

# Least sum of squares of trimmed residuals regression

Yijun Zuo and Hanwen Zuo

*Department of Statistics and Probability and Department of Computer Science  
Michigan State University, East Lansing, MI 48824, USA  
e-mail: [zuohanwe@msu.edu](mailto:zuohanwe@msu.edu); [zuo@msu.edu](mailto:zuo@msu.edu)*

**Abstract:** In the famous least sum of trimmed squares (LTS) estimator [21], residuals are first squared and then trimmed. In this article, we first trim residuals – using a depth trimming scheme – and then square the remaining of residuals. The estimator that minimizes the sum of trimmed and squared residuals, is called an LST estimator.

Not only is the LST a robust alternative to the classic least sum of squares (LS) estimator. It also has a high finite sample breakdown point and can resist, asymptotically, up to 50% contamination without breakdown – in sharp contrast to the 0% of the LS estimator.

The population version of the LST is Fisher consistent, and the sample version is strong, root- $n$  consistent, and asymptotically normal. We propose approximate algorithms for computing the LST and test on synthetic and real data sets. Despite being approximate, one of the algorithms compute the LST estimator quickly with relatively small variances in contrast to the famous LTS estimator. Thus, evidence suggests the LST serves as a robust alternative to the LS estimator and is feasible even in high dimension data sets with contamination and outliers.

**MSC2020 subject classifications:** Primary 62J05, 62G36; secondary 62J99, 62G99.

**Keywords and phrases:** Trimmed residuals, robust regression, finite sample breakdown point, consistency, approximate computation algorithm.

Received November 2022.

## Contents

1	Introduction . . . . .	2417
2	Least sum of squares of trimmed residuals estimator . . . . .	2420
	2.1 Trimming schemes . . . . .	2420
	2.2 Definition and properties of the LST . . . . .	2421
	2.3 Existence, uniqueness and equivariance . . . . .	2423
3	Robustness of LST . . . . .	2424
	3.1 Finite sample breakdown point . . . . .	2424
	3.2 Influence function . . . . .	2426
4	Consistency . . . . .	2428
	4.1 Fisher consistency . . . . .	2428
	4.2 Strong consistency . . . . .	2429

---

arXiv: [2202.10329](https://arxiv.org/abs/2202.10329)

4.3	$\sqrt{n}$ -consistency . . . . .	2431
5	Asymptotic normality . . . . .	2432
6	Computation . . . . .	2433
6.1	A procedure based Theorem 2.1 . . . . .	2434
6.2	A subsampling procedure . . . . .	2435
7	Examples and comparison . . . . .	2436
8	Final discussions . . . . .	2441
	Acknowledgments . . . . .	2442
	Funding . . . . .	2443
	Supplementary Material . . . . .	2443
	References . . . . .	2443

## 1. Introduction

In the classical regression analysis, we assume that there is a relationship for a given data set  $\{(\mathbf{x}'_i, y_i)', i \in \{1, 2, \dots, n\}\}$ :

$$y_i = (1, \mathbf{x}'_i)\boldsymbol{\beta}_0 + e_i, \quad i \in \{1, \dots, n\} \quad (1.1)$$

where  $y_i \in \mathbb{R}^1$ ,  $'$  stands for the transpose,  $\boldsymbol{\beta}_0 = (\beta_{01}, \dots, \beta_{0p})'$  (the true unknown parameter) in  $\mathbb{R}^p$  and  $\mathbf{x}_i = (x_{i1}, \dots, x_{i(p-1)})'$  in  $\mathbb{R}^{p-1}$  ( $p > 1$ ),  $e_i \in \mathbb{R}^1$  is called an error term (or random fluctuation/disturbances, which is usually assumed to have zero mean and variance  $\sigma^2$  in classic regression theory). That is,  $\beta_{01}$  is the intercept term of the model. Write  $\mathbf{w}_i = (1, \mathbf{x}'_i)'$ , then one has  $y_i = \mathbf{w}'_i\boldsymbol{\beta}_0 + e_i$ , which will be used interchangeably with model (1.1).

One wants to estimate the  $\boldsymbol{\beta}_0$  based on a given sample  $\mathbf{Z}^{(n)} := \{(\mathbf{x}'_i, y_i)', i \in \{1, \dots, n\}\}$  from the model  $y = (1, \mathbf{x}')\boldsymbol{\beta}_0 + e$ . We call the difference between  $y_i$  and  $\mathbf{w}'_i\boldsymbol{\beta}$  the  $i$ th residual,  $r_i(\boldsymbol{\beta})$ , for a candidate coefficient vector  $\boldsymbol{\beta}$  (which is often suppressed). That is,

$$r_i := r_i(\boldsymbol{\beta}) = y_i - \mathbf{w}'_i\boldsymbol{\beta}. \quad (1.2)$$

To estimate  $\boldsymbol{\beta}_0$ , the classic *least squares* (LS) minimizes the sum of squares of residuals,

$$\widehat{\boldsymbol{\beta}}_{ls} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^n r_i^2.$$

Alternatively, one can replace the square above by the absolute value to obtain the least absolute deviations estimator (the  $L_1$  estimator, in contrast to the  $L_2$  (LS) estimator).

The LS estimator is very popular in practice across a broader spectrum of disciplines due to its great computability and optimal properties when the error  $e_i$  follows a normal  $\mathcal{N}(\mathbf{0}, \sigma^2)$  distribution. However, it can behave badly when the error distribution departs even slightly from the normal distribution, particularly when the errors are heavy-tailed or contain outliers.

Robust alternatives to the  $\widehat{\boldsymbol{\beta}}_{ls}$  have been present in the literature for a long time. The most popular ones are, among others, the M-estimators [14], least

median squares (LMS) and least trimmed squares (LTS) estimators [21], S-estimators [27], MM-estimators [46],  $\tau$ -estimators [47], and maximum depth estimators ([22, 52], and [53]). For more related discussions, see Sects. 1.2 and 4.4 of [23], and Sect. 5.14 of [17].

In practice, the LTS is the most common estimator used across multiple disciplines. Its idea is simple, it orders the squared residuals and then trims the larger ones keeping at least  $\lceil n/2 \rceil$  squared residuals, where  $\lceil \cdot \rceil$  is the ceiling function; finally, the minimizer of the sum of those *trimmed squared residuals* is called the LTS estimator:

$$\widehat{\beta}_{lts} := \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^h (r^2)_{i:n},$$

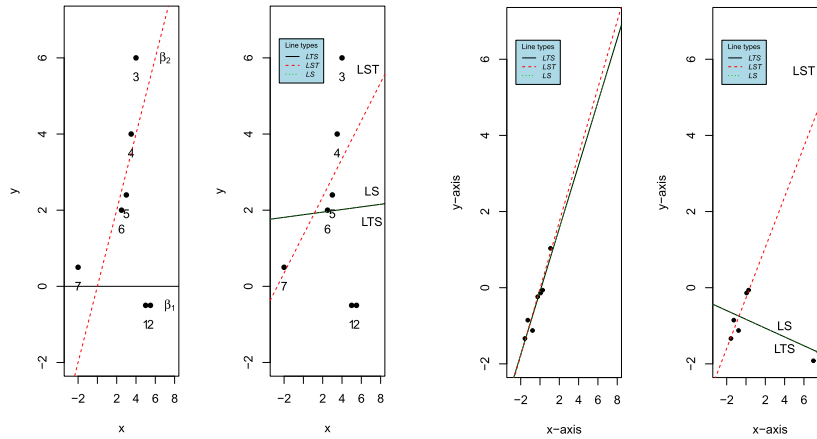
where  $(r^2)_{1:n} \leq (r^2)_{2:n} \leq \dots, (r^2)_{n:n}$  are the ordered squared residuals and constant  $h$  satisfies  $\lceil n/2 \rceil \leq h \leq n$ .

Naturally we wonder, what if we first trim (employing the scheme given in Sect. 2) the residuals and then minimize the sum of *squares of trimmed residuals*? Thus, the minimizer will be called LST. Is there any difference between the two procedures? Outlying (large or small) original residuals are trimmed after squaring in the LTS – those residuals are certainly trimmed in the LST. However, the outlying residuals which have a small squared magnitude will not be trimmed in the LTS and are trimmed in the LST (see (a) of Fig. 1). Before formally introducing the LST in Sect. 2, let us first appreciate the difference between the two procedures.

**Example 1.1.** We constructed a small data set in  $\mathbb{R}^2$  with  $\mathbf{x} = (5, 5.5, 4, 3.5, 3, 2.5, -2)$  and  $\mathbf{y} = (-.5, -.5, 6, 4, 2.4, 2, .5)$ , they are plotted in the left panel of the (a) of Fig. 1. We also provide two candidate regression lines  $\beta_1$  ( $y = 0$ ) and  $\beta_2$  ( $y = x$ ). Which better represents the overall pattern of the data set?

If we use the number  $h = \lfloor n/2 \rfloor + \lfloor (p+1)/2 \rfloor$  given on page 132 of [23] to achieve the maximum possible breakdown point (see Sect. 3 for definition) for the LTS estimator, or employing the four smallest squared residuals, then the LTS prefers  $\beta_1$  (using residuals from points 1, 2, 6, and 7) to  $\beta_2$  (using points 4, 5, 6, and 7), whereas for the LST,  $\beta_2$  (using residuals from points 4, 5, 6, and 7) is preferred. Some would argue that this is not representative since the LTS searches all possible (not just two) lines and outputs the best one, but if one utilized the R function `ltsReg`, then it produces the solid (black) line whereas the line based on algorithms (see Sect. 5) for the LST is the dashed (red) one in the right panel of the (a) of Fig. 1. For benchmark purposes, the LS line dotted (green) is also given, which overlaps with the LTS line. From this instance, we can appreciate the difference between trimming schemes of the LTS and the LST. Of course, some might argue that the data set in (a) is purely synthetic and fixed.

So, in (b) of Fig. 1, we generated seven highly correlated normal points (with correlation 0.88 between  $x$  and  $y$ ), when there is no contamination the LTS



(a) Left panel: plot of seven artificial points and two candidate lines ( $\beta_1$  and  $\beta_2$ ), which line would you pick? Sheerly based on the trimming scheme and objective function value, if one uses the number  $h = \lfloor n/2 \rfloor + \lfloor (p+1)/2 \rfloor$  given on page 132 of [23], that is, employing four smallest squared residuals, then the LTS prefers  $\beta_1$  to  $\beta_2$  whereas the LST reverses the preference.

Right panel: the same seven points are fitted by the LTS, the LST, and the LS (benchmark). A solid black line is the LTS given by `ltsReg`. Red dashed line is given by the LST, and green dotted line is given by the LS - which is identical to the LTS line in this case.

(b) Left panel: plot of seven highly correlated normal points (with mean being the zero vector and covariance matrix with diagonal entries being one and off-diagonal entries being 0.88) and three lines given by the LST, the LTS, and the LS. The LS line is identical to the LTS line again.

Right panel: The LTS line (solid black) and the LST line (dashed red), and the LS (dotted green) for the same seven highly correlated normal points but with two points contaminated nevertheless. The LS line is identical to the LTS line due to the attributes in the R function `ltsReg` that is based on [26]).

FIG 1. (a) Difference between the two procedures: the LST and the LTS. (b) Performance difference between the LST and the LTS when there are contaminated points (x-axis leverage points).

(identical to the LS again) and the LST pick out the linear pattern. If there are two contaminated points (note that the LTS allows  $m := \lfloor (n - p)/2 \rfloor = 2$  contaminated points in this case in light of Theorem 6 on page 132 of [23]), the LTS line changes drastically in this instance becoming identical to the LS again.

For examples with an increased sample size, see Sect. 6. Incidentally, the instability of the LMS (not the LTS) was already documented in [13].  $\square$

The idea of trimming residuals and then doing regression has appeared in the literature for quite some time. The trimming idea was first introduced in

location setting and later extended to regression, see, [15, 2, 28, 44], and [23], among others. The trimmed mean has been used in practice for more than two centuries (see [8], page 34, and is attributed to “Anonymous” (1821) [1] (Gergonne, see [33]), or [18]. Tukey [37, 4] was an outstanding advocator for the trimmed mean in the last century.

However, trimming residuals based on depth or outlyingness employed in this article (see Sect. 2) is novel and has never been utilized before. A more recent study on the topic is given by Johansen and Nielsen (2013), where the authors used an iterated one-step approximation to the Huber-skip estimator to detect outliers in regression and theoretical justification for the approximation is provided. Their Huber-skip estimator defined on page 56 is closely related to our LST, but there are two essential differences (i) their estimator more resembles the least winsorized squares regression (see page 135 of [23]), (ii) residuals in their estimator are not centered by the median of residuals.

In light of [52], both the LTS and the LST could be regarded as the deepest estimator (aka regression median) with respect to the corresponding objective function type of regression depth (see Sect. 2.3.1 of [52] and Sect. 4).

The rest of the article is organized as follows. Section 2 introduces trimming schemes and the least sum of squares of trimmed (LST) residuals estimator and establishes the existence and equivariance properties. Section 3 investigates the robustness of the LST in terms of its finite sample breakdown point and its influence function. Section 4 establishes the Fisher as well as the strong and the root-n consistency. The asymptotic normality is derived from stochastic equicontinuity in Sect. 5. Section 6 introduces two approximate computation algorithms of the LST. Section 7 presents examples of simulated and real data and carries out a comparison against the leading regression estimators, the LTS and the LMS. Section 8 consists of concluding discussions. Long proofs are deferred to Appendix.

