}_1^T \mathbf{R} = \begin{bmatrix} \times \\ \mathbf{R}^- \\ \vdots \end{bmatrix}. \quad (2.83)$$

Thus result is

$$\mathbf{P}\mathbf{A} = (\mathbf{P}\mathbf{Q}\mathbf{Q}^T(1, 2) \dots \mathbf{Q}^T(n-1, n))(\mathbf{Q}(n-1, n) \dots \mathbf{Q}(1, 2) \mathbf{q}_1^T \mathbf{R}) = \begin{bmatrix} 1 & 0 \\ 0 & \mathbf{Q}^- \end{bmatrix} \begin{bmatrix} \times \\ \mathbf{R}^- \end{bmatrix}, \quad (2.84)$$

so that

$$\mathbf{A}^- = \mathbf{Q}^- \mathbf{R}^- \quad (2.85)$$

Algorithm 7: QR delete

Input: $\mathbf{Q}, \mathbf{R}, i$

Output: $\mathbf{Q}^-, \mathbf{R}^-$

1 return $\mathbf{Q}^-, \mathbf{R}^-$;

Computing \mathbf{R}^- is $\mathcal{O}(p^2)$ and computing \mathbf{Q}^- is $\mathcal{O}(np)$.

2.6.0.2 Calculation of improved FSA using QR

So let's focus on our problem. As we said using decomposition is better than calculating inversion primarily due to higher numerical stability. We'll show that both FSA as well as improved FSA can be calculated using QR decomposition.

Let's start with 2.38. Here inversion is need to calculate $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T, \mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T$ and $\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T$. This can also be done without inversion, only using the decomposition. We can write this equation

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

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

$$\mathbf{R}^T \mathbf{v} = \mathbf{x}_i^T. \quad (2.87)$$

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

$$\mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_j^T = \mathbf{x}_j \mathbf{u} \iff \mathbf{x}_i (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{x}_j^T = \mathbf{x}_j^T \mathbf{u} \quad (2.88)$$

Where column vector \mathbf{u} is defined by 2.19. We can see that

$$\mathbf{u} = \mathbf{Z}^{-1} \mathbf{x}_i^T \iff (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{x}_i^T = \mathbf{u} \quad (2.89)$$

napsat pseudokod.
opet se mi moc
necche :-D

nevim jestli zminovat ze muzu delat i SVD dekompozici ktera je nejstabilnejši ze vseh...

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

$$\mathbf{R}\mathbf{u} = \mathbf{v} \quad (2.90)$$

where \mathbf{v} is solution of 2.87. Other quantities of 2.38 does not require \mathbf{Z}^{-1} .

When optimal exchange $\mathbf{z}_i, \mathbf{z}_j$ is found then we need to update RSS which can be done same way by 2.37. Updating $\hat{\mathbf{w}}$ to $\bar{\hat{\mathbf{w}}}$ by 2.22 requires \mathbf{u} which in this case we calculate by 2.90 and b which requires $\mathbf{x}_j \mathbf{Z}^{-1} \mathbf{x}_j^T$ but that we've also already covered by idea 2.86. Analogous operations can be used to update $\bar{\hat{\mathbf{w}}}$ to $\bar{\bar{\hat{\mathbf{w}}}}$ by means of 2.31. Only thing we need to realize that we can express 2.35 as $\mathbf{x}_j \bar{\mathbf{Z}}^{-1} \mathbf{x}_j^T = \mathbf{x}_i \mathbf{Z}^{-1} \mathbf{x}_i^T + b(\mathbf{u} \mathbf{x}_j^2)$.

We can't use updating inversion because we don't have one so we need to somehow update our QR decomposition. This can be done by both Householder rotation as well as by Givens rotation, but because our \mathbf{R} is already sparse matrix, Givens rotations are ideal tool for this. We've already described how this can be done. First we can use 6 to add row \mathbf{x}_i to the decomposition and consequently 7 to remove \mathbf{x}_j from the decomposition. We can see that this is slower than updating inversion directly. On the other hand this solution is numerically stable and requires less time than computing factorization again from scratch.

