



An algorithm for computing exact least-trimmed squares estimate of simple linear regression with constraints

Lei M. Li*

*Department of Computational Biology and Mathematics, University of Southern California,
CA 90089-1113, USA, 1042 West 36th Place, DRB 289, Los Angeles, CA 900891113, USA*

Received 19 November 2003; received in revised form 2 April 2004; accepted 3 April 2004

Abstract

The least-trimmed squares estimation (LTS) is a robust solution for regression problems. On the one hand, it can achieve any given breakdown value by setting a proper trimming fraction. On the other hand, it has \sqrt{n} -consistency and asymptotic normality under some conditions. In addition, the LTS estimator is regression, scale, and affine equivariant. In practical regression problems, we often need to impose constraints on slopes. In this paper, we describe a stable algorithm to compute the exact LTS solution for simple linear regression with constraints on the slope parameter. Without constraints, the overall complexity of the algorithm is $O(n^2 \log n)$ in time and $O(n^2)$ in storage. According to our numerical tests, constraints can reduce computing load substantially. In order to achieve stability, we design the algorithm in such a way that we can take advantage of well-developed sorting algorithms and softwares. We illustrate the algorithm by some examples.

© 2004 Elsevier B.V. All rights reserved.

Keywords: Least-trimmed squares; Simple regression; Robust; Breakdown value; Constraint

1. Introduction

Statisticians routinely apply regression analysis to fit models to observations. To deal with outliers, we seek for robust and resistant regression procedures. Quite some number of perspectives exist in the literature regarding the definition of robustness.

* Tel.: +1-213-7402407; fax: +1-213-7402437.

E-mail address: lilei@hto.usc.edu (L.M. Li).

For example, Huber (1964) studied robustness from the point of view of minimax variance. Hampel (1971, 1974) proposed the idea of influence function as an asymptotic tool to study robustness. Breakdown point is another important notion in robust analysis. Donoho and Huber (1983) defined a finite-sample version of breakdown point. Consider the classical linear model

$$y = X\beta + \varepsilon, \quad (1)$$

where $y = (y_1, \dots, y_n)'$, $\beta = (\beta_1, \dots, \beta_p)'$, $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)'$, and $X = (x_{ij})_{i=1, \dots, n, j=1, \dots, p}$. For a set of parameter β_0 , we define the residuals by $r(\beta_0) = y - X\beta_0$. The least-squares estimator minimizes the sum of squares $\sum_{i=1}^n r_i^2(\beta)$ over β . The breakdown value of least-squares estimator is $1/n \rightarrow 0$ as $n \rightarrow \infty$. On the other hand, the highest possible breakdown point is 50%. One solution that reaches this bound of breakdown point is the least median of squares (LMS) estimator, cf. Rousseeuw (1984), which minimizes the median of squared residuals: $\min_{\beta} [\text{med}_i r^2(\beta)_i]$. Unfortunately, the asymptotic efficiency of LMS is unsatisfactory because its convergence rate is only of the order $n^{-1/3}$. Another robust solution is the least-trimmed squares (LTS) estimator, which takes as its objective function the sum of smallest squared residuals; see Rousseeuw (1984). We denote the squared residuals in the ascending order by $|r^2(\beta)|_{(1)} \leq |r^2(\beta)|_{(2)} \leq \dots \leq |r^2(\beta)|_{(n)}$. Then the LTS estimate of coverage h is obtained by

$$\min_{\beta} \sum_{i=1}^h |r^2(\beta)|_{(i)}.$$

This definition implies that observations with the largest residuals will not affect the estimate. The LTS estimator is regression, scale, and affine equivariant; see Rousseeuw and Leroy (1987, Lemma 3, Chapter 3). In terms of robustness, we can roughly achieve a breakdown point of ρ by setting $h = [n(1 - \rho)] + 1$. In terms of efficiency, \sqrt{n} -consistency and asymptotic normality similar to M-estimator exist for LTS under some conditions; see Věšek (1996, 2000) for example. Despite its good properties, the computation of LTS remains a problem.

The problem of computing the LTS estimate of β is equivalent to searching for the size- h subset(s) whose least-squares solution achieves the minimum of trimmed squares. The total number of size- h subsets in a sample of size n is $\binom{n}{h}$. A full search through all size- h subsets is impossible unless the sample size is small. Several ideas have been proposed to compute approximate solutions. First, instead of exhaustive search we can randomly sample size- h subsets. Second, Rousseeuw and Van Driessen (1999) proposed a so-called C-step technique (C stands for “concentration”). That is, having selected a size- h subset, we apply the least-squares estimator to them. Next, for the estimated regression coefficients, we evaluate residuals for all observations. Then a new size- h subset with the smallest squared residuals is selected. This step can be iterated starting from any subset. In the case of estimating a location parameter, Rousseeuw and Leroy (1987, pp. 171–172), described a procedure to compute the exact LTS solution. Rousseeuw and Van Driessen (1999) applied this idea to adjust the intercept in the

regression case. The idea of subset search is also relevant for LMS; see Hawkins (1993). Hössjer (1995) described an algorithm for computing the exact LTS estimate of simple linear regression, which requires $O(n^3 \log n)$ computations and $O(n^2)$ storage. A refined algorithm, which requires $O(n^2 \log n)$ computations and $O(n)$ storage, was also sketched. However, Hössjer remarked that the refined algorithm is not stable.