## 2. Least sum of squares of trimmed residuals estimator

### 2.1. *Trimming schemes*

**Rank based trimming** This scheme is based on the ranks of data points, usually trimming an equal number of points at both tails of a data set (that is, lower or higher rank points are trimmed) and can trim points one-sided if needed (such as when all data points lie on the positive side of the number axis).

This scheme is related to the trimmed mean, which keeps a good balance between robustness and efficiency, alleviating the extreme sensitivity of the sample mean and enhancing the efficiency of the sample median.

Rank-based trimming focuses only on the relative position of points with respect to others and ignores the magnitude of the point and the relative distance between points. [49] and [45] discuss an alternative trimming scheme, which catches these two important attributes (magnitude and relative distance). It

orders data from a center (the median) outward and trims the points that are far away from the center. This is known as depth-based trimming.

**Depth (or outlyingness) based trimming** In other words, the depth-based trimming scheme trims points that lie on the outskirts (i.e. points that are less deep, or outlying). The outlyingness (or, equivalently, depth) of a point  $x$  is defined to be (strictly speaking,  $\text{depth} = 1/(1 + \text{outlyingness})$ ) in [48])

$$D(x, X^{(n)}) = |x - \text{Med}(X^{(n)})| / \text{MAD}(X^{(n)}), \quad (2.1)$$

where  $X^{(n)} = \{x_1, \dots, x_n\}$  is a data set in  $\mathbb{R}^1$ ,  $\text{Med}(X^{(n)}) = \text{median}(X^{(n)})$  is the median of the data points, and  $\text{MAD}(X^{(n)}) = \text{Med}(\{|x_i - \text{Med}(X^{(n)})|, i \in \{1, 2, \dots, n\}\})$  is the median of absolute deviations to the center (median). It is readily seen that  $D(x, X^{(n)})$  is a generalized standard deviation, or equivalent to the one-dimensional projection depth/outlyingness (see [55] and [48, 49] for a high dimensional version). For notion of outlyingness, cf. [32, 5], and [6].

The LTS essentially employs one-sided rank based trimming scheme (w.r.t. squared residuals), whereas depth based trimming is utilized in the LST which is introduced next.

## 2.2. Definition and properties of the LST

**Definition.** For a given sample  $\mathbf{Z}^{(n)} = \{(x'_i, y_i)', i \in \{1, 2, \dots, n\}\}$  in  $\mathbb{R}^p$  from  $y = \mathbf{w}'\beta_0 + e$  and a  $\beta \in \mathbb{R}^p$ , define

$$m_n(\beta) := m(\mathbf{Z}^{(n)}, \beta) = \text{Med}_i\{r_i\}, \quad (2.2)$$

$$\sigma_n(\beta) := \sigma(\mathbf{Z}^{(n)}, \beta) = \text{MAD}_i\{r_i\}, \quad (2.3)$$

where operators  $\text{Med}$  and  $\text{MAD}$  are used for discrete data sets (and distributions as well) and  $r_i$  defined in (1.2). For a constant  $\alpha$  in the depth trimming scheme, consider the quantity

$$Q(\mathbf{Z}^{(n)}, \beta, \alpha) := \sum_{i=1}^n r_i^2 \mathbb{1} \left( \frac{|r_i - m(\mathbf{Z}^{(n)}, \beta)|}{\sigma(\mathbf{Z}^{(n)}, \beta)} \leq \alpha \right), \quad (2.4)$$

where  $\mathbb{1}(A)$  is the indicator of  $A$  (i.e., it is one if  $A$  holds and zero otherwise). Namely, residuals with their outlyingness (or equivalently reciprocal of depth minus one) greater than  $\alpha$  will be trimmed. When there is a majority ( $\geq \lfloor (n+1)/2 \rfloor$ ) of identical  $r_i$ s, we define  $\sigma(\mathbf{Z}^{(n)}, \beta) = 1$  (since those  $r_i$  lie in the deepest position (or are the least outlying points)).

Minimizing  $Q(\mathbf{Z}^{(n)}, \beta, \alpha)$ , We get the *least sum of squares of trimmed* (LST) residuals estimator,

$$\widehat{\beta}_{lst}^n := \widehat{\beta}_{lst}(\mathbf{Z}^{(n)}, \alpha) = \arg \min_{\beta \in \mathbb{R}^p} Q(\mathbf{Z}^{(n)}, \beta, \alpha). \quad (2.5)$$

Does the right-hand side (RHS) of (2.5) always have a minimizer? If it exists, is it unique? We treat this problem formally next.

Hereafter we will assume that  $\alpha \geq 1$ . That is, we will keep the residuals that are no greater than one MAD away from the center (the median of residuals) untrimmed. For a given  $\alpha$ ,  $\beta$ , and  $\mathbf{Z}^{(n)}$ , define a set of indexes for  $1 \leq i \leq n$

$$I(\beta) = \left\{ i : \frac{|r_i - m(\mathbf{Z}^{(n)}, \beta)|}{\sigma(\mathbf{Z}^{(n)}, \beta)} \leq \alpha \right\}. \quad (2.6)$$

Namely, the set of subscripts so that the outlyingness (see (2.1)) of the corresponding residuals are no greater than  $\alpha$ . It depends on  $\mathbf{Z}^{(n)}$  and  $\alpha$ , which are suppressed in the notation. Following the convention, we denote the cardinality of set  $A$  by  $|A|$ . We have

**Lemma 2.1.** *For any  $\beta \in \mathbb{R}^p$  and the given  $\mathbf{Z}^{(n)}$  and  $\alpha$ ,  $|I(\beta)| \geq \lfloor (n+1)/2 \rfloor$ .*

*Proof.* By the definition of MAD (the median of the absolute deviations to the center (median)), it is readily seen that

$$\begin{aligned} |I(\beta)| &= \sum_{i=1}^n \mathbb{1} \left( \frac{|r_i - m(\mathbf{Z}^{(n)}, \beta)|}{\sigma(\mathbf{Z}^{(n)}, \beta)} \leq \alpha \right) \\ &\geq \sum_{i=1}^n \mathbb{1} \left( \frac{|r_i - m(\mathbf{Z}^{(n)}, \beta)|}{\sigma(\mathbf{Z}^{(n)}, \beta)} \leq 1 \right) = \lfloor (n+1)/2 \rfloor, \end{aligned}$$

This completes the proof.  $\square$

The lemma implies that the RHS of (2.4) sums a majority of squared residuals.

**Properties of the objective function** Write  $D_i := D(r_i, \beta) = |r_i - m(\mathbf{Z}^{(n)}, \beta)| / \sigma(\mathbf{Z}^{(n)}, \beta)$  for a given  $\mathbf{Z}^{(n)}$  and  $\beta$ . Let  $i_1, \dots, i_K$  be in  $I(\beta)$  such that  $D_{i_1} \leq D_{i_2} \leq \dots \leq D_{i_K}$  (i.e. ordered depth values of residuals),  $K := |I(\beta)|$ . Both  $i_j$  and  $D_{i_j}$  clearly depend on  $\beta$  and  $\mathbf{Z}^{(n)}$ .

Generally, the inequalities between the  $D_i$ 's cannot be strict unless we assume **(A00)**:  $r(\beta) := y - \mathbf{w}'\beta$  has a density for any  $\beta \in \mathbb{R}^p$ . Hereafter, we assume that **(A00)** holds, then the strict inequalities hold almost surely (a.s.), i.e.,  $D_{i_1} < D_{i_2} < \dots < D_{i_K}$  (a.s.). Define for any  $\beta^1 \in \mathbb{R}^p$  and a given  $\mathbf{Z}^{(n)}$

$$R_{\beta^1} = \{ \beta \in \mathbb{R}^p : I(\beta) = I(\beta^1), D_{i_1}(\beta) < D_{i_2}(\beta) < \dots < D_{i_K}(\beta) \}. \quad (2.7)$$

That is, the set of all  $\beta$ s that share the same index set  $I(\beta^1)$  of  $\beta^1$ . If **(A00)** holds, then  $R_{\beta^1} \neq \emptyset$  (a.s.). For a fixed  $n$ , there are at most finitely many  $R_{\beta^k}$ s,  $\beta^k \in \mathbb{R}^p$ ,  $1 \leq k \leq L := \lfloor \binom{n}{\lfloor (n+1)/2 \rfloor} \rfloor$  such that  $\cup_{k=1}^L \overline{R_{\beta^k}} = \mathbb{R}^p$ , where  $R_{\beta^k}$  is defined similarly to (2.7) and  $\overline{A}$  stands for the closure of the set  $A$ . For any  $\beta \in \mathbb{R}^p$ , either there is  $R_\eta$  and  $\beta \in R_\eta$  or there is  $R_\xi$ , such that  $\beta \notin R_\eta \cup R_\xi$  and  $\beta \in \overline{R_\eta} \cap \overline{R_\xi}$ . In the latter case, there are  $i_k, i_l \in I(\beta)$   $i_k \neq i_l$ , such that  $D_{i_k} = D_{i_l}$ .

For a given sample  $\mathbf{Z}^{(n)}$ , write  $Q^n(\beta)$  for  $Q(\mathbf{Z}^{(n)}, \beta, \alpha)$ ,  $B(\boldsymbol{\eta}, \delta)$  for an open ball in  $\mathbb{R}^p$  centered at  $\boldsymbol{\eta}$  with a radius  $\delta > 0$ , and  $\mathbb{1}_i$ , which depends on  $\beta$ , for  $\mathbb{1}(|r_i - m_n(\beta)|/\sigma_n(\beta) \leq \alpha)$ . Let  $\mathbf{X}_n = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n)'$ ,  $\mathbf{Y}_n = (y_1, \dots, y_n)'$ , and  $\mathbf{M}_n := \mathbf{M}(\mathbf{Y}_n, \mathbf{X}_n, \beta, \alpha) = \sum_{i=1}^n \mathbf{w}_i \mathbf{w}_i' \mathbb{1}_i = \sum_{i \in I(\beta)} \mathbf{w}_i \mathbf{w}_i'$ . Assume **(A0\*)**:  $\mathbf{X}_n$  and any its  $K := |I(\beta)|$  sub-rows ( $K > p$ ) have a full rank  $p$ .

**Lemma 2.2.** *Assume that **(A00)** and **(A0\*)** hold, then*

- (i) *For a given  $\mathbf{Z}^{(n)}$  and  $\alpha$ , for any  $1 \leq k \leq L$  and any  $\boldsymbol{\eta} \in R_{\beta^k}$ , there exists a  $B(\boldsymbol{\eta}, \delta)$  such that for any  $\beta \in B(\boldsymbol{\eta}, \delta)$ ,  $\beta \in R_{\beta^k}$ , i.e.,*

$$Q^n(\beta) = \sum_{i \in I(\beta^k)} r_i^2,$$

- (ii) *For any  $1 \leq k \leq L$ ,  $R_{\beta^k}$  is open,*
- (iii)  *$Q^n(\beta)$  is continuous in  $\beta \in \mathbb{R}^p$ ,*
- (iv) *Over each  $R_{\beta^k}$ ,  $1 \leq k \leq L$ ,  $Q^n(\beta)$  is twice differentiable and convex, and strictly convex if the rank of  $\mathbf{X}_n$  is  $p$ .*

*Proof.* See the Appendix. □

**Remark 2.1.** (i) By discussions above and Lemma 2.2, we see that the domain of  $Q^n(\beta)$  (the parameter space) is partitioned into at most  $L$  pieces and over each piece the graph of  $Q^n(\beta)$  is that of the quadratic function of the sum of squared residuals. Hence the graph of  $Q^n(\beta)$  is composed of at most  $L$  of those components.

- (ii) The continuity deduced from  $Q^n(\beta)$  being the sum of some squared residuals without (i) of Lemma 2.2 might not be flawless. The unified expression for  $Q^n(\beta)$  around the small neighborhood of  $\beta$  such as the one given in (i) of the Lemma 2.2 is indispensable. □

### 2.3. Existence, uniqueness and equivariance

**Theorem 2.1.** *Assume that **(A00)** and **(A0\*)** hold, then*

- (i)  $\widehat{\beta}_{lst}^n$  exists and is the unique local minimum of  $Q^n(\beta)$  over  $R_{\beta^{k_0}}$  for some  $k_0$  ( $1 \leq k_0 \leq L$ ).
- (ii) Over  $R_{\beta^{k_0}}$ ,  $\widehat{\beta}_{lst}^n$  is the solution of the system of equations

$$\sum_{i=1}^n (y_i - \mathbf{w}_i' \beta) \mathbf{w}_i \mathbb{1}_i = \mathbf{0}. \tag{2.8}$$

- (iii) Over  $R_{\beta^{k_0}}$ , the unique solution is

$$\widehat{\beta}_{lst}^n = \mathbf{M}_n(\mathbf{Y}_n, \mathbf{X}_n, \widehat{\beta}_{lst}^n, \alpha)^{-1} \sum_{i \in I(\beta^{k_0})} y_i \mathbf{w}_i. \tag{2.9}$$

*Proof.* See the Appendix. □



Note that **(A0\*)** is sufficient for the matrix in the theorem to be invertible. The existence could also be established as follows. In the sequel, we will assume that

**(A0)** there is no vertical hyperplane which contains at least  $\lfloor (n+1)/2 \rfloor$  points of  $\mathbf{Z}^{(n)}$ .

