

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 8, 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 8, 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

<b>Introduction</b>	<b>1</b>
<b>1 Linear Regression</b>	<b>3</b>
1.1 Description . . . . .	3
1.2 Computation . . . . .	3
1.3 Downfalls . . . . .	3
<b>2 The Least trimmed squares</b>	<b>5</b>
<b>3 Algorithms</b>	<b>7</b>
3.1 FAST-LTS . . . . .	7
3.2 Exact algorithm . . . . .	13
3.3 Feasible solution . . . . .	13
3.4 MMEA . . . . .	13
3.5 Branch and bound . . . . .	13
3.6 Adding row . . . . .	13
<b>4 Experiments</b>	<b>15</b>
4.1 Data . . . . .	15
4.2 Results . . . . .	15
4.3 Outlier detection . . . . .	15
<b>Conclusion</b>	<b>17</b>
<b>Bibliography</b>	<b>19</b>
<b>A Datasets</b>	<b>21</b>
<b>B Contents of enclosed CD</b>	<b>23</b>



---

## List of Figures



---

# Introduction





# Linear Regression

- 1.1 Description
- 1.2 Computation
- 1.3 Downfalls



# **The Least trimmed squares**

## **2.0.1 Objective function**

### **2.0.1.1 Problems**



# Algorithms

## 3.1 FAST-LTS

In this section we will introduce FAST-LTS algorithm[1]. 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.

### 3.1.1 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 1.** 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

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

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

### 3. ALGORITHMS

---

$\sum_{i \in H_1} (r_0(i))^2$ . 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 2.** 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 3.** Time complexity of algorithm C-step 1 is the same as time complexity as OLS. Thus  $O(p^2n)$  **TODO**

**Lemma 4.** 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  $\mathcal{O}(p^2n + p^3)$ .

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

*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  $\mathcal{O}(p^2n + p^3)$  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 5.** 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 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 hezky 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. [1]

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  $RSS(\hat{\mathbf{w}})$ .

### 3. ALGORITHMS

---

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  $H_1$ .

#### 3.1.2 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  $H_1$  subset. Lets start with most trivial one.

##### 3.1.2.1 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 6.** 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

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

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

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

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

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

$$\begin{aligned} \text{Because } n \rightarrow \infty \Rightarrow (1 - \epsilon)^h \rightarrow 0 \Rightarrow 1 - (1 - \epsilon)^h \rightarrow 1 \Rightarrow (1 - (1 - \epsilon)^h)^m \rightarrow \\ 1 \Rightarrow 1 - (1 - (1 - \epsilon)^h)^m \rightarrow 0 \end{aligned} \quad \square$$

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.

##### 3.1.2.2 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 7.** 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. □

Note that there are other possibilities of choosing  $H_1$  subset other than these presented in [1]. We'll properly discuss them in chapter **TODO**.

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.

### 3.1.3 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.

#### 3.1.3.1 Selective iteration

---

**Input:** A finite set  $A = \{a_1, a_2, \dots, a_n\}$  of integers  
**Output:** The largest element in the set

```

1  $max \leftarrow a_1$ 
2 for  $i \leftarrow 2$  to  $n$  do
3   | if  $a_i > max$  then
4   |   |  $max \leftarrow a_i$ 
5   | end
6 end
7 return  $max$ 
```

---

Algorithm 8 is a greedy change-making algorithm (Slide 19 in Class Slides).

Algorithm 14 and Algorithm 13 will find the first duplicate element in a sequence of integers.

### 3. ALGORITHMS

---

---

---

**Input:** A set  $C = \{c_1, c_2, \dots, c_r\}$  of denominations of coins, where  $c_1 > c_2 > \dots > c_r$  and a positive number  $n$

**Output:** A list of coins  $d_1, d_2, \dots, d_k$ , such that  $\sum_{i=1}^k d_i = n$  and  $k$  is minimized

```
1  $C \leftarrow \emptyset$ 
2 for  $i \leftarrow 1$  to  $r$  do
3   while  $n \geq c_i$  do
4      $C \leftarrow C \cup \{c_i\}$ 
5      $n \leftarrow n - c_i$ 
6   end
7 end
8 return  $C$ 
```

---

---

---

**Input:** A sequence of integers  $\langle a_1, a_2, \dots, a_n \rangle$

**Output:** The index of first location with the same value as in a previous location in the sequence

```
1  $location \leftarrow 0$ 
2  $i \leftarrow 2$ 
3 while  $i \leq n$  and  $location = 0$  do
4    $j \leftarrow 1$ 
5   while  $j < i$  and  $location = 0$  do
6     if  $a_i = a_j$  then
7        $location \leftarrow i$ 
8     else
9        $j \leftarrow j + 1$ 
10    end
11  end
12   $i \leftarrow i + 1$ 
13 end
14 return  $location$ 
```

---

---

**Input:** A sequence of integers  $\langle a_1, a_2, \dots, a_n \rangle$

**Output:** The index of first location with the same value as in a previous location in the sequence

```
1 location  $\leftarrow$  0
2 i  $\leftarrow$  2
3 while  $i \leq n \wedge \textit{location} = 0$  do
4   j  $\leftarrow$  1
5   while  $j < i \wedge \textit{location} = 0$  do
6     if  $a_i = a_j$  then location  $\leftarrow$  i
7   else j  $\leftarrow$  j + 1
8   end
9   i  $\leftarrow$  i + 1
10 end
11 return location
```

---

In this section we will describe FAST-LTS algorithm and its main properties. The main idea of this algorithm is based on the fact that from one approximation of the algorithm we can compute another which can have lower objective function. TA DAAAAAAAAAAAAAAAAAAAA Theorem 1: [2] Let  $w_0$  ...  $w_p$  be the LTS estimate. for each data sample we can compute  $\|y - wx\|$

## 3.2 Exact algorithm

## 3.3 Feasible solution

## 3.4 MMEA

## 3.5 Branch and bound

## 3.6 Adding row



# Experiments

- 4.1 Data
- 4.2 Results
- 4.3 Outlier detection



---

## Conclusion





---

## Bibliography

- [1] 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.
- [2] Rybicka, J. *LaTeX pro začátečníky*. Brno: Konvoj, third edition, ISBN 80-7302-049-1.



## 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 $\text{\LaTeX}$ 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