In this paper, we describe an algorithm that computes the exact solution of LTS for simple linear regression with constraints on the slope. The idea is to divide the range of slope into regions in such a way that within each region the order of residuals is unchanged. As a result, the search for LTS within each such region becomes a problem of linear complexity. Hössjer (1995) considered a similar idea from the perspective of dual plots and absolute dual plots. In comparison, we take a direct treatment to deal with constraints by checking Kuhn–Tucker conditions and ties of slopes by defining equivalent classes. Without constraints, the overall complexity of the algorithm is $O(n^2 \log n)$ in time and $O(n^2)$ in storage. In practical regression problems, we often need to impose constraints on slopes. According to our numerical tests, these constraints can reduce computing load substantially. In order to achieve stability, we design the algorithm in such a way that we can take advantage of well-developed sorting algorithms and softwares.

We arrange the materials in this paper as follows. In Section 2, we describe the main results. In Section 3, first we briefly describe the color-correction problem in DNA sequencing that motivates our work. Then we apply our algorithm to a simulated example. In Section 4, we discuss some relevant issues.

2. The main result and algorithm

We consider a simple linear regression problem with a response variable y and an explanatory variable x . Let (x_i, y_i) , $i = 1, \dots, n$, be data.

Proposition 1. *Consider a simple linear regression with restricted slope range, namely, we impose a constraint on the slope: $\beta \in [b_1, b_2]$. Denote the unrestricted least-squares solution by $(\hat{\alpha}, \hat{\beta})$.*

- *If $\hat{\beta} \in [b_1, b_2]$, then the unrestricted solution is also the restricted solution.*
- *If $\hat{\beta} \notin [b_1, b_2]$, then the restricted solution of the slope is either b_1 or b_2 .*

Proof. If $\hat{\beta} \in [b_1, b_2]$, then the unrestricted solution satisfies the requirement. Otherwise, according to the Kuhn–Tucker condition, the solution must be on the boundary; see Luenberger (1984) and Lawson and Hanson (1974). Stark and Parker (1995) provided another algorithm for the problem of bounded-variable least squares. \square

Hereafter, we refer to any subset of the data by the set of their indices. For the sake of completeness and convenience, we write down the formulas of the least-squares

estimate for any size- h subset H :

$$\left\{ \begin{array}{l} \bar{x}_H = \frac{1}{h} \sum_{i \in H} x_i, \\ \bar{y}_H = \frac{1}{h} \sum_{i \in H} y_i, \\ \overline{xx}_H = \frac{1}{h} \sum_{i \in H} x_i^2, \\ \overline{yy}_H = \frac{1}{h} \sum_{i \in H} y_i^2, \\ \overline{xy}_H = \frac{1}{h} \sum_{i \in H} x_i y_i, \\ \overline{cx\bar{y}}_H = \overline{xy}_H - \bar{x}_H \bar{y}_H, \\ \hat{\beta}_H = \overline{cx\bar{y}}_H / [\overline{xx}_H - \bar{x}_H^2], \\ \hat{\alpha}_H = \bar{y}_H - \hat{\beta}_H \bar{x}_H, \\ SS_H = [\overline{yy}_H - \bar{y}_H^2] - \hat{\beta}_H \overline{cx\bar{y}}_H. \end{array} \right. \quad (2)$$

The last quantity is the averaged sum of squares. $\overline{cx\bar{y}}_H$ represents the centerized version of \overline{xy}_H . To proceed, we need some notation and definitions.

Definition 1. For any (α, β) ,

1. their residuals are defined by $r(\alpha, \beta)_i = y_i - (\alpha + \beta x_i)$;
2. $H_{(\alpha, \beta)}$ is a size- h index set that satisfies the following property: $|r(\alpha, \beta)_i| \leq |r(\alpha, \beta)_j|$, for any $i \in H_{(\alpha, \beta)}$ and $j \notin H_{(\alpha, \beta)}$.

On the other hand, given any size- h subset H , we can find the least-squares estimate (α_H, β_H) for the subsample $\{(x_i, y_i), i \in H\}$. Thus, we have established a correspondence between the parameter space and the collection of all size- h subsets. This leads to our following definition of LTS. Please notice that we consider a general case with possible constraints on the slope. For convenience, hereafter in this section we use coverage h instead of a trimming fraction. The breakdown value for coverage h is $(n - h)/n$.

Definition 2. For a given coverage h , the LTS estimate subject to $b_1 \leq \beta \leq b_2$, where $-\infty \leq b_1 < b_2 \leq \infty$, is defined by

1. (α, β) that minimizes $\sum_{i \in H_{(\alpha, \beta)}} r(\alpha, \beta)_i^2$ subject to $b_1 \leq \beta \leq b_2$;

2. or equivalently, the least-squares estimate (α_H, β_H) of a size- h subset H that minimizes SS_H subject to $b_1 \leq \beta_H \leq b_2$; cf. (2).