This holds true with probability one if  $(\mathbf{x}', y)'$  has a joint density or holds if  $\mathbf{Z}^{(n)}$  is *in general position* (see Sect. 3 for definition) (assume that  $n > 2p + 1$  hereafter).

**Theorem 2.2.** *The minimizer  $\widehat{\beta}_{lst}^n$  of  $Q(\mathbf{Z}^{(n)}, \beta, \alpha)$  defined in (2.4) over  $\beta \in \mathbb{R}^p$  always exists for a given  $\mathbf{Z}^{(n)}$  and an  $\alpha$  provided that **(A0)** holds.*

*Proof.* See the Appendix. □

**Equivariance** A regression estimator  $\mathbf{T}$  is called *regression, scale, and affine equivariant* if, respectively (see page 116 of [23]) with  $i \in \mathbb{N} := \{1, 2, \dots, n\}$

$$\begin{aligned} \mathbf{T}(\{\mathbf{w}'_i, y_i + \mathbf{w}'_i \mathbf{b}\}') &= \mathbf{T}(\{\mathbf{w}'_i, y_i\}') + \mathbf{b}, \quad \forall \mathbf{b} \in \mathbb{R}^p \\ \mathbf{T}(\{\mathbf{w}'_i, sy_i\}') &= s\mathbf{T}(\{\mathbf{w}'_i, y_i\}'), \quad \forall s \in \mathbb{R}^1 \\ \mathbf{T}(\{(A' \mathbf{w}_i)', y_i\}') &= A^{-1} \mathbf{T}(\{\mathbf{w}'_i, y_i\}'), \quad \forall \text{nonsingular } A \in \mathbb{R}^{p \times p} \end{aligned}$$

**Theorem 2.3.**  *$\widehat{\beta}_{lst}^n$  is regression, scale, and affine equivariant.*

*Proof.* We have the identities

$$\begin{aligned} y_i + \mathbf{w}'_i \mathbf{b} - \mathbf{w}'_i (\beta + \mathbf{b}) &= y_i - \mathbf{w}'_i \beta, \quad \forall \mathbf{b} \in \mathbb{R}^p \\ sy_i - \mathbf{w}'_i (s\beta) &= s(y_i - \mathbf{w}'_i \beta), \quad \forall s \in \mathbb{R}^1 \\ y_i - (A' \mathbf{w}_i)' A^{-1} \beta &= y_i - \mathbf{w}'_i \beta, \quad \forall \text{nonsingular } A \in \mathbb{R}^{p \times p}. \end{aligned}$$

The desired result follows by these identities and the (regression, scale, and affine) invariance (see page 148 of [52] for definition) of  $\frac{|r_i - m(\mathbf{Z}^{(n)}, \beta)|}{\sigma(\mathbf{Z}^{(n)}, \beta)}$ . □

### 3. Robustness of LST

#### 3.1. Finite sample breakdown point

As an alternative to the least-squares, is the LST estimator more robust? The most prevailing quantitative measure of global robustness of any location or regression estimators in the finite sample practice is the *finite sample breakdown point* (FSBP), introduced by [7].

Roughly speaking, the FSBP is the minimum fraction of ‘bad’ (or contaminated) data points that can force the estimator beyond any bound (becoming useless). For example, in the context of estimating the center of a data set, the sample mean has a breakdown point of  $1/n$  (or 0%), because even one bad observation can change the mean by an arbitrary amount; in contrast, the sample median has a breakdown point of  $\lfloor (n+1)/2 \rfloor / n$  (or 50%).

**Definition 3.1** ([7]). The finite sample *replacement breakdown point* (RBP) of a regression estimator  $\mathbf{T}$  at the given sample  $\mathbf{Z}^{(n)} = \{Z_1, Z_2, \dots, Z_n\}$ , where  $Z_i := (\mathbf{x}'_i, y_i)'$ , is defined as

$$\text{RBP}(\mathbf{T}, \mathbf{Z}^{(n)}) = \min_{1 \leq m \leq n, m \in \mathbb{N}} \left\{ \frac{m}{n} : \sup_{\mathbf{Z}_m^{(n)}} \|\mathbf{T}(\mathbf{Z}_m^{(n)}) - \mathbf{T}(\mathbf{Z}^{(n)})\| = \infty \right\}, \quad (3.1)$$

where  $\mathbf{Z}_m^{(n)}$  denotes an arbitrary contaminated sample by replacing  $m$  original sample points in  $\mathbf{Z}^{(n)}$  with arbitrary points in  $\mathbb{R}^p$ . Namely, the RBP of an estimator is the minimum replacement fraction that could drive the estimator beyond any bound. It turns out that both the  $L_1$  (least absolute deviations) and the  $L_2$  (least squares) estimators have RBP  $1/n$  (or 0%), the lowest possible value whereas the LTS can have  $(\lfloor (n-p)/2 \rfloor + 1)/n$  (or 50%), the highest possible value for any regression equivariant estimators (see pages 124–125 of [23]).

We shall say  $\mathbf{Z}^{(n)}$  is *in general position* when any  $p$  of observations in  $\mathbf{Z}^{(n)}$  gives a unique determination of  $\beta$ . In other words, any  $(p-1)$  dimensional subspace of the space  $(\mathbf{x}', y)'$  contains at most  $p$  observations of  $\mathbf{Z}^{(n)}$ . When the observations come from continuous distributions, the event ( $\mathbf{Z}^{(n)}$  being in general position) happens with probability one.

**Theorem 3.1.** For  $\hat{\beta}_{lst}^n$  defined in (2.5) and  $\mathbf{Z}^{(n)}$  in general position, we have

$$\text{RBP}(\hat{\beta}_{lst}^n, \mathbf{Z}^{(n)}) = \begin{cases} \lfloor (n+1)/2 \rfloor / n, & \text{if } p = 1, \\ (\lfloor n/2 \rfloor - p + 2) / n, & \text{if } p > 1. \end{cases} \quad (3.2)$$

*Proof.* See the Appendix. □

**Remark 3.1.**

- (i) The assumption that  $\mathbf{Z}^{(n)}$  is in general position seems to play a central role in the proof. But actually, one can drop it and introduce an index:  $c(\mathbf{Z}^{(n)})$  (which is the maximum number of observations from  $\mathbf{Z}^{(n)}$  contained in any  $(p-1)$  dimensional subspace/hyperplane) to replace  $p$  in the derivation of the proof and the final RBP result (when  $p > 1$ ).
- (ii) Asymptotically speaking (i.e. as  $n \rightarrow \infty$ ),  $\hat{\beta}_{lst}^n$  has the best possible asymptotic breakdown point (ABP) 50%, the same as that of the LTS. The RBP of  $\hat{\beta}_{lst}^n$ , albeit very high (indeed as high as that of the LMS), is slightly less than that of the LTS (with the best choice of  $h$ ). However, it can be improved to attain the best possible value if one modifies  $\alpha$  so that it is the  $h$ th quantile of the  $n$  outlyingness of residuals with  $h = \lfloor n/2 \rfloor + \lfloor (p+1)/2 \rfloor$  to include exact  $h$  squares of residuals in the sum of the RHS of (2.4). □

### 3.2. Influence function

Throughout  $F_{\mathbf{z}}$  stands for the distribution of random vector  $\mathbf{z}$  unless otherwise stated. Write  $F_{(\mathbf{x}', y)}$  for the joint distribution of  $\mathbf{x}'$  and  $y$  in (1.1),  $r := r(F_{(\mathbf{x}', y)}, \beta) = y - (1, \mathbf{x}')\beta := y - \mathbf{w}'\beta$ .

$$\begin{aligned} m &:= m(F_{(\mathbf{x}', y)}, \beta) = \text{Med}(F_r), \\ \sigma &:= \sigma(F_{(\mathbf{x}', y)}, \beta) = \text{MAD}(F_r), \end{aligned}$$

hereafter we assume that  $m$  and  $\sigma$  exist uniquely. The population counterparts of (2.4) and (2.5) are respectively:

$$Q(F_{(\mathbf{x}', y)}, \beta, \alpha) := \int (y - \mathbf{w}'\beta)^2 \mathbf{1} \left( \frac{|y - \mathbf{w}'\beta - m|}{\sigma} \leq \alpha \right) dF_{(\mathbf{x}', y)}, \quad (3.3)$$

$$\beta_{lst}(F_{(\mathbf{x}', y)}, \alpha) := \arg \min_{\beta \in \mathbb{R}^p} Q(F_{(\mathbf{x}', y)}, \beta, \alpha). \quad (3.4)$$

The RBP gauges the global robustness of an estimator at finite sample practice. To assess the local robustness at the population setting, one can use the influence function approach (see [8]), which depicts the local robustness of a functional with an infinitesimal point-mass contamination at a single point  $\mathbf{z} \in \mathbb{R}^p$ .

For a given distribution  $F$  defined on  $\mathbb{R}^p$  and an  $\varepsilon > 0$ , the version of  $F$  contaminated by an  $\varepsilon$  amount of an *arbitrary distribution*  $G$  on  $\mathbb{R}^p$  is denoted by  $F(\varepsilon, G) = (1 - \varepsilon)F + \varepsilon G$  (an  $\varepsilon$  amount deviation from the assumed  $F$ ). Hereafter it is assumed that  $\varepsilon < 1/2$ , otherwise  $F(\varepsilon, G) = G((1 - \varepsilon), F)$ , which means we cannot distinguish which one is contaminated and which is not

**Definition 3.2** ([8]). The *influence function* (IF) of a functional  $T$  at a given point  $\mathbf{z} \in \mathbb{R}^p$  for a given  $F$  is defined as

$$\text{IF}(\mathbf{z}; T, F) = \lim_{\varepsilon \rightarrow 0^+} \frac{T(F(\varepsilon, \delta_{\mathbf{z}})) - T(F)}{\varepsilon}, \quad (3.5)$$

where  $\delta_{\mathbf{z}}$  is the point-mass probability measure at  $\mathbf{z} \in \mathbb{R}^p$ .

The function  $\text{IF}(\mathbf{z}; T, F)$  describes the relative influence on  $T$  of an infinitesimal point-mass contamination at  $\mathbf{z}$  and gauges the local robustness of  $T$ .

It is desirable that a regression estimating functional has a bounded influence function. This, however, does not hold for an arbitrary regression estimating functional (such as the classical least squares functional). Now we investigate this for the functional of the least sum of squares of trimmed residuals,  $\beta_{lst}(F_{(\mathbf{x}', y)}, \alpha)$ . Put

$$\begin{aligned} F_{\varepsilon}(\mathbf{z}) &:= F(\varepsilon, \delta_{\mathbf{z}}) = (1 - \varepsilon)F_{(\mathbf{x}', y)} + \varepsilon\delta_{\mathbf{z}}, \\ m_{\varepsilon}(\mathbf{z}) &:= m(F_{\varepsilon}(\mathbf{z}), \beta) = \text{Med}(F_{R_{\varepsilon}(\mathbf{z})}), \\ \sigma_{\varepsilon}(\mathbf{z}) &:= \sigma(F_{\varepsilon}(\mathbf{z}), \beta) = \text{MAD}(F_{R_{\varepsilon}(\mathbf{z})}), \end{aligned}$$

where  $R_\varepsilon(\mathbf{z}) = r(F_\varepsilon(\mathbf{z}), \beta) = t - (1, \mathbf{s}')\beta$ , and  $F_\varepsilon(\mathbf{z}) =: F_{\mathbf{u}}(\mathbf{z})$  with a random vector  $\mathbf{u} = (\mathbf{s}', t)' \in \mathbb{R}^p$ ,  $\mathbf{s} \in \mathbb{R}^{p-1}$ , and  $t \in \mathbb{R}^1$  (i.e.,  $\mathbf{u}$  is the random vector that has the CDF  $F_\varepsilon(\mathbf{z})$ ). Hereafter we assume that  $m_\varepsilon(\mathbf{z})$  and  $\sigma_\varepsilon(\mathbf{z})$  uniquely exist. The versions of (3.3) and (3.4) at the contaminated distribution  $F_\varepsilon(\mathbf{z})$  are respectively

$$Q(F_\varepsilon(\mathbf{z}), \beta, \alpha) := \int (t - (1, \mathbf{s}')\beta)^2 \mathbf{1} \left( \frac{|(t - (1, \mathbf{s}')\beta) - m_\varepsilon(\mathbf{z})|}{\sigma_\varepsilon(\mathbf{z})} \leq \alpha \right) dF_{\mathbf{u}}(\mathbf{s}', t), \tag{3.6}$$

$$\beta_{lst}(F_\varepsilon(\mathbf{z}), \alpha) := \arg \min_{\beta \in \mathbb{R}^p} Q(F_\varepsilon(\mathbf{z}), \beta, \alpha). \tag{3.7}$$

**Lemma 3.1.**  $\beta_{lst} := \beta_{lst}(F_{(\mathbf{x}', y)}, \alpha)$  is regression, scale, and affine equivariant (see [52] for definition).

*Proof.* It is trivial (analogous to that of Theorem 2.3). □

To investigate the influence function of  $\beta_{lst}$  especially the consistency of its sample version in the next section, we first need to establish its existence and uniqueness. We need assumptions: **(A1)**  $y$  has a density, and **(A2)** the distribution  $F_r$  with  $r = y - \mathbf{w}'\beta$  is non-flat around  $m = \text{Med}(F_r)$  and  $\sigma = \text{MAD}(F_r)$  for any  $\beta \in \mathbb{R}^p$ .

Write  $Q(\beta)$  for  $Q(F_{(\mathbf{x}', y)}, \beta, \alpha)$  in (3.3). We have a population counterpart of Lemma 2.2.

