

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.

March 16, 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 March 16, 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	7
1.1 Linear regression model	7
1.2 Ordinary least squares	8
1.3 Robust statistics	10
1.4 Least trimmed squares	12
2 The Least trimmed squares	13
3 Algorithms	15
3.1 Computing LTS	15
3.2 FAST-LTS	15
3.3 Exact algorithm	23
3.4 Feasible solution	23
3.5 MMEA	23
3.6 Branch and bound	23
3.7 Adding row	23
4 Experiments	25
4.1 Data	25
4.2 Results	25
4.3 Outlier detection	25
Conclusion	27
Bibliography	29
A Datasets	31

B Contents of enclosed CD

33

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

Todo list

proof it is a local minimum	9
ols estimate is BLUE if we assume this and that	10

Least trimmed squares

In this chapter will be introduced one of the most common regression analysis models which is linear regression model as well as ordinary least squares method for estimating is variables. This model tries to model relationship between one variable which is considered dependent and one or more variables which are considered to be explanatory. Relationship is based on model function with parameters which are not known in advance and are estimated from data. We will also introduce one of the most commons methods of finding those parameters of this model, namely ordinary least squares.

1.1 Linear regression model

Definition 1. Linear regression model is

$$y_i = \mathbf{x}_i^T \mathbf{w} + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (1.1)$$

where $y_i \in \mathbb{R}$ is random variable which we denote as *dependent variable*, vector $\mathbf{x}_i^T = (x_1, x_2, \dots, x_p)$ is column vector of *explanatory variables* and $\varepsilon \in \mathbb{R}$ random variable called *noise* or *error*. Vector $\mathbf{w} = (w_1, w_2, \dots, w_p)$ is vector of parameters called *regression coefficients*. It is common to write whole model in matrix form.

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

where

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

We consider $\mathbf{x}_{i1} = 1$ to be a constant. Corresponding w_1 we call *intercept*. We refer this model as *model with intercept*. It is possible to have model without intercept but then hyperplane generated by this model goes through origin. That is usually rare case so we'll always use intercept unless we'll explicitly mention it. Due to this we can assume that in model with intercept expected value $\mathbb{E}(\varepsilon_i) = 0$.

1.1.1 Prediction with linear regression model

Linear regression model contains vector of real regression coefficients which we don't know and which we need to estimate in order to be able to use model for prediction. So let us assume that we already have estimated regression coefficients which we'll mark $\hat{\mathbf{w}}$. Then we're able to predict value of y by

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

\hat{y} denotes that it is predicted value. Real value of y is given by

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

Because we assume linear dependence between dependent variable y and explanatory variables \mathbf{x} then what makes y random variable is actually random variable ε . Because we use model with intercept thus with $\mathbb{E}(e) = 0$ we can see that

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

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

1.2 Ordinary least squares

We want to estimate \mathbf{w} so that error of the model will be minimal. Measurement of this error is most often done by *loss function*

$$L : \mathbb{R}^2 \rightarrow \mathbb{R} \quad (1.6)$$

Which in case of ordinary least squares is quadratic loss function $L(y, \hat{y}) := (y - \hat{y})^2$. We refer to $r(\hat{\mathbf{w}}) = y - \hat{y} = y - \hat{\mathbf{w}}^T \mathbf{x}$ as to residual.

So the idea of this method lies in fact that we want to minimize error given by sum of squared residuals commonly known as residual sum of squares *RSS*

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

So if we think about RSS as the function of $\hat{\mathbf{w}}$ we can mark

$$\hat{\mathbf{w}} = \arg \min_{\tilde{\mathbf{w}} \in \mathbb{R}^p} RSS(\tilde{\mathbf{w}}) = \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (y_i - \tilde{\mathbf{w}}^T \mathbf{x}_i)^2 \quad (1.8)$$

which is basically a definition of ordinary least squares estimator.

Definition 2. Estimate of regression parameters using ordinary least squares (OLS) on n observations is defined as

$$\hat{\mathbf{w}}^{OLS,n} = \sum_{i=1}^n (y_i - \tilde{\mathbf{w}}^T \mathbf{x}_i)^2 \quad (1.9)$$

Let's now talk about finding solution of this. If we want to find minimum of this function, first we need to find gradient by calculating all partial derivatives

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

By this we obtain gradient

$$\nabla RSS = - \sum_{i=1}^n 2(y_i - \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i. \quad (1.11)$$

As a *normal equation* we'll mark gradient equal to zero.

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

Which we can write in matrix form as

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

Let's now construct hessian matrix using second-order partial derivatives.

