

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 6, 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 6, 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
1 Linear Regression	3
1.1 Description	3
1.2 Computation	3
1.3 Downfalls	3
2 The Least trimmed squares	5
3 Algorithms	7
3.1 FAST-LTS	7
3.2 Exact algorithm	9
3.3 Feasible solution	9
3.4 MMEA	9
3.5 Branch and bound	9
3.6 Adding row	9
4 Experiments	11
4.1 Data	11
4.2 Results	11
4.3 Outlier detection	11
Conclusion	13
Bibliography	15
A Datasets	17
B Contents of enclosed CD	19

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'll introduce FAST-LTS algorithm[1]. It's, as well as in other cases, iterative algorithm. We'll discuss all main components of the algorithm starting with its core idea called concentration step which authors simply calls C-step.

3.1.1 C-step

We'll 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 old one. Based on this property we'll 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 $RSS(\hat{\mathbf{w}}_1) = \sum_{i \in H_1} (r_i^1)^2$ sum of least squares. Then

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

Proof. Because we took 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 $RSS(\hat{\mathbf{w}}_1) = \sum_{i \in H_1} (r_i^1)^2 \leq$

3. ALGORITHMS

$\sum_{i \in H_1} (r_i^0)^2$. Together we get $RSS(\hat{\mathbf{w}}_1) = \sum_{i \in H_1} (r_i^1)^2 \leq \sum_{i \in H_1} (r_0(i))^2 \leq \sum_{i \in H_0} (r_0(i))^2 = RSS(\hat{\mathbf{w}}_0)$ \square

$$F(x) := \frac{a_0}{2} + \sum_{k=1}^{\infty} a_k \cos \frac{k\pi x}{T} + b_k \sin \frac{k\pi x}{T}, \quad \forall x \in \mathbb{R},$$

Data: this text

Result: how to write algorithm with L^AT_EX2e

```

1 initialization;
2 while not at end of this document do
3     read current;
4     if understand then
5         go to next section;
6         current section becomes this one;
7     else
8         go back to the beginning of current section;
9     end
10 end
```

In this section we'll describe FAST-LTS algorithm and it's 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 Thoeerem 1: [2] Let $w_0 \dots w_p$ be the LTS estimate. for each data sample we can compute $-y-wx-$

Hlavní myšlenka tohoto algoritmu spočívá ve faktu,

V $\check{\text{A}}\text{Desk}\check{\text{A}}\text{I}$ variant $\check{\text{A}}\check{\text{Z}}$ naleznete $\check{\text{A}}\check{\text{a}}\text{ablony}$ v souborech pojmenovan $\check{\text{A}} \cdot \text{ch}$ ve form $\check{\text{A}}\check{\text{a}}\text{tu}$ pr $\check{\text{A}}\check{\text{a}}\text{ce.k}\check{\text{A}}\check{\text{s}}\text{dov}\check{\text{A}}\check{\text{a}}\text{n}\check{\text{A}}\check{\text{n}}\text{.tex}$. Typ pr $\check{\text{A}}\check{\text{a}}\text{ce}$ m $\check{\text{A}}\check{\text{r}}\check{\text{A}}\text{"e}$ b $\check{\text{A}} \cdot \text{t}$:

BP bakal $\check{\text{A}}\check{\text{a}}\check{\text{Z}}\text{sk}\check{\text{A}}\check{\text{a}}$ pr $\check{\text{A}}\check{\text{a}}\text{ce}$,

DP diplomov $\check{\text{A}}\check{\text{a}}$ (magistersk $\check{\text{A}}\check{\text{a}}$) pr $\check{\text{A}}\check{\text{a}}\text{ce}$.

UTF-8 k $\check{\text{A}}\check{\text{s}}\text{dov}\check{\text{A}}\check{\text{a}}\text{n}\check{\text{A}}\check{\text{n}}$ Unicode,

ISO-8859-2 latin2,

Windows-1250 znakov $\check{\text{A}}\check{\text{a}}$ sada 1250 Windows.

V p $\check{\text{A}}\check{\text{Z}}\check{\text{A}}\text{npad}\check{\text{A}}\check{\text{Z}}$ nejistoty ohledn $\check{\text{A}}\check{\text{Z}}$ k $\check{\text{A}}\check{\text{s}}\text{dov}\check{\text{A}}\check{\text{a}}\text{n}\check{\text{A}}\check{\text{n}}$ doporu $\check{\text{A}}\check{\text{D}}$ ujeme n $\check{\text{A}}\check{\text{a}}\text{sle-}$ duj $\check{\text{A}}\check{\text{n}}\text{c}\check{\text{A}}\check{\text{n}}$ postup:

1. V opa $\check{\text{A}}\check{\text{D}}\text{n}\check{\text{A}}\check{\text{I}}\text{m}$ p $\check{\text{A}}\check{\text{Z}}\check{\text{A}}\text{npad}\check{\text{A}}\check{\text{Z}}$ postupujte d $\check{\text{A}}\check{\text{a}}\text{le}$ podle toho, jak $\check{\text{A}} \cdot \text{op-}$ era $\check{\text{A}}\check{\text{D}}\text{n}\check{\text{A}}\check{\text{n}}$ syst $\check{\text{A}}\check{\text{I}}\text{m}$ pou $\check{\text{A}}\text{"}\check{\text{A}}\text{n}\check{\text{v}}\check{\text{A}}\text{te}$:
 - v p $\check{\text{A}}\check{\text{Z}}\check{\text{A}}\text{npad}\check{\text{A}}\check{\text{Z}}$ Windows pou $\check{\text{A}}\text{"}$ ijte $\check{\text{A}}\check{\text{a}}\text{blonu}$ pro k $\check{\text{A}}\check{\text{s}}\text{dov}\check{\text{A}}\check{\text{a}}\text{n}\check{\text{A}}\check{\text{n}}$ Windows-1250,
 - jinak zkuste pou $\check{\text{A}}\text{"}\check{\text{A}}\text{n}\check{\text{t}}$ $\check{\text{A}}\check{\text{a}}\text{blonu}$ pro k $\check{\text{A}}\check{\text{s}}\text{dov}\check{\text{A}}\check{\text{a}}\text{n}\check{\text{A}}\check{\text{n}}$ ISO-8859-2.

V anglických variantách jsou úkoly pojmenovány podle typu práce, možnosti jsou:

bachelors bakalářské práce,

masters diplomová (magisterská) práce.

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 \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