**Lemma 3.2.** Assume **(A1)**–**(A2)** hold. Then  $Q(\beta)$

- (i) is continuous in  $\beta \in \mathbb{R}^p$ ;
- (ii) is twice differentiable in  $\beta \in \mathbb{R}^p$  with (assume that  $E(\mathbf{x}\mathbf{x}')$  exists)

$$\partial^2 Q(\beta) / \partial \beta^2 = 2E\mathbf{w}\mathbf{w}' \mathbf{1} (|y - \mathbf{w}'\beta - m| / \sigma \leq \alpha);$$

- (iii) is convex in  $\beta \in \mathbb{R}^p$  and strictly convex if  $E\mathbf{w}\mathbf{w}' \mathbf{1} (|y - \mathbf{w}'\beta - m| / \sigma \leq \alpha)$  is invertible.

*Proof.* See the Appendix. □

**Theorem 3.2.** Assume that **(A1)**–**(A2)** hold and  $m(F_\varepsilon(\mathbf{z}), \beta)$  and  $\sigma(F_\varepsilon(\mathbf{z}), \beta)$  are continuous in  $\beta$  around a small neighborhood of  $\beta_{lst}(F_\varepsilon(\mathbf{z}), \alpha)$ . Write  $\mathbf{v}' = (1, \mathbf{s}')$  and let  $\mathbf{u}$  be the random variable with CDF  $F_\varepsilon(\mathbf{z})$ . We have

- (i)  $\beta_{lts}(F_{(\mathbf{x}', y)}, \alpha)$  and  $\beta_{lts}(F_\varepsilon(\mathbf{z}), \alpha)$  exist.
- (ii) Furthermore, they are the solution of system of equations, respectively

$$\int (y - \mathbf{w}'\beta) \mathbf{w} \mathbf{1} (|y - \mathbf{w}'\beta - m| / \sigma \leq \alpha) dF_{(\mathbf{x}', y)}(\mathbf{x}, y) = \mathbf{0}, \tag{3.8}$$

$$\int (t - \mathbf{v}'\beta) \mathbf{v} \mathbf{1} (|(t - \mathbf{v}'\beta) - m_\varepsilon(\mathbf{z})| / \sigma_\varepsilon(\mathbf{z}) \leq \alpha) dF_{\mathbf{u}}(\mathbf{s}, t) = \mathbf{0}. \tag{3.9}$$

(iii)  $\beta_{lts}(F_{(\mathbf{x}', y)}, \alpha)$  and  $\beta_{lts}(F_\varepsilon(\mathbf{z}), \alpha)$  are unique provided that

$$\int \mathbf{w}\mathbf{w}' \mathbf{1} (|y - \mathbf{w}'\beta - m|/\sigma \leq \alpha) dF_{(\mathbf{x}', y)}(\mathbf{x}, y), \quad (3.10)$$

$$\int \mathbf{v}\mathbf{v}' \mathbf{1} (|(t - \mathbf{v}')\beta - m_\varepsilon(\mathbf{z})|/\sigma_\varepsilon(\mathbf{z}) \leq \alpha) dF_{\mathbf{u}}(\mathbf{s}, t) \quad (3.11)$$

are respectively invertible.

*Proof.* See the Appendix.  $\square$

**Theorem 3.3.** *If assumptions in theorem 3.2 hold, then for any  $\mathbf{z}_0 := (\mathbf{s}'_0, t_0)' \in \mathbb{R}^p$ , we have*

$$\dot{\beta}_{lts}(\mathbf{z}_0, F_{(\mathbf{x}', y)}) = \begin{cases} \mathbf{0}, & \text{if } t_0 - (1, \mathbf{s}'_0)\beta_{lts} \notin [m(\beta_{lts}) \pm \alpha\sigma(\beta_{lts})] \\ (t_0 - (1, \mathbf{s}'_0)\beta_{lts})M^{-1}(1, \mathbf{s}'_0)', & \text{otherwise,} \end{cases}$$

where  $\dot{\beta}_{lts}(\mathbf{z}_0, F_{(\mathbf{x}', y)})$  stands for the IF( $\mathbf{z}_0; \beta_{lts}, F_{(\mathbf{x}', y)}$ ),  $M^{-1}$  stands for the inverse of the matrix  $E(\mathbf{w}\mathbf{w}' \mathbf{1} (|r(\beta) - m(F_{r(\beta)})|/\sigma(F_{r(\beta)}) \leq \alpha))$  with  $\beta = \beta_{lts}$ , and  $[a \pm b]$  stands for  $[a - b, a + b]$ .

*Proof.* See the Appendix.  $\square$

**Remark 3.2.** See the Appendix.  $\square$

Overall, we see that the LST is globally robust with the best possible ABP of 50% and robust locally against point-mass contamination when there are vertical and bad leverage outliers.

Besides robustness, does the  $\beta_{lts}(F_{(\mathbf{x}', y)}, \alpha)$  really catch the true parameter (i.e. is it Fisher consistent)? Does the sample  $\beta_{lts}(Z^{(n)})$  converge to  $\beta_{lts}$  (or the true parameter  $\beta_0$ ) (i.e. strong or root-n consistency), and how fast does it converge? We answer these questions next.

## 4. Consistency

### 4.1. Fisher consistency

Before establishing strong or root-n consistency, we like to first show that the population version of the LST,  $\beta_{lts}(F_{(\mathbf{x}', y)}, \alpha)$ , is consistent with (or rather identical to) the true unknown parameter  $\beta_0$  under some assumptions – which is called Fisher consistency of the estimation functional. To that end, let us first recall our general model:

$$y = (1, \mathbf{x}')\beta_0 + e, \quad (4.1)$$

with its sample version given in model (1.1). In addition to the assumptions given in Theorem 3.2 for the existence and uniqueness of  $\beta_{lts}$ , we need one more assumption:

**(A3)**  $\mathbf{x}$  and  $e$  are independent and  $E_{(\mathbf{x}',y)}(e\mathbf{1}(|e - m(F_e)|/\sigma(F_e) \leq \alpha)) = 0$ . Hereafter we assume that  $m(F_e)$  and  $\sigma(F_e)$  exist uniquely.

The independence assumption between  $\mathbf{x}$  and  $e$  is typical in the traditional regression analysis. However, we can drop it here by modifying the integration appropriately (see the proof below), and it is unnecessary for  $\mathbf{x}$  to be non-random covariate (carrier). The assumption that integration equals to zero is very mild, and it automatically holds under the common assumption that the  $e$  is symmetric with respect to 0 (that is,  $e \stackrel{d}{=} -e$ ). We have

**Theorem 4.1.** *Under assumptions (A1)-(A3),  $\beta_{lst}(F_{(\mathbf{x}',y)}, \alpha) = \beta_0$  (i.e. it is Fisher consistent) provided that  $E\mathbf{w}\mathbf{w}'\mathbf{1}(|e - m(F_e)|/\sigma(F_e) \leq \alpha)$  is invertible.*

*Proof.* Notice that  $y - \mathbf{w}'\beta = \mathbf{w}'(\beta_0 - \beta) + e$ . This in conjunction with equation (3.8) yields,

$$\int (\mathbf{w}'(\beta_0 - \beta) + e)\mathbf{w}\mathbf{1}(|(\mathbf{w}'(\beta_0 - \beta) + e) - m|/\sigma \leq \alpha) dF_{(\mathbf{x}',y)} = \mathbf{0},$$

one sees that  $\beta = \beta_0$  indeed is one solution of the equation system by virtue of (A3). In light of Theorem 3.2 and the uniqueness of the solution, the desired result follows.  $\square$

#### 4.2. Strong consistency

To establish the strong consistency of  $\widehat{\beta}_{lst}(\mathbf{Z}^{(n)}, \alpha)$  for the  $\beta_{lst}(F_{(\mathbf{x}',y)}, \alpha)$ , write  $\widehat{\beta}_{lst}(F_{\mathbf{Z}}^n) := \widehat{\beta}_{lst}(\mathbf{Z}^{(n)}, \alpha)$ ,  $\beta_{lst}(F_{\mathbf{Z}}) := \beta_{lst}(F_{(\mathbf{x}',y)}, \alpha)$ ,  $Q(F_{\mathbf{Z}}^n, \beta) := Q(\mathbf{Z}^{(n)}, \beta, \alpha)$ , and  $Q(F_{\mathbf{Z}}, \beta) := Q(F_{(\mathbf{x}',y)}, \beta, \alpha)$ , for notation simplicity. where  $F_{\mathbf{Z}}^n$  is the sample version of  $F_{\mathbf{Z}} := F_{(\mathbf{x}',y)}$ , corresponding to  $\mathbf{Z}^{(n)}$  and  $\alpha$  are suppressed.

We will follow the approach in [51] and treat the problem in a more general setting. To that end, we introduce the regression depth functions  $D(F_{\mathbf{Z}}^n, \beta) = (1 + Q(F_{\mathbf{Z}}^n, \beta))^{-1}$  and  $D(F_{\mathbf{Z}}, \beta) = (1 + Q(F_{\mathbf{Z}}, \beta))^{-1}$  (see page 144 of [52] for the objective function approach). The original minimization issue becomes a maximization problem.

Let  $M_n$  be stochastic processes indexed by a metric space  $\Theta$  of  $\theta$ , and  $M: \Theta \rightarrow \mathbb{R}$  be a deterministic function of  $\theta$  which has its maximum at a point  $\theta_0$ .

The sufficient conditions for the consistency of this type of problem were given in [38] and [39], they are:

- C1:**  $\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| = o_p(1)$ ;
- C2:**  $\sup_{\{\theta: d(\theta, \theta_0) \geq \delta\}} M(\theta) < M(\theta_0)$ , for any  $\delta > 0$  and the metric  $d$  on  $\Theta$ ;

Then any sequence  $\theta_n$  is consistent for  $\theta_0$  providing that it satisfies

- C3:**  $M_n(\theta_n) \geq M_n(\theta_0) - o_p(1)$ .

**Lemma 4.1** ([38]). *If C1 and C2 hold, then any  $\theta_n$  satisfying C3 is consistent for  $\theta_0$ .*

**Remark 4.1.**

- (i) **C1** requires that the  $M_n(\boldsymbol{\theta})$  converges to  $M(\boldsymbol{\theta})$  in probability uniformly in  $\boldsymbol{\theta}$ . For the depth process  $D(F_{\mathbf{Z}}^n, \boldsymbol{\beta})$  and  $D(F_{\mathbf{Z}}, \boldsymbol{\beta})$ , it holds true (the convergence here is almost surely (a.s.) and uniformly in  $\boldsymbol{\beta}$  as shown in Lemma 4.2 below).
- (ii) **C2** essentially demands that the unique maximizer  $\boldsymbol{\theta}_0$  is well separated. This holds true for  $D(F_{\mathbf{Z}}, \boldsymbol{\beta})$  as shown in Lemma 4.3 below.
- (iii) **C3** asks that  $\boldsymbol{\theta}_n$  is very close to  $\boldsymbol{\theta}_0$  in the sense that the difference of images of the two at  $M_n$  is within  $o_p(1)$ . In [10] and [39] a stronger version of **C3** is required:

**C3\***:  $M_n(\boldsymbol{\theta}_n) \geq \sup_{\boldsymbol{\theta} \in \Theta} M_n(\boldsymbol{\theta}) - o_p(1)$ , which implies **C3**. This strong version mandates that  $\boldsymbol{\theta}_n$  nearly maximizes  $M_n(\boldsymbol{\theta})$ . Our maximum regression depth estimator  $\hat{\boldsymbol{\beta}}_{lst}(F_{\mathbf{Z}}^n, \alpha)(:= \boldsymbol{\theta}_n)$  is defined to be the maximizer of  $M_n(\boldsymbol{\theta}) := D(F_{\mathbf{Z}}^n, \boldsymbol{\beta})$ , hence **C3\*** (and thus **C3**) holds automatically.  $\square$

In light of above, we have

**Corollary 4.1.**  $\hat{\boldsymbol{\beta}}_{lst}(F_{\mathbf{Z}}^n)$  induced from  $D(F_{\mathbf{Z}}^n, \boldsymbol{\beta})$  is consistent for  $\boldsymbol{\beta}_{lst}(F_{\mathbf{Z}})$ .  $\square$

But, we can have more.

Based on the proofs of Theorems 2.2 and 3.2 and in light of Theorem 4.1, under assumptions **(A0)**–**(A3)**, we assume without loss of generality (w.l.o.g.) that  $\hat{\boldsymbol{\beta}}_{lst}(F_{\mathbf{Z}}^n) \in B(\boldsymbol{\beta}_0, r)$  and  $\boldsymbol{\beta}_{lst}(F_{\mathbf{Z}}) \in B(\boldsymbol{\beta}_0, r)$ , where  $B(\boldsymbol{\beta}_0, r)$  is a ball centered at  $\boldsymbol{\beta}_0$  with radius  $r$ . Now  $B(\boldsymbol{\beta}_0, r)$  can serve, w.l.o.g., as out parameter space  $\Theta$  of  $\boldsymbol{\beta}$  in the sequel.

**Lemma 4.2.** Under assumption **(A2)**, (a)  $\sup_{\boldsymbol{\beta} \in \Theta} |Q(F_{\mathbf{Z}}^n, \boldsymbol{\beta}) - Q(F_{\mathbf{Z}}, \boldsymbol{\beta})| = o(1)$ , a.s. and (b)  $\sup_{\boldsymbol{\beta} \in \Theta} |D(F_{\mathbf{Z}}^n, \boldsymbol{\beta}) - D(F_{\mathbf{Z}}, \boldsymbol{\beta})| = o(1)$ , a.s.