With

$$\frac{\partial^2 RSS}{\partial w_h \partial w_j} = \sum_{i=1}^n 2(-x_{ih})(-x_{ij}), \quad h \in \{1, 2, \dots, p\}. \quad (1.14)$$

we get

$$\mathbf{H}_{RSS} = 2\mathbf{X}^T \mathbf{X}. \quad (1.15)$$

We can see that hessian \mathbf{H}_{RSS} is always positive semi-definite because for all $\mathbf{s} \in \mathbb{R}^p$

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

It's easy to proof that twice differentiable function is convex if and only if the hessian of such function is positive semi-definite. Due to that we can say that solution of (1.13) gives us not only local minimum but global minimum.

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

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

proof it is a local minimum

1. LEAST TRIMMED SQUARES

where (OLS, n) denotes that $\hat{\mathbf{W}}$ is estimate of \mathbf{w} calculated using ordinary least squares on n observations.

Moreover we can see that if $\mathbf{X}^T \mathbf{X}$ is regular matrix, then a hessian \mathbf{H}_{RSS} is positive definite because for all nonzero $\mathbf{s} \in \mathbb{R}^p$

$$\mathbf{s}^T 2\mathbf{X}^T \mathbf{X} \mathbf{s} > 0 \quad (1.18)$$

thus $\hat{\mathbf{w}}^{(OLS, n)}$ in (1.17) is strict global minimum.

ols estimate is
BLUE if we as-
sume this and that

1.3 Robust statistics

Standard statistic methods to work properly expects some assumptions and fail if those assumptions are not met. Robust statistic are statistics that produce acceptable results even if 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 real world we are often forced to work with such data. Classical statistics methods provide poor results with such data.

Ordinary least squares method has such assumptions about data. It expects that

$$\varepsilon_i \sim \mathcal{N}(\mu, \sigma^2), \quad i \in \{1, 2, \dots, n\}. \quad (1.19)$$

Condition of expected value to be zero is met, if we use model with intercept. But we can't be sure if errors are normally distributed. Moreover we expect that that variance stays the same for all $\varepsilon_i, i \in \{1, 2, \dots, n\}$. That's what we call *homoscedasticity*. Beside assumptions about ε we expect that explanatory variables x_1, x_2, \dots, x_p are not linear dependent, or in terms of statistics, uncorrelated.

Before we'll explain what happens if those conditions are not met or are met only partially let's describe one of the most common reason are assumptions are false.

1.3.1 Outliers

We stated a lot of assumptions that are required for ordinary least squares to give good estimate of $\hat{\mathbf{w}}$. Unfortunately in real conditions these assumptions are often false so that ordinary least squares won't produce acceptable results. One of the most common reasons of false assumptions are abnormal observations called outliers.

Outliers are very common and they are for instance erroneous measurements such as transmission errors or noise. Other very common reason is that nowadays we are mostly given data which were automatically processed by computers. Sometimes we are also presented data which are heterogenous in the sense that they contain data from multiple regression models. In certain

sense outliers are inevitable. One would say that we would 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 simply too big to check and even more often we don't know how the data should look. In higher dimensions it is very difficult to find outliers. There are some methods that tries 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 presence of outliers but even help to identify data which seems to be outliers.

We have some terminology to describe certain types of outliers. We'll use terminology from [1]. Let's 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 liner 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 couple of tools to measure robustness of statistics. Most popular one is called *breakdown point*. Then there are *empirical influence function* and *influence function* and *sensitivity curve*. For sake of simplicity we'll talk only about breakdown point right now.

Definition 3. Let T be a functional, $\mathbf{x} = (x_1, x_2, \dots, x_n), x_i \in \mathcal{X}$ be an n -dimensional random sample and $T_n(\mathbf{x})$ value of this functional with parameter \mathbf{x} . The breakdown point of T at sample \mathbf{x} we can define using sample $\mathbf{x}_{new}^{(k)}$ where we exchange k points from original sample \mathbf{x} with random values x_i . we get $T_n(\mathbf{x}_{new}^{(k)})$. Then the *breakdown point* is

$$\text{bdpoint}(T, \mathbf{x}_n) = \frac{1}{n} \min S_{T, \mathbf{x}_n} \quad (1.20)$$

where

$$S_{T, \mathbf{x}_n, D} = \{k \in \{1, 2, \dots, n\} : \sup_{\mathbf{x}_{new}^{(k)}} \|T_n(\mathbf{x}), T_n(\mathbf{x}_{new}^{(k)})\| = \infty\} \quad (1.21)$$

This definition is very general but if we specify it to our linear regression problem it basically says that breakdown point is proportion of minimal number of observations needed to be switched for some others say some really bad ones so then the estimator will give incorrect results. More robust estimators have higher breakdown point.