Finally let's talk about matrix $\tilde{\mathbf{Z}}$ which we used 2.8 for derivation of our equations. If we use 2.64 observation then we realize that if we make QR factorization of $\mathbf{A} = (\mathbf{X}, \mathbf{y})$ then

$$\tilde{\mathbf{Z}} = \begin{bmatrix} \mathbf{Z} & \mathbf{X}^T \mathbf{y} \\ \mathbf{y}^T \mathbf{X} & \mathbf{y}^T \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{R}^T & 0 \\ \phi^T & r \end{bmatrix} \begin{bmatrix} \mathbf{R} & \phi \\ 0 & r \end{bmatrix} \quad (2.91)$$

where

$$\tilde{\mathbf{R}} = \begin{bmatrix} \mathbf{R} & \phi \\ 0 & r \end{bmatrix} = \tilde{\mathbf{Q}}^T \mathbf{A} \quad (2.92)$$

is matrix from QR factorization of \mathbf{A} . Next we can realize that $\mathbf{R} \in \mathbb{R}^{p \times p}$ is actually 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

$$\phi = \mathbf{Q}^T \mathbf{y} \quad (2.93)$$

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

$$\mathbf{R}\hat{\mathbf{w}} = \phi \quad (2.94)$$

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

$$r^2 = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{\mathbf{w}} = RSS \quad (2.95)$$

Remark 27. We can see that all quantities we use in the algorithm can be obtained from matrix $\tilde{\mathbf{R}}$ moreover if we didn't need to remove rows from \mathbf{X} then matrix $\tilde{\mathbf{Q}}$ would not need 26, thus speedup while inserting and rows would be possible.

Now we've described all properties of improved FSA algorithm. We can see that two different methods of computation are possible. One, using inversion \mathbf{Z}^{-1} is faster, but numerically less stable, second using QR decomposition is slower due to updating matrix \mathbf{R} and matrix \mathbf{Q} . We've also shown that using matrix $\tilde{\mathbf{R}}$ is useful because it contains all necessary quantities for the improved FSA algorithm, although it won't improve the speed of the algorithm. Later we'll introduce different algorithm where we won't need to remove rows from \mathbf{X} so then we'll be able to use this observation in our favor.

time complexity
QR qwertý

time complexity
tildeZQR is the
same as QR

no taak. to je na
stranku. to das

zminit podobnost
s: Hoffmann, Kon-
toghiorghes, 2010
(Matrix strategies
SUR)

2.6.1 Min-Max-FSA aka. Min Max exchange (MMEA)

2.7 Branch and bound aka. BAB

2.8 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

$OF(ols, x, y)(w)$ $OF^{(OLS, X, y)}(w)$
 $\hat{w} (OLS, n)$ $\hat{w}^{(OLS, n)}$
 $OF(lts, x, y)(w)$ $\hat{w}^{(LTS, X, y)}(w)$ NOT USED
 $OF(lts, J(w))$ $J(m)$
 $\hat{w} (LTS, h, n)$ $\hat{w}^{(LTS, h, n)}$
 $VEC\ m$ $m_{LTS} \in Q^{(n, h)}$ $m \in Q^{(n, h)}$ $O\ Z$
 $M = diag(m)$ $MX\ MY$ $OF^{(OLS, MX, MY)}$
 h subset g subset
 $J\ J$
 $RSS\ RSS = \sum_{i=1}^n r_i^2 = (y_i - w^T x_i)^2$

Conclusion

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] Bai, E.-W. A random least-trimmed-squares identification algorithm. *Automatica*, volume 39, no. 9, 2003: pp. 1651–1659.
- [7] Rousseeuw, P. J.; Leroy, A. M. *Robust regression and outlier detection*. John wiley & sons, 1987.
- [8] 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.

- [9] 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.
- [10] Hawkins, D. M. The feasible solution algorithm for least trimmed squares regression. *Computational statistics & data analysis*, volume 17, no. 2, 1994: pp. 185–196.
- [11] 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.
- [12] Bartlett, M. S. An inverse matrix adjustment arising in discriminant analysis. *The Annals of Mathematical Statistics*, volume 22, no. 1, 1951: pp. 107–111.
- [13] Krishnamoorthy, A.; Menon, D. Matrix inversion using Cholesky decomposition. In *2013 Signal Processing: Algorithms, Architectures, Arrangements, and Applications (SPA)*, IEEE, 2013, pp. 70–72.
- [14] Hammarling, S.; Lucas, C. Updating the QR factorization and the least squares problem. 2008.

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