*Proof.* See the Appendix.  $\square$

**Lemma 4.3.** Assume that a regression (or location) depth function  $D(\boldsymbol{\beta}; F_{\mathbf{Z}})$  is continuous in  $\boldsymbol{\beta}$  and  $\boldsymbol{\beta} \in \Theta$  is bounded. Let  $\boldsymbol{\eta} \in \Theta$  be the unique point with  $\boldsymbol{\eta} = \arg \max_{\boldsymbol{\beta} \in \Theta} D(\boldsymbol{\beta}; F_{\mathbf{Z}})$  and  $D(\boldsymbol{\eta}; F_{\mathbf{Z}}) > 0$ . Then  $\sup_{\boldsymbol{\beta} \in N_{\varepsilon}^c(\boldsymbol{\eta})} D(\boldsymbol{\beta}; F_{\mathbf{Z}}) < D(\boldsymbol{\eta}; F_{\mathbf{Z}})$ , for any  $\varepsilon > 0$ , where  $N_{\varepsilon}^c(\boldsymbol{\eta}) = \{\boldsymbol{\beta} \in \Theta : \|\boldsymbol{\beta} - \boldsymbol{\eta}\| \geq \varepsilon\}$  and “ $A^c$ ” stands for “complement” of the set  $A$ .

*Proof.* See the Appendix.  $\square$

**Theorem 4.2.** Under assumptions **(A1)**–**(A3)**,  $\hat{\boldsymbol{\beta}}_{lst}(F_{\mathbf{Z}}^n)$  is strongly consistent for  $\boldsymbol{\beta}_{lst}(F_{\mathbf{Z}})$  (i.e.,  $\hat{\boldsymbol{\beta}}_{lst}^n - \boldsymbol{\beta}_{lst} = o(1)$  a.s.).

*Proof.* The proof for the consistency of Lemma 4.1 could easily extend to the strong consistency with a strengthened version of **C1**

**C1\***:  $\sup_{\boldsymbol{\theta} \in \Theta} |M_n(\boldsymbol{\theta}) - M(\boldsymbol{\theta})| = o(1)$ , a.s.

In the light of the proof of Lemma 4.1, we need only verify the sufficient conditions **C1\*** and **C2–C3**. By (III) of Remark 4.1, **C3** holds automatically,

so we need to verify **C1\*** and **C2**. **C1\*** follows from Lemma 4.2. So the only item left is to verify **C2** for  $D(F_Z, \beta)$  which is guaranteed by Lemma 4.3.  $\square$

**Remark 4.2.** (i) The approach utilizing a generalized Glivenko-Cantelli theorem over a class of functions with polynomial discrimination in the proof of lemma 4.2 is very powerful and applicable to many regression estimators to obtain the strong consistency result. It is certainly applicable to the least trimmed squares (LTS).

(ii) The consistency (not the strong version in Theorem 4.2) of the LTS has been obtained in [40] using standard analysis (under many assumptions on non-random  $\mathbf{x}_i$  and on the distribution of  $e$ ) which, is difficult, lengthy (an entire article in it of itself), and tedious. The approach here is different, concise and the estimator (LST) is, of course, different to the LTS.  $\square$

Consistency does not reveal the speed of convergence of sample  $\widehat{\beta}_{lst}(F_Z^n)$  to its population counterpart  $\beta_{lst}(F_Z)$ . Standard speed of  $O_p(1/\sqrt{n})$  is desirable and expected for  $\widehat{\beta}_{lst}(F_Z^n)$ . We investigate this issue next.

### 4.3. $\sqrt{n}$ -consistency

To establish the root-n consistency we need one more assumption:

(A4)  $E(e) = 0$  and  $E(\mathbf{x}\mathbf{x}')$  exists.

$E(e) = 0$  is commonly required in traditional regression analysis. The existence of covariance (and the mean) of  $\mathbf{x}$  is sufficient for the existence of  $E(\mathbf{x}\mathbf{x}')$ .

In the following, we will employ big  $O$  and little  $o$  notations for the vectors or matrices.

**Definition 4.1.** For a sequence of random vectors or matrices  $\mathbf{X}_n$ , we say

$$\begin{aligned}\mathbf{X}_n = o_p(1) &\text{ means } \|\mathbf{X}_n\| \xrightarrow{p} 0; \\ \mathbf{X}_n = O_p(1) &\text{ means } \|\mathbf{X}_n\| = O_p(1),\end{aligned}$$

where norm of a matrix  $A_{m \times n}$  is defined as  $\|A\| := \sup_{\mathbf{x} \neq 0 \in \mathbb{R}^n} \|A\mathbf{x}\|_p / \|\mathbf{x}\|_p$ ,  $p$  could be 1, 2, or  $\infty$  (see page 82 of [3]).  $\square$

**Theorem 4.3.** Assume that assumptions in Theorem 4.1 and (A4) hold, then  $\widehat{\beta}_{lst}^n - \beta_{lst} = \widehat{\beta}_{lst}^n - \beta_0 = O_p(1/\sqrt{n})$ .

*Proof.* See the Appendix.  $\square$

**Remark 4.3.** (i) The root-n consistency of an arg max estimator could be established by a general approach given in [30, 31, Theorem 1]. With the depth process introduced in the Sect. 4.2, we are unable to verify the second requirement in that theorem.

(ii) The approach here for the root-n consistency of the LST is analogous to what is given in [41] for the LTS. However, the latter is lengthy and needs a twenty-two pages article.  $\square$



## 5. Asymptotic normality

The root- $n$  consistency above could be obtained as a by-product of the asymptotic normality which will be established in the following via stochastic equicontinuity (see page 139 of [20], or the supplementary of [51]).

*Stochastic equicontinuity* refers to a sequence of stochastic processes  $\{Z_n(t) : t \in T\}$  whose shared index set  $T$  comes equipped with a semi metric  $d(\cdot, \cdot)$ .

**Definition 5.1** (IIV. 1, Def. 2 of [20]). Call  $Z_n$  stochastically equicontinuous at  $t_0$  if for each  $\eta > 0$  and  $\epsilon > 0$  there exists a neighborhood  $U$  of  $t_0$  for which

$$\limsup P \left( \sup_U |Z_n(t) - Z_n(t_0)| > \eta \right) < \epsilon. \quad (5.1)$$

□

If  $\tau_n$  is a sequence of random elements of  $T$  that converges in probability to  $t_0$ , then

$$Z_n(\tau_n) - Z_n(t_0) \rightarrow 0 \text{ in probability,} \quad (5.2)$$

because, with probability tending to one,  $\tau_n$  will belong to each  $U$ . The form above will be easier to apply, especially when behavior of a particular  $\tau_n$  sequence is under investigation.

Suppose  $\mathcal{F} = \{f(\cdot, t) : t \in T\}$ , with  $T$  a subset of  $\mathbb{R}^k$ , is a collection of real,  $P$ -integrable functions on the set  $S$  where  $P$  (probability measure) lives. Denote by  $P_n$  the empirical measure formed from  $n$  independent observations on  $P$ , and define the empirical process  $E_n$  as the signed measure  $n^{1/2}(P_n - P)$ . Define

$$\begin{aligned} F(t) &= Pf(\cdot, t), \\ F_n(t) &= P_n f(\cdot, t). \end{aligned}$$

Suppose  $f(\cdot, t)$  has a linear approximation near the  $t_0$  at which  $F(\cdot)$  takes on its minimum value:

$$f(\cdot, t) = f(\cdot, t_0) + (t - t_0)' \nabla(\cdot) + |t - t_0| r(\cdot, t). \quad (5.3)$$

For completeness set  $r(\cdot, t_0) = 0$ , where  $\nabla$  (differential operator) is a vector of  $k$  real functions on  $S$ . We cite theorem 5 of IIV.1 of [20] (page 141) for the asymptotic normality of  $\tau_n$ .

**Lemma 5.1.** *Suppose  $\{\tau_n\}$  is a sequence of random vectors converging in probability to the value  $t_0$  at which  $F(\cdot)$  has its minimum. Define  $r(\cdot, t)$  and the vector of functions  $\nabla(\cdot)$  by (5.3). If*

- (i)  $t_0$  is an interior point of the parameter set  $T$ ;
- (ii)  $F(\cdot)$  has a non-singular second derivative matrix  $V$  at  $t_0$ ;
- (iii)  $F_n(\tau_n) = o_p(n^{-1}) + \inf_t F_n(t)$ ;
- (iv) the components of  $\nabla(\cdot)$  all belong to  $\mathcal{L}^2(P)$ ;
- (v) the sequence  $\{E_n r(\cdot, t)\}$  is stochastically equicontinuous at  $t_0$ ;

then

$$n^{1/2}(\tau_n - t_0) \xrightarrow{d} \mathcal{N}(O, V^{-1}[P(\nabla\nabla') - (P\nabla)(P\nabla)']V^{-1}).$$

**Theorem 5.1.** Assume that

- (i) the uniqueness assumptions for  $\widehat{\beta}_{lst}^n$  and  $\beta_{lst}$  in theorems 2.1 and 3.2 hold;
- (ii)  $P(x_i^2)$  exists;

then

$$n^{1/2}(\widehat{\beta}_{lst}^n - \beta_{lst}) \xrightarrow{d} \mathcal{N}(O, V^{-1}[P(\nabla\nabla') - (P\nabla)(P\nabla)']V^{-1}),$$

where  $\beta$  in  $V$  and  $\nabla$  is replaced by  $\beta_{lst}$  (which could be assumed to be zero).

*Proof.* See the Appendix. □

Assume that  $\mathbf{z} = (\mathbf{x}', y)'$  follows elliptical distributions  $E(g; \boldsymbol{\mu}, \boldsymbol{\Sigma})$  with density

$$f_{\mathbf{z}}(\mathbf{x}', y) = \frac{g(((\mathbf{x}', y)' - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} ((\mathbf{x}', y)' - \boldsymbol{\mu}))}{\sqrt{\det(\boldsymbol{\Sigma})}}, \tag{5.4}$$

where  $\boldsymbol{\mu} \in \mathbb{R}^p$  and  $\boldsymbol{\Sigma}$  a positive definite matrix of size  $p$  which is proportional to the covariance matrix if the latter exists. We assume the function  $g$  to have a strictly negative derivative, so that the  $f_{\mathbf{z}}$  is unimodal.

In light of Lemma 3.1 and under some transformations (see the Appendix in the supplementary material), we can assume, w.l.o.g. that  $(\mathbf{x}', y)$  follows an  $E(g; \mathbf{0}, \mathbf{I}_{p \times p})$  distribution and  $\mathbf{I}_{p \times p}$  is the covariance matrix of  $(\mathbf{x}', y)$  hereafter.

**Corollary 5.1.** Assume that

- (i) assumptions of Theorem 5.1 hold;
- (ii)  $e \sim \mathcal{N}(0, \sigma^2)$  and  $\mathbf{x}$  are independent.

Then

- (1)  $P\nabla = \mathbf{0}$  and  $P(\nabla\nabla') = 8\sigma^2 C \mathbf{I}_{p \times p}$ , with  $C = \Gamma(1/2, 1)(\alpha c / \sigma)$ , where  $c = \sigma \Phi^{-1}(3/4)$ ,  $\Gamma(1/2, 1)(x)$  is the cumulative distribution function (CDF) of random variable  $\Gamma(a, b)$  which has a pdf:  $\frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}$ , and  $\Phi(x)$  is the CDF of  $\mathcal{N}(0, 1)$ .
- (2)  $\mathbf{V} = 2C_1 \mathbf{I}_{p \times p}$  with  $C_1 = 2 * \Phi(\alpha c / \sigma) - 1$ .
- (3)  $n^{1/2}(\widehat{\beta}_{lst}^n - \beta_{lst}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \frac{2C\sigma^2}{C_1^2} \mathbf{I}_{p \times p})$ .

*Proof.* By Theorem 4.1 and Lemma 3.1, we can assume, w.l.o.g., that  $\beta_{lst} = \beta_0 = \mathbf{0}$ . Utilizing the independence between  $e$  and  $\mathbf{x}$  and Theorem 5.1, a straightforward calculation leads to the results. □

## 6. Computation

Now we address one of the most important topics on robust regression estimation, that is, the computation of the estimator. Unlike the LS estimator, which

has an analytical formula for the computation, for the LST estimator, we do not have such a formula. The formula given in (2.9) can not serve our purpose (due to the circular dependency: the RHS depends on the LHS). For small sample size  $n$  and dimension  $p$ , one can compute the LST exactly (the  $L$  in Theorem 2.1 is not a big number), but that is not affordable for large  $n$  and  $p$ . In general, we use approximate algorithms (AAs).

### 6.1. A procedure based Theorem 2.1

In light of Theorem 2.1, if we discover all  $R_{\beta^k}$ s for  $1 \leq k \leq L$ , then we can get the exact result. But in practice, this might not be computationally affordable. However, we can simply search as many  $R_{\beta^k}$ s as possible to get a good approximation of the estimate  $\widehat{\beta}_{l_{st}}^n$ .

