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 5, 2019

# Acknowledgements

# 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 5, 2019                    . . . . . . . . . . . . . . . . . . . . .

**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, lineÃąrnÃ regrese, optimalizace, nejmenÅąÃ usekanÃľ ÄŊtvrece, metoda nejmenÅąÃŋch ÄŊtvercÅŕ, 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 regressin, robust estimator, least trimmed squares, ordinary least squares, outliers, outliers detection

# Contents

# 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 $\hat{\boldsymbol{w}}_{old}$ we can construct new LTS estimate $\hat{\boldsymbol{w}}_{new}$ 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:** Let's have dataset consisting of $x_1, x_2...x_n$ explanatory variables and its corresponding $y_1, y_2...y_n$ response variables where $x_i$

**Theorem 2:** BuÄŔ $f$ funkce, ktera ma[1] po ÄDÃ ąstech spojitou derivaci na intervalu $\langle -T, T \rangle$. Potom Fourierova ÅŹada funkce $f$ na intervalu $(-T, T)$ konverguje na celÃľ mnoÅ"inÄŻ $\mathbb{R}$. OznaÄDme $F$ jejÃŋ souÄDtovou funkci, tzn.

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

kde posloupnosti $(a_n)_{n=0}^{\infty}$ a $(b_n)_{n=1}^{\infty}$ jsou urÄDeny vztahy . Potom platÃŋ:

(i) $F$ je periodickÃ ą funkce s periodou $2T$.

(ii) $F(x) = \dfrac{f(x+) + f(x-)}{2}$ pro kaÅ"dÃľ $x \in (-T, T)$.

(iii) $F(T) = F(-T) = \dfrac{f(-T+) + f(T-)}{2}$.

*Proof.* NeuvÃądÃŋme.                                                                    □

> **Data:** this text
> **Result:** how to write algorithm with LaTeX2e

**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 algorthm is based on the fact that from one approximation of the algorithm we can compute another which can have lower objective function. TA DAAAAAAAAAAAAAAAA Thoerem 1: [2] Let w0 ... wp be the LTS estimate. for each data sample we can compute —y-wx—

Hlavni myslenka tohoto algoritmu spociva ve faktu,

V ÄŊeskÃĬ variantÄŻ naleznete Åąablony v souborech pojmenovanÃ·ch ve formÃątu prÃące_kÃşdovÃąnÃŋ.tex. Typ prÃące mÃĕÅĕÅ"e bÃ·t:

**BP** bakalÃąÅŹskÃą prÃące,

**DP** diplomovÃą (magisterskÃą) prÃące.

**UTF-8** kÃşdovÃąnÃŋ Unicode,

**ISO-8859-2** latin2,

**Windows-1250** znakovÃą sada 1250 Windows.

V přÃŋpadÄŻ nejistoty ohlednÄŻ kÃşdovÃąnÃŋ doporuÄŊujeme nÃąsledujÃŋcÃŋ postup:

1. V opaÄŊnÃĬm přÃŋpadÄŻ postupujte dÃąle podle toho, jakÃ· operaÄŊnÃŋ systÃĬm pouÅ"ÃŋvÃąte:

   - v přÃŋpadÄŻ Windows pouÅ"ijte Åąablonu pro kÃşdovÃąnÃŋ Windows-1250,
   - jinak zkuste pouÅ"Ãŋt Åąablonu pro kÃşdovÃąnÃŋ ISO-8859-2.

V anglickÃĬm variantÄŻ jsou Åąablony pojmenovanÃĬ podle typu prÃące, moÅ"nosti jsou:

**bachelors** bakalÃąÅŹskÃą prÃące,

**masters** diplomovÃą (magisterskÃą) prÃące.

## 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