It's intuitive that breakdown point cannot be higher than 0.5 [2] because if we exchange more than 50% of the data, we wouldn't be able to distinguish

between the original and exchanged data. Moreover the original data would become minority over exchanged.

In context of ordinary least squares estimator (OLS) it's easy to show that one outlier is enough to increase value of T to any desired value thus

$$\text{bdpoint}(OLS, \mathbf{x}_n) = \frac{1}{n} . \quad (1.22)$$

For 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

Least trimmed squares (LTS) estimator builds on OLS but is more robust. In this section we'll define LTS estimator and show that it's breakdown point is variable and can go up to maximum possible value of breakdown point, thus 0.5.

Definition 4. Estimate of regression parameters using least trimmed squares (LTS) on h observations out of n is defined as

$$\hat{\mathbf{w}}^{LTS, h, n} = \sum_{i=1}^h r_{i:n}^2(\mathbf{w}) \quad (1.23)$$

where h , $n/2 \leq h \leq n$ is number of used observations and $r_{1:n}(\mathbf{w}), r_{2:n}(\mathbf{w}), \dots, r_{h:n}(\mathbf{w})$ are h smallest residuals.

Even though that estimation using LTS seems similar to the estimation using OLS the computation is much more complex because smallest residuals are dependent on \mathbf{w} . We'll show that this fact makes from finding LTS estimate non-convex optimization problem and finding global minimum is NP-hard.

The Least trimmed squares

2.0.1 Objective function

2.0.1.1 Problems

Algorithms

In this chapter we'll introduce currently popular algorithms of computing FAST-LTS estimate.

3.1 Computing LTS

Note 5. 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 .

3.2 FAST-LTS

In this section we will introduce FAST-LTS algorithm[3]. 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.

3.2.1 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 6. Consider dataset consisting of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ explanatory variables where $\mathbf{x}_i \in \mathbb{R}^p, \forall \mathbf{x}_i = (fx_1^i, x_2^i, \dots, x_p^i)$ where $x_1^i = 1$ and its corre-

3. ALGORITHMS

sponding 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

$$RSS(\hat{\mathbf{w}}_1) \leq RSS(\hat{\mathbf{w}}_0) \quad (3.1)$$

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 7. 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** $\hat{\mathbf{w}}_{new}, H_{new}$;
-

Observation 8. Time complexity of algorithm C-step 1 is the same as time complexity as OLS. Thus $O(p^2 n)$ **TODO**

Lemma 9. Time complexity of OLS on $\mathbf{X}^{n \times p}$ and $\mathbf{Y}^{n \times 1}$ is $O(p^2 n)$.

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 $\mathbf{A} = \mathbf{X}^T \mathbf{X} \sim \mathcal{O}(p^2 n)$ 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 $\mathcal{O}(p^2 n + pn + p^3 + p^2)$. Because $\mathcal{O}(p^2 n)$ and $\mathcal{O}(p^3)$ asymptotically dominates over $\mathcal{O}(p^2)$ and $\mathcal{O}(pn)$ we can write $\mathcal{O}(p^2 n + p^3)$.

TODO CO zo toho je vic? Neni casove narocnejši vynasobení $\mathbf{X}^T \mathbf{X}$ nez inverze, když bereme v úvahu $n \gg p$??? \square

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. reference: **TODO** 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 $\mathcal{O}(p^2 n + p^3)$ which is asymptotically dominant against previous steps which are $\mathcal{O}(np + n)$. \square

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 10. 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 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

3. ALGORITHMS

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. *TODO* nějaký hezký grafík který to ukazuje.... 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. [3]

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 .

3.2.2 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.

3.2.2.1 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 11. 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

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

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.

3.2.2.2 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 $\hat{\mathbf{w}}_0 = \text{OLS}(J)$ 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 12. 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 [3]. We'll properly discuss them in chapter **TODO**.*

3. ALGORITHMS

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.

3.2.3 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.

3.2.3.1 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 loose 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 $RSS(\hat{\mathbf{w}}_3)$ or $RSS(\hat{\mathbf{w}}_4)$ 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 12. 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 **TODO**. So we can say that with $p < 10$ choosing $m = 500$ is usually safe starting point.*

3.2.3.2 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 $H_{P_{i1}}$ 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 $H_{P_{i3}}$ 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.

3.2.3.3 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_iiteration* and *threshold* which will limit number of C-steps iterations in some rare cases when convergence is too slow. Parameter *max_iiteration* 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

3. 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  $\hat{\mathbf{w}}_{final}, H_{final};$ 

```

3.3 Exact algorithm

3.4 Feasible solution

3.5 MMEA

3.6 Branch and bound

3.7 Adding row

Experiments

- 4.1 Data
- 4.2 Results
- 4.3 Outlier detection

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] 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.

Datasets

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