To identify  $R_{\beta^k}$  is equivalent to identifying  $i_1, \dots, i_K$  so that  $D_{i_1} < D_{i_2} < \dots, D_{i_K}$  in light to (2.7), where  $K = |I(\beta^k)|$ . The latter is equivalent to finding a  $\beta \in R_{\beta^k}$ , then we get the desired  $i_1, \dots, i_K$ . To find the desired  $\beta$ , one way is to find a  $\bar{\beta}$  on the common boundary of  $R_{\beta^k}$  and  $R_{\beta^l}$  so that there are  $i \neq j$ ,  $D_i = D_j$  for some  $1 \leq l \neq k \leq L$  and  $1 \leq i, j \leq n$ . A small perturbation of the coordinates of the  $\bar{\beta} = (\beta_1, \dots, \beta_p)'$  leads to more than one  $\beta$ s ( $\beta = (\beta_1, \dots, \beta_j \pm \delta, \dots, \beta_p)'$  (for some  $1 \leq j \leq p$  and  $\delta > 0$ ) that belong to  $R_{\beta^k}$  or  $R_{\beta^l}$ .

Now we address the way to find out  $\bar{\beta}$ . In light of (2.7), there are  $i \neq j$ ,  $D_i = D_j$  for some  $1 \leq i, j \leq n$ . The equality  $D_i = D_j$  implies that (i)  $r_i = r_j$  or (ii)  $(r_i + r_j)/2 = m_n(\beta)$ . Both equalities could lead to some  $\bar{\beta}$ s, but the first equality  $r_i = r_j$  is more convenient.

We focus on the first equality which amounts to  $y_i - y_j = (\mathbf{w}_i - \mathbf{w}_j)'\beta = (\mathbf{x}_i - \mathbf{x}_j)'(\beta_2, \dots, \beta_p)'$ , where  $\mathbf{w}' = (1, \mathbf{x}')$ ,  $\beta = (\beta_1, \dots, \beta_p)'$ . Assume that  $\mathbf{x}_i \neq \mathbf{x}_j$  for  $i \neq j$ . If  $y_i = y_j$ , then,  $\beta = (\beta_1, \mathbf{0}'_{p-1})'$  is one of solutions, otherwise, from this equation, we see that (i)  $\beta_1$  could be any number in  $\mathbb{R}^1$ , (ii) the equation defines a  $(p-1)$ -dimensional hyperplane. Consequently, all  $\beta = (\beta_1, 0, \dots, 0, \frac{y_i - y_j}{x_{ik} - x_{jk}}, 0, \dots, 0) \in \mathbb{R}^p$  are solutions, where  $\beta_1 \in \mathbb{R}^1$  and  $x_{ik} \neq x_{jk}$ ,  $1 \leq k \leq (p-1)$ . Simple choices for  $\beta_1$  could be 0 and 1 or any constant. From here we obtain at least two  $\beta$ s that lie on the common boundary.

With the small perturbation ( $\pm\delta$ ) to the  $i$ th coordinate of the  $\beta$ s above we could obtain  $4p$  new  $\beta$ s. For each such  $\beta$ , we first obtain  $i_1, \dots, i_K$  with  $K = |I(\beta)|$  and then check if the strict inequalities in (2.7) hold.

If they do not hold, then move to the next  $\beta$ . Otherwise, check if the  $K$  indices already appear before, if it has, then do nothing, else update the data structure that stores the indices, and obtain the least square solution  $\beta_{l_s}$ -new based on the sub-data set with the  $K$  subscripts ( $I(\beta)$ ) and the sum of squared residuals. If the latter is smaller than SS-min, then set it to be the SS-min and update  $\widehat{\beta}_{l_{st}}^n$  with  $\beta_{l_s}$ -new. Increase  $T_{l_s}$ , which is the counter for the number of LS calculations, by one. Move to the next  $\beta$  until all  $4p$   $\beta$ s are exhausted. Then repeat the entire process with a new pair  $(i, j)$ . Summarizing discussions so far, we have

**AA1 – pseudocode for computing the LST based on Theorem 2.1**

**Input:** A data set  $\mathbf{Z}^{(n)} = \{(\mathbf{x}'_i, y_i)', i = 1, 2, \dots, n\}$ , a fixed  $\alpha$ . Assume that  $\mathbf{x}_i \neq \mathbf{x}_j$  if  $i \neq j$ .

- (1) Sample two indices  $i$  and  $j$  from  $\{1, \dots, n\}$ , assume that  $x_{ik} \neq x_{jk}$ ,  $1 \leq k \leq (p-1)$  (i.e. the  $k$ th coordinates of  $\mathbf{x}_i$  and  $\mathbf{x}_j$  do not equal). Consider

$$\boldsymbol{\beta}^0 = (0, 0, \dots, 0, b_{k+1}, 0, \dots, 0)', \boldsymbol{\beta}^1 = (1, 0, \dots, 0, b_{k+1}, 0, \dots, 0)' \text{ in } \mathbb{R}^p$$

Both have the same  $(k+1)$ th coordinate,  $b_{k+1} := (y_i - y_j)/(x_{ik} - x_{jk})$ .

- (2) Write  $\boldsymbol{\beta}^j(l, \pm\delta)$  for the perturbed  $\boldsymbol{\beta}^j$  with its  $l$ th coordinate adding or subtracting a  $\delta > 0$ . Define a set

$$S(\boldsymbol{\beta}) = \bigcup_{l=1}^p \{\boldsymbol{\beta}^0(l, \pm\delta)\} \bigcup_{l=1}^p \{\boldsymbol{\beta}^1(l, \pm\delta)\}.$$

- (3) For each  $\boldsymbol{\beta}$  of  $4p$   $\boldsymbol{\beta}$ s in the set  $S(\boldsymbol{\beta})$ ,
- (a) obtain  $i_1, \dots, i_K$  with  $K = |I(\boldsymbol{\beta})|$  and check to see if the strict inequalities in (2.7) hold.
    - (a1) If not, move to the next  $\boldsymbol{\beta}$ ; else
    - (a2) check if the  $K$  indices already appear in a structure  $S_{ind}$ 
      - (i) if yes, then move to the next  $\boldsymbol{\beta}$ ; else
      - (ii) update  $S_{ind}$  by storing the  $K$  indices in the structure  $S_{ind}$  and calculate the LS estimate  $\boldsymbol{\beta}_{ls}$ -new based on the sub-data set with index in  $I(\boldsymbol{\beta})$  and obtain the sum of  $|I(\boldsymbol{\beta})|$  squared residuals,  $SS(\boldsymbol{\beta}_{ls}$ -new).
      - (iii) Update  $SS_{\min}$  if it is greater than  $SS(\boldsymbol{\beta}_{ls}$ -new) and update  $\hat{\boldsymbol{\beta}}_{lst}^n$  with  $\boldsymbol{\beta}_{ls}$ -new. Update the counter for the total number  $T_{ls}$  of LS calculations, if the latter is less than  $N$  (the total number of LS calculations decided to perform), then continue the loop (go to (3)), else break the stop.
  - (b) If  $T_{ls} < N$ , then go to (1), else break the loop.

**Output:**  $\hat{\boldsymbol{\beta}}_{lst}^n$

**Remark 6.1.** See the Appendix. □

**6.2. A subsampling procedure**

Many robust regression estimators use subsampling procedures in practice (see [23, 11, 12, 24, 43, 25, 26, 50, 54], among others).

The basic idea is straightforward: (1) draw a sub-sample of size  $m$  from data set  $\mathbf{Z}^{(n)} = \{(\mathbf{x}'_i, y_i)' \in \mathbb{R}^p, \mathbf{x}_i \in \mathbb{R}^{p-1}, i \in \{1, 2, \dots, n\}\}$ . (2) compute an estimate based on the sub-sample and obtain the objective function value. (3) if

the objective function value can be further improved (reduced), then go to (1), otherwise, stop and output the final step estimate.

Natural questions for the above procedure include (1) how do we guarantee the convergence of the procedure and the final answer is the global minimum? (2) what is the exact size  $m$  and what is the relationship with  $n$  and dimension  $p$ ? To better address these matters, we first propose the corresponding procedure for our LST.

### AA2 pseudocode for a sub-sampling procedure for LST

**Input:** A data set  $\mathbf{Z}^{(n)} = \{\mathbf{Z}_1, \dots, \mathbf{Z}_n\} = \{(\mathbf{x}'_i, y_i)', i = 1, 2, \dots, n\} \in \mathbb{R}^p$  (assume that  $p \geq 2$ ) and an  $\alpha \geq 1$  (default is one).

(a) **Initialization:**  $N = \min\{\binom{n}{p}, 300(p-1)\}$ ,  $R=0$ ,  $Q_{old} = 10^8$ ,  $\beta_{old} = \mathbf{0}$  (or an LS (or LTS) estimate).

(b) **Iteration:** while ( $R \leq N$ )

keep sampling  $p$  indices  $\{i_1, \dots, i_p\}$  from  $\{1, 2, \dots, n\}$  (without replacement) until  $M'_x := (\mathbf{w}_{i_1}, \dots, \mathbf{w}_{i_p})$  becomes invertible. Let  $\beta_{new} = (M_x)^{-1}(y_{i_1}, \dots, y_{i_p})'$ .

(1) Calculate  $I(\beta_{new})$  (based on (2.6)) and  $Q_{new} := Q^n(\beta_{new})$  (based on (2.4)).

(2) \* If  $Q_{new} < Q_{old}$ , then  $Q_{old} = Q_{new}$ ,  $\beta_{old} = \beta_{new}$ . Get an LS estimator  $\beta_{ls}$  based on the data points of  $\mathbf{Z}^{(n)}$  with subscripts from  $I(\beta_{new})$ . Go to (1) with  $\beta_{new} = \beta_{ls}$ .

\* Else if  $Q_{new} = Q_{old}$  break  
else  $R=R+1$ , go to (b)

**Output:**  $\beta_{new}$ . □

**Remark 6.2.** See the Appendix. □

## 7. Examples and comparison

This section investigates the performance of the AAs and compares it with that of the benchmark LTS. First, we like to give some guidance for selection among the two AAs.

**Example 7.1** (Performance of the two AAs). There are two AAs, so which of them should we recommend? This example tries to answer this by examining the speed and accuracy of the two AAs.

We generate 1000 samples  $\mathbf{Z}^{(n)} = \{(\mathbf{x}'_i, y_i)', i \in \{1, \dots, n\}, \mathbf{x}_i \in \mathbb{R}^{p-1}\}$  from the standard Gaussian distribution for various sample size  $n$  and dimension  $p$ . For the speed, we calculate the *total time* consumed for all 1000 samples (dividing it by 1000, we get the average time consumed per sample) by different AAs. For accuracy (or variance, or efficiency), we will compute their empirical mean squared error (EMSE).

TABLE 1

Total computation time for all 1000 samples (seconds) and empirical mean squared error (EMSE) of different AAs for various  $ns$  and  $ps$ .

Table entries (a, b) are: a := empirical mean squared error, b := total time consumed

$n$	$p$	AA1	AA2
50	3	(0.3499, 566.49)	(0.5290, 651.25)
	5	(0.5817, 457.49)	(0.7645, 861.75)
	10	(0.5390, 682.41)	(1.7177, 1016.6)
100	3	(0.1755, 573.07)	(0.3619, 879.01)
	5	(0.2023, 638.76)	(0.4528, 1042.6)
	10	(0.2576, 702.02)	(0.7000, 1071.5)
200	3	(0.0825, 619.75)	(0.3025, 1309.7)
	5	(0.1055, 676.63)	(0.3501, 1285.6)
	10	(0.1283, 698.14)	(0.4178, 1310.2)

For a general estimator  $\mathbf{T}$ , if it is regression equivariant, then we can assume w.l.o.g. that the true parameter  $\beta_0 = \mathbf{0} \in \mathbb{R}^p$ . We calculate  $\text{EMSE} := \sum_{i=1}^R \|\mathbf{T}_i - \beta_0\|^2 / R$ , the empirical mean squared error (EMSE) for  $\mathbf{T}$ , where  $R = 1000$ ,  $\beta_0 = (0, \dots, 0)' \in \mathbb{R}^p$ , and  $\mathbf{T}_i$  is the realization of  $\mathbf{T}$  obtained from the  $i$ th sample with size  $n$  and dimension  $p$ . The EMSE and the total time consumed (in seconds) by different AAs are listed in Table 1.

Inspecting Table 1 immediately reveals that (i) AA2 is not only the slowest but is the most inaccurate (with the largest EMSEs) in all cases considered. (ii) AA1 has both speed and accuracy advantages for all cases considered.

Overall, we recommend AA1. That does not exclude the potential of improvement of AA2 via the idea in [26].  $\square$

All R code for simulation and examples as well as figures in this article (downloadable via <https://github.com/left-github-4-codes/LST>) were run on a desktop Intel(R)Core(TM) 21 i7-2600 CPU @ 3.40 GHz.

The data points in the example above are standard normal and hence not realistic. In the following, we will investigate the performance of AA1 versus the LTS for contaminated standard normal data sets and for moderate as well as large  $ns$  and  $ps$ .

**Example 7.2** (Multiple regression with contaminated normal data sets). Now we consider data with contamination, which is typical for big data sets in the big-data era.

We consider the contaminated highly correlated normal data points scheme. We generate 1000 samples  $\mathbf{Z}_i = (\mathbf{x}_i', y_i)'$  with various  $ns$  from the normal distribution  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\mu}$  is a zero-vector in  $\mathbb{R}^p$ , and  $\boldsymbol{\Sigma}$  is a  $p$  by  $p$  matrix with diagonal entries being 1 and off-diagonal entries being 0.9. Then  $\varepsilon\%$  of them are contaminated by normal points with  $\boldsymbol{\mu}$  being the  $p$ -vector with all elements being 7 except the last one being  $-2$  and the covariance matrix being diagonal with diagonal being 0.1 and off-diagonal being zero. The results are listed in Table 2.

Inspecting the table reveals that (i) in terms of EMSE, AA1 is the overall

TABLE 2

Total computation time for all 1000 samples (seconds) and empirical mean squared error (EMSE) of the LST(AA1) versus the LTS(ltsReg) for various  $n$ ,  $p$ , and contaminations.

Normal data sets, each with  $\varepsilon\%$  contamination

Table entries (a, b) are: a := empirical mean squared error, b := total time consumed

$p$	$n$	$\varepsilon = 5\%$		$\varepsilon = 10\%$	
		AA1	ltsReg	AA1	ltsReg
5	100	(0.2971, 9.6581)	(0.3010, 22.867)	(0.2843, 494.01)	(0.2942, 25.289)
	200	(0.2503, 26.045)	(0.2650, 41.861)	(0.2517, 26.629)	(0.2630, 43.504)
	300	(0.2396, 54.100)	(0.2551, 63.639)	(0.2366, 54.885)	(0.2534, 63.522)
10	400	(0.1335, 1085.6)	(0.1394, 181.18)	(0.1340, 1056.2)	(0.1382, 175.92)
	500	(0.1280, 1207.7)	(0.1321, 222.81)	(0.1289, 1178.5)	(0.1321, 218.94)
	600	(0.1247, 1308.4)	(0.1285, 152.47)	(0.1253, 1273.6)	(0.1276, 149.99)
20	700	(0.0815, 2044.9)	(0.0885, 549.61)	(0.0838, 1994.0)	(0.0882, 547.53)
	800	(0.0776, 2261.7)	(0.0837, 620.63)	(0.0796, 2177.0)	(0.0837, 616.87)
	900	(0.0748, 2436.1)	(0.0804, 541.20)	(0.0761, 2353.7)	(0.0795, 538.43)
40		$\varepsilon = 30\%$		$\varepsilon = 40\%$	
	300	(0.4347, 53.248)	(1.9236, 1635.1)	(0.4352, 56.430)	(1.3517, 1712.8)
	400	(0.3362, 100.04)	(1.2604, 2401.5)	(0.3314, 102.81)	(0.8995, 2399.5)
50	500	(0.2594, 147.66)	(0.9514, 2963.4)	(0.2873, 146.67)	(0.6851, 2787.7)
	300	(0.5242, 58.736)	(2.7826, 2861.8)	(0.5700, 59.903)	(1.9808, 2896.3)
	400	(0.4085, 89.897)	(1.7562, 3292.0)	(0.4539, 108.88)	(1.2547, 3925.5)
	500	(0.3107, 145.84)	(1.2870, 4510.5)	(0.3406, 145.75)	(0.9086, 4419.6)

winner (with the smallest EMSE in all cases considered), the LTS has the largest EMSE in all the cases; (ii) in terms of speed, the LTS (or rather ltsReg) is the winner when  $p = 10$  or  $20$ . The AA1 is the winner for all other  $p$ 's, except when  $p = 5$ ,  $n = 100$  and  $\varepsilon = 10\%$ . For the latter case, AA1 can still be faster by tuning  $T_{ls}$  to be 1, then we get (0.2986, 10.396) for AA1 versus (0.2948, 23.133) for ltsReg (suffering a slight increase in EMSE).  $\square$

The LTS (or ltsReg) demonstrates its well-known speedy advantage, which is partially due to its background computation via Fortran subroutine and the computation scheme proposed in [26]. The AA1 (a pure R programming procedure), on the other hand, has the potential to speed up via Rcpp or even Fortran.

**Remark 7.1. (i) Parameters tuning.** Two parameters in AA1 that can be tuned. The  $T_{ls}$  is set to 300 for better EMSE (as in the  $p = 5$ ,  $n = 100$ , and  $\varepsilon = 10\%$  case). If tuning it to be 1, we get a much faster AA1 (as in the cases  $p = 30$ ,  $40$ , and  $p = 5$ , except when  $n = 100$ , and  $\varepsilon = 10\%$ ). For the  $\alpha$  in the definition of the LST, it is set to 1 (default value) in Table 1, it is set to 3 as in Table 2 when there are contaminations (or outliers). Note that theoretically speaking, both the LST and the LTS can resist 50% contamination without breakdown. So the 40% contamination rate in Table 2 is relevant and is also employed in [26].

(ii) The LTS estimate is obtained via the R package ltsReg,  $h$  is the default value  $\lfloor (n + p + 1)/2 \rfloor$ , we could tune this  $h$  to get better performance from the

LTS. However, this will decrease LTS’s finite sample breakdown value. This is not the case for LST with the  $\alpha$  (see Theorem 3.1).  $\square$

So far we have assumed that the true  $\beta_0$  is the zero vector based on the regression equivariance. One might not be used to this assumption.

**Example 7.3** (Performance of the LST and the LTS with respect to a given  $\beta_0$ ). Now we examine the performance of three regression estimators the LST, the LTS, and LMS in a slightly different setting. We generate 1000 samples  $\{(\mathbf{x}'_i, y_i) \in \mathbb{R}^p\}$  with a fixed sample size 100 from an assumed model:  $y_i = \beta_0' \mathbf{x}_i + e_i$ , where  $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ip-1})'$  and  $\beta_0 = (\beta_0, \dots, \beta_{p-1})'$  are in  $\mathbb{R}^p$  and  $x_{ij}$  and  $e_i$  are independently from either the Cauchy or  $\mathcal{N}(0, 1)$  distribution.

We list the total time consumed (in seconds) and the EMSE (the same formula as before but the true  $\beta_0$  is the given one no longer being the zero vector) for the three methods with respect to different  $\beta_0$ ’s in Table 3.

**Case I**  $\beta_0 = (-2, 0.1, 1)'$ , all  $x_{ij}$  and  $e_i$  are from  $\mathcal{N}(0, 1)$  distribution.

**Case II**  $\beta_0 = (-2, 0.1, 1, 5)'$ ,  $x_{i1}$ ,  $x_{i2}$ , and  $e_i$  are from  $\mathcal{N}(0, 1)$  and  $x_{i3}$  is from Cauchy distribution.

**Case III**  $\beta_0 = (50, 0.1, -2, 15, 100)'$ , all  $x_{ij}$  and  $e_i$  are from  $\mathcal{N}(0, 1)$ .

TABLE 3  
Performance of the LST, the LTS, and the LMS for three true  $\beta_0$ ’s.

Replication 1000 times, $n = 100$			
Performance criteria	LST(AA1)	LMS(lmsreg)	LTS(ltsReg)
<b>Case I</b>			
$p = 3$			
EMSE	3.525451	4.204053	3.806951
Total time consumed	11.53858	10.49865	17.81713
<b>Case II</b>			
$p = 4$			
EMSE	29.91539	30.23814	29.97682
Total time consumed	9.919189	6.087584	10.31606
<b>Case III</b>			
$p = 5$			
EMSE	12724.32	12726.87	12724.74
Total time consumed	14.54680	17.42145	22.08751

Inspecting the Table reveals that (i) in terms of EMSE, the LST (AA1) is the overall winner (has the smallest EMSE in all cases) whereas the LMS is the loser; (ii) in terms of speed, there is no overall winner. In two respective cases, the LMS is the fastest whereas the LST is fastest in  $p = 5$  case and the LTS is the slowest in all cases.  $\square$

Up to this point, we have dealt with synthetic data sets. Next we examine the performance of the LST, the LTS and the LMS with respect to real data sets in high dimensions.



**Example 7.4** (Textbook size real data sets). We first look at real data sets with relatively small sample size  $n$  and moderate dimension  $p$ . For a description of data sets, see [23]. Since some of methods might depend on randomness, So we run the computation  $R = 1000$  times to alleviate the randomness. We then calculate the *total* time consumed (in seconds) by different methods for all replications, and the EMSE (with true  $\beta_0$  being replaced by the sample mean of 1000  $\hat{\beta}$ s), which is the sample variance of all  $\hat{\beta}$ s up to a factor 1000/999. The results are reported in Table 4, where the parameters  $\alpha$  and  $T_{l_s}$  in AA1 are tuned.

TABLE 4  
Total time consumed (in seconds) and sample variance in 1000 replications by the LST (AA1), the LTS (ltsReg), and the LMS (lmsreg) for various real data sets.

Table entries (a, b) are: a := empirical variance of  $\hat{\beta}$ s, b := total time consumed

data set	$(n, p)$	AA1	ltsReg	lmsreg
salinity	(28, 4)	(0.0, 2.3290)	(0.0, 8.8385)	(1.3719, 4.9425)
phosphor	(18, 3)	(0.0, 4.9218)	(0.0, 8.3902)	(0.0000, 1.5153)
wood	(20, 6)	(0.0, 4.8013)	(0.0, 10.343)	(2.6470, 8.3714)
coleman	(20, 6)	(0.0, 14.585)	(0.0, 10.159)	(243.11, 8.3560)

Inspecting the Table reveals that (i) in terms of the EMSE (or rather empirical variance), AA1 and ltsReg are the overall winners for all cases considered (no randomness) and LMS has the largest sample variance. (ii) in terms of computation speed, there is no overall winner, but AA1 is faster than ltsReg in three out of four cases. The LMS is the fastest in one case.  $\square$

The limitation of this example is that the data sets are still relatively small and not in high dimensions. We examine a high dimension and large sample dataset next.

**Example 7.5** (A large real data set). Boston housing is a famous data set [9] and studied by many authors with different emphasizes (transformation, quantile, nonparametric regression, etc.) in the literature. For a more detailed description of the data set, see <http://lib.stat.cmu.edu/datasets/>.

The analysis reported here did not include any of the previous results, but consisted of just a straight linear regression of the dependent variable (median price of a house) on the thirteen explanatory variables as might be used in an initial exploratory analysis of a data set. We have sample size  $n = 506$  and dimension  $p = 14$ .

We assess the performance of the LST, the LTS, and the LMS as follows: (i) we sample  $m$  points (without replacement) ( $m = 506$ , entire data set, or  $m = 200, 250, 300, 350$ ) from the entire data set, and compute the  $\hat{\beta}$ s with different methods, we do this RepN times, where replication number RepN varies with respect to different  $m$ s. (ii) we calculate the total time consumed (in seconds) by different methods for all replications, and the EMSE (with true  $\beta_0$  being replaced by the sample mean of RepN  $\hat{\beta}$ s from (i)), which is the sample variance of all  $\hat{\beta}$ s up to a factor  $RepN/(RepN - 1)$ . The results are reported in Table 5.

TABLE 5

Total time consumed (in seconds) and sample variance in  $RepN$  replications by the LTS (*ltsReg*), the LST (*AA1*), and the LMS (*lmsreg*) for real data sets with various sample size  $m$ 's and  $p = 14$ .

Table entries (a, b) are: a := empirical variance of  $\hat{\beta}$ s, b := total time consumed

$p$	$m$	RepN	LST(AA1)	LTS( <i>ltsReg</i> )	LMS( <i>lmsreg</i> )
	200	$10^4$	(195.3379, 595.7677)	(220.8644, 480.0612)	(847.2457, 472.4671)
	250	$10^4$	(164.4042, 723.5861)	(169.5725, 597.2802)	(791.2557, 555.3318)
14	300	$10^4$	(461.5653, 514.8522)	(126.7703, 683.3362)	(754.2416, 623.5828)
	350	$10^4$	(453.3266, 695.9286)	(97.86377, 821.1486)	(724.2104, 732.2517)
	506	$10^3$	(0.000000, 142.4225)	(42.58697, 116.5830)	(703.7999, 101.0454)

Inspecting the Table reveals that (i) the LMS has the largest EMSEs while being faster than the LTS in all cases; (ii) the LST has smallest EMSE in three cases among the five (in those cases it is slower than the LTS) (in the other two cases the LTS is slower); (iii) in the entire data-set case, the LST returned the same estimate every replication whereas the LTS and the LMS did not.

## 8. Final discussions

**The difference between the LTS and the LST** The least sum of squares of trimmed (LST) residuals estimator has the proven best asymptotic breakdown point of 50% and is another robust alternative to the classical least sum of squares (LS) of residuals estimator. The latter keeps all squared residuals whereas the former trims some residuals and then squares what is left. Trimming is also utilized in the popular least sum of trimmed squares (LTS) of the residuals estimator. However, the two trimming schemes are quite different, the one used in the LTS is a one-sided trim (only large squared residuals are trimmed, of course, it also might be regarded as a two-sided trim with respect to the un-squared residuals) whereas the one utilized in the LST is a depth-based (or outlyingness-based) trim (see [49] and [45] for more discussions on trimming schemes) which can trim both ends of un-squared residuals and does not trim a fixed number of residuals.

Besides the trimming scheme difference, there is another difference between the LTS and the LST, that is, the order of trimming and squaring. In the LTS, squaring is first, followed by trimming whereas, in the LST, the order is reversed. All these difference leads to an unexpected performance contrast between in the LTS and the LST as demonstrated in the last section.

**Fairness of performance criteria** For comparison of the performance between the LST and the LTS, we have focused on the variance (accuracy, efficiency, or EMSE) and the computation speed of the algorithms for the estimators. The asymptotic efficiency (AE) of the LTS has been reported to be just 7% in [34] or 8% in [17] (page 132), the AE of the LST is yet to be discovered; however, it is expected to be better than 8%. This assertion is supported by

the experimental results in the last section (Tables 2, and 3 indicate that the LST is more efficient than the LTS). Furthermore, this was also supported by the results of [45] for various trimming schemes in the case of  $p = 1$ .

The computation speed comparison of the LTS versus the LST in the last section is somewhat biased in favor of the LTS. It is essentially a speed comparison of pure R versus R plus Fortran since the Fortran subroutine (rfltsreg) is called in ltsReg (similarly lmsreg also call a Fortran subroutine). Even with that, ltsReg does not have an overwhelming advantage in speed over AA1. However, there is still room for improvement in AA1 by utilizing Fortran or better Rcpp.

**Parameters tuning and finite sample breakdown point** There are two parameters  $h$  in the LTS and  $\alpha$  in the LST which can be tuned in the program for computation. Their values have a connection with the finite sample breakdown point. For example, when  $h$  takes its default value  $\lfloor (n + p + 1)/2 \rfloor$ , then the FSBP of the LTS is  $(n - h + 1)/n$  which will decrease from the best FSBP result  $(\lfloor (n - p)/2 \rfloor + 1)/n$  (see pages 125, 132 of [23]) when  $h$  increases. For the parameter  $\alpha$  in LST, as long as  $\alpha \geq 1$  then the high FSBP in theorem 3.1 remains valid. This is due to the difference in the trimming schemes (see [45]).

**Open and future problems** By simply switching the order of trimming and squaring and adopting a depth-based trimming scheme, the LTS and the LST can have different performances. One naturally wonders what if one does the same thing with respect to the famous the LMS introduced also by [21] (i.e. the least square of the median (LSM) of residuals estimator). It turns out, this is not a good idea since there is a universal solution, it is  $\hat{\beta} = (\text{Med}\{y_i\}, 0, \dots, 0) \in R^p$ .

One interesting problem that remains is the investigation of the least sum of squares of trimmed residuals with yet another trimming scheme such as the winsorized version given in [45], that is, replacing the residuals beyond the cutoff values at the two ends with just the cutoff values or even a more generalized weighted (trimming) scheme which includes the hard 0 and 1 trimming scheme. Other challenging open topics that deserve to be pursued elsewhere include (i) providing a finite sample estimation error analysis (non-asymptotic analysis) and (ii) regularized regression based on the LST to handle variable selection and model interpretation issues when dimension  $p$  is much larger than sample size  $n$ .

## Acknowledgments

The authors thank Denis Selyuzhitsky, Nadav Langberg, and Professors Wei Shao and Yimin Xiao for their stimulating discussions and the authors thank the Co-Editor-in-Chiefs, Professors Grace Y. Yi and Gang Li and the anonymous referees for their insightful and constructive comments. All of this feedback has significantly improved the manuscript.

## Funding

Authors declare that there is no funding received for this study.

## Conflicts of interests/Competing interests

Authors declare that there is no conflict of interests/Competing interests.

## Supplementary Material

### Supplementary to “Least sum of squares of trimmed residuals regression”

(doi: [10.1214/23-EJS2164SUPP](https://doi.org/10.1214/23-EJS2164SUPP); .pdf).

## References

- [1] Anonymous (1821), “Dissertation sur la recherche du milieu le plus probable, entre les rbsultats de plusieurs observations ou experiences”, *Ann. Math. Pures Appl.*, 12, 181–204. [MR1556114](#)
- [2] Bickel, P. J. (1975), “One-step Huber estimates in the linear model”, *J. Am. Statist. Assoc.*, 70, 428–434. [MR0386168](#)
- [3] Boyd, S. and Vandenberghe, L. (2004), “*Convex Optimization*”, Cambridge University Press. [MR2061575](#)
- [4] Dixon, W. J. and Tukey, J. W. (1968), “Approximate behavior of the distribution of Winsorized  $t$  (Trimming/Winsorization 2)”, *Technometrics*, 10(1), 83–98. [MR0221700](#)
- [5] Donoho, D. L. “Breakdown properties of multivariate location estimators”. PhD Qualifying paper, Harvard Univ. (1982).
- [6] Donoho, D. L., and Gasko, M. (1992), “Breakdown properties of multivariate location parameters and dispersion matrices”, *Ann. Statist.*, 20, 1803–1827. [MR1193313](#)
- [7] Donoho, D. L., and Huber, P. J. (1983), “The notion of breakdown point”, in: P. J. Bickel, K. A. Doksum and J. L. Hodges, Jr., eds. *A Festschrift for Erich L. Lehmann* (Wadsworth, Belmont, CA), pp. 157–184. [MR0689745](#)
- [8] Hampel, F. R., Ronchetti, E. M., Rousseeuw, P. J., and Stahel, W. A. (1986), “*Robust Statistics: The Approach Based on Influence Functions*”, John Wiley & Sons, New York. [MR0829458](#)
- [9] Harrison, D. and Rubinfeld, D.L. (1987), “Hedonic prices and the demand for clean air”, *J. Environ. Economics and Management*, 5, 81–102.
- [10] Kim, J. and Pollard, D. (1990), “Cube root asymptotics”. *Ann. Statist.*, 18, 191–219. [MR1041391](#)
- [11] Hawkins, D. M. (1994), “The feasible solution algorithm for least trimmed squares regression”, *Computational Statistics & Data Analysis*, 17, 185–196.
- [12] Hawkins, D. M. and Olive, D. J. (1999), “Improved feasible solution algorithms for high breakdown estimation”, *Computational Statistics & Data Analysis*, 30(1), 1–11. [MR1681450](#)

- [13] Hettmansperger, T. P. and Sheather, S. J. (1992), “A cautionary note on the method of least median squares”, *The American Statistician*, 46(2), 79–83. [MR1165565](#)
- [14] Huber, P. J. (1964), “Robust estimation of a location parameter”, *Ann. Math. Statist.*, 35, 73–101. [MR0161415](#)
- [15] Huber, P. J. (1973), “Robust regression,” *Ann. Statist.*, 1, 799–821. [MR0356373](#)
- [16] Johansen, S., and Nielsen, B. (2013), “Outlier detection in regression using an iterated one-step approximation to the Huber-Skip estimator”, *Econometrics*, 1, 53–70.
- [17] Maronna, R. A., Martin, R. D., and Yohai, V. J. (2006), “*Robust Statistics: Theory and Methods*”, John Wiley & Sons. [MR2238141](#)
- [18] Mendeleev, D. I. (1895), “Course of work on the renewal of prototypes or standard measures of lengths and weights (Russian)”, *Vremennik Glavnoi Palaty Mer i Vesu* 2, 157–185. Reprinted 1950: *Collected Writings (Socheneniya)*, Izdat. Akad. Nauk, SSSR, Leningrad-Moscow, Vol. 22, pp. 175–213.
- [19] Öllerer, V., Croux, C., and Alfons, A. (2015), “The influence function of penalized regression estimators”, *Statistics*, 49(4), 741–765. [MR3367721](#)
- [20] Pollard, D. (1984), “*Convergence of Stochastic Processes*”, Springer, Berlin. [MR0762984](#)
- [21] Rousseeuw, P. J. (1984), “Least median of squares regression”, *J. Amer. Statist. Assoc.*, 79, 871–880. [MR0770281](#)
- [22] Rousseeuw, P. J., and Hubert, M. (1999), “Regression depth (with discussion)”, *J. Amer. Statist. Assoc.*, 94, 388–433. [MR1702314](#)
- [23] Rousseeuw, P. J., and Leroy, A. (1987), “*Robust Regression and Outlier Detection*”, Wiley New York. [MR0914792](#)
- [24] Rousseeuw, P. J., Struyf, A. (1998), “Computing location depth and regression depth in higher dimensions”, *Statistics and Computing*, 8, 193–203. [MR1702314](#)
- [25] Rousseeuw, P. J. and Van Driessen, K. (1999), “A fast algorithm for the minimum covariance determinant estimator”, *Technometrics*, 41(3), 212–223.
- [26] Rousseeuw, P. J. and Van Driessen, K. (2006), “Computing LTS regression for large data sets”, *Data Mining and Knowledge Discovery*, 12, 29–45. [MR2225526](#)
- [27] Rousseeuw, P. J. and Yohai, V. J. (1984), “*Robust regression by means of S-estimators*”, in: *Robust and Nonlinear Time Series Analysis*, Lecture Notes in Statist., vol. 26, pp. 256–272, Springer, New York. [MR0786313](#)
- [28] Ruppert, D. and Carroll, R. J. (1980), “Trimmed least squares estimation in the linear model”, *J. Amer. Statist. Assoc.*, 75, 828–838. [MR0600964](#)
- [29] Serfling, R. J. (1980), ‘*Approximation Theorems of Mathematical Statistics*’, Wiley, New York. [MR0595165](#)
- [30] Sherman, R. P. (1993), “The limiting distribution of the maximum rank correlation estimator”, *Econometrica*, 61(1), 123–137. [MR1201705](#)
- [31] Sherman, R. P. (1994), “Maximal inequalities for degenerate U-processes

- with applications to optimization estimators”, *Ann. Statist.*, 22(1), 439–459. [MR1272092](#)
- [32] Stahel, W. A. (1981), “Robuste Schätzungen: Infinitesimale Optimalität und Schätzungen von Kovarianzmatrizen”, Ph.D. dissertation, ETH, Zurich.
- [33] Stigler, S. M. (1976), “The anonymous Professor Gergonne”, *Hist. Math.*, 3, 71–74. [MR0497722](#)
- [34] Stromberg, A. J., Hawkins, D. M., and Hössjer, O. (2000), “The least trimmed differences regression estimator and alternatives”, *J. Amer. Statist. Assoc.*, 95, 853–864. [MR1804444](#)
- [35] Transtrum, M. K., Machta, B. B., and Sethna, J. P. (2011), “Geometry of nonlinear least squares with applications to sloppy models and optimization”, *Phys. Rev. E*, 83, 036701.
- [36] Tableman, M. (1994), “The influence functions for the least trimmed squares and the least trimmed absolute deviations estimators”, *Statistics & Probability Letters*, 19, 329–337. [MR1278670](#)
- [37] Tukey, J. W. and McLaughlin, D. H. (1963), “Less vulnerable confidence and significance procedures for location based on a single sample: trimming/winsorization 1”, *Sankhyā: The Indian Journal of Statistics, Series A*, 25(3), 331–352. [MR0169354](#)
- [38] Van Der Vaart, A. W. (1998), “*Asymptotic Statistics*”, Cambridge University Press. [MR1652247](#)
- [39] Van Der Vaart, A. W. and Wellner, J. A. (1996), “*Weak Convergence and Empirical Processes with Applications to Statistics*”, Springer, New York. [MR1385671](#)
- [40] Víšek, J. Á. (2006a), “The least trimmed squares. Part I: Consistency”, *Kybernetika*, 42, 1–36. [MR2208518](#)
- [41] Víšek, J. Á. (2006b), “The least trimmed squares. Part II:  $\sqrt{n}$ -consistency”, *Kybernetika*, 42, 181–202. [MR2241784](#)
- [42] Víšek, J. Á. (2006c), “The least trimmed squares. Part III: Asymptotic normality”, *Kybernetika*, 42, 203–224. [MR2241785](#)
- [43] Víšek, J. Á. (2001), “Regression with high breakdown point”, in: *RO-BUST’2000*, pp. 324–356.
- [44] Welsh, A. H. (1987), “The trimmed mean in the linear model”, *Ann. Statist.*, 15(1), 20–36. [MR0885722](#)
- [45] Wu, M., and Zuo, Y. (2009), “Trimmed and Winsorized means based on a scaled deviation”, *J. Statist. Plann. Inference*, 139(2), 350–365. [MR2474011](#)
- [46] Yohai, V. J. (1987), “High breakdown-point and high efficiency estimates for regression”, *Ann. Statist.*, 15, 642–656. [MR0888431](#)
- [47] Yohai, V. J. and Zamar, R.H. (1988), “High breakdown estimates of regression by means of the minimization of an efficient scale”, *J. Amer. Statist. Assoc.*, 83, 406–413. [MR0971366](#)
- [48] Zuo, Y. (2003), “Projection-based depth functions and associated medians”, *Ann. Statist.*, 31, 1460–1490. [MR2012822](#)
- [49] Zuo, Y. (2006), “Multi-dimensional trimming based on projection depth”, *Ann. Statist.*, 34(5), 2211–2251. [MR2291498](#)

- [50] Zuo, Y. (2018), “A new approach for the computation of halfspace depth in high dimensions”, *Communications in Statistics – Simulation and Computation*, 48(3), 900–921. [MR3938062](#)
- [51] Zuo, Y. (2020), “Large sample properties of the regression depth induced median”, *Statistics and Probability Letters*, November 2020, 166, [arXiv:1809.09896](#). [MR4129102](#)
- [52] Zuo, Y. (2021a), “On general notions of depth for regression”, *Statistical Science*, 36(1), 142–157, [arXiv:1805.02046](#). [MR4194208](#)
- [53] Zuo, Y. (2021b), “Robustness of the deepest projection regression depth functional”, *Statistical Papers*, 62(3), 1167–1193. [MR4262190](#)
- [54] Zuo, Y. (2021c), “Computation of projection regression depth and its induced median”, *Computational Statistics and Data Analysis*, 158, 107184. [MR4215773](#)
- [55] Zuo, Y., Serfling, R. (2000), “General notions of statistical depth function”, *Ann. Statist.*, 28, 461–482. [MR1790005](#)