Namely, we can either search in the parameter space, or search in the space of size- h subsets. This dual form of the LTS solution is the key to the development of the algorithms in this paper.

Proposition 2. For any parameter (α, β) , we sort its residuals in the ascending order, namely, $r(\alpha, \beta)_{\pi(1)} \leq r(\alpha, \beta)_{\pi(2)} \leq \dots \leq r(\alpha, \beta)_{\pi(n)}$, where $\{\pi(1), \pi(2), \dots, \pi(n)\}$ is a permutation of $\{1, 2, \dots, n\}$. Then $H_{(\alpha, \beta)}$ must be from the $n - h + 1$ size- h subsets: $\{\pi(k), \pi(k+1), \dots, \pi(k+h-1)\}$, $k = 1, \dots, n - h + 1$.

Proof. This is true because of the following fact: after taking absolute values, one or more residuals become the smallest. As moving away from the smallest, residuals increase monotonically leftwards or rightwards. \square

The significance of the result is that the complexity for the search of $H_{(\alpha, \beta)}$ reduces to $O(n)$ from $\binom{n}{h}$. We notice that for fixed β , the order of residuals keeps unchanged whatever intercept α is chosen. This leads to the following definition.

Definition 3. For a fixed β , denote $r(\beta)_i = y_i - \beta x_i$. Let $r(\beta)_{(i)}$, $i = 1, \dots, n$, be their ordered residuals. Their partial-ordered averages and partial sums of squares are defined by

$$\begin{cases} \bar{r}(\beta)_{(l:m)} = \frac{1}{m-l+1} \sum_{k=l}^m r(\beta)_{(k)}, \\ SS(\beta)_{(l:m)} = \frac{1}{m-l+1} \sum_{k=l}^m (r(\beta)_{(k)} - \bar{r}(\beta)_{(l:m)})^2, \quad 1 \leq l \leq m \leq n. \end{cases} \quad (3)$$

This structure and Proposition 2 are the basis of the following algorithm for computing the LTS estimate of a location parameter; see pp. 171–172 in [Rousseeuw and Leroy \(1987\)](#).

Definition 4. For $\{y_i, i = 1, \dots, n\}$, we define partial-ordered averages and partial sums of squares, respectively, by

$$\begin{cases} \text{bary}_{(l:l+h-1)} = \frac{\sum_{k=l}^{l+h-1} y_{(k)}}{h}, \\ SS_{(l:l+h-1)} = \frac{\sum_{k=l}^{l+h-1} [y_{(k)} - \bar{y}_{(l:l+h-1)}]^2}{h}. \end{cases}$$

Corollary 1. *The LTS estimate of the location of $\{y_i, i = 1, \dots, n\}$ is given by the partial ordered average(s) that achieves the smallest among $SS_{(l:l+h-1)}$, $l = 1, \dots, n - h + 1$. In case there are ties, the solution is not unique.*

This case corresponds to $\beta = 0$ in the simple linear regression. Proposition 2 implies that we only need to check $n - h + 1$ partial ordered averages for searching the LTS estimate of the location. Partial ordered averages and sums of squares can be calculated in a recursive way; see [Rousseeuw and Leroy \(1987\)](#). Go back to the case of simple linear regression. The next result is the basis of our algorithms.

Proposition 3. *Suppose that for $\beta \in [b_1, b_2]$, the order of $\{r(\beta)_i\}$ keeps unchanged. Namely, there exists a permutation $\{\pi(1), \pi(2), \dots, \pi(n)\}$, such that $r(\beta)_{\pi(1)} \leq r(\beta)_{\pi(2)} \leq \dots \leq r(\beta)_{\pi(n)}$ for any $\beta \in [b_1, b_2]$.*

- Then the size- h subset(s) for the LTS must be from the $n - h + 1$ subsamples: $\{\pi(k), \pi(k + 1), \dots, \pi(k + h - 1)\}$, where $k = 1, \dots, n - h + 1$.
- Moreover, for each subsample $\{\pi(k), \pi(k + 1), \dots, \pi(k + h - 1)\}$, we compute its regular least-squares estimate denoted by $\hat{\beta}_{(k:k+h-1)}$. Then the LTS estimate of β subject to the constraint $\beta \in [b_1, b_2]$ must be from the following: (1) $\{b_1, b_2\}$; (2) or $\{\hat{\beta}_{(k:k+h-1)}, \text{ satisfying } b_1 < \hat{\beta}_{(k:k+h-1)} < b_2, k = 1, \dots, n - h + 1\}$.

Proof. Let $(\hat{\alpha}, \hat{\beta})$ be one solution. According to the assumption, we have $r(\hat{\beta})_{\pi(1)} \leq r(\hat{\beta})_{\pi(2)} \leq \dots \leq r(\hat{\beta})_{\pi(n)}$. Using the same argument leading to Proposition 2, we conclude that the size- h subset $H_{(\hat{\alpha}, \hat{\beta})}$ must be from the $n - h + 1$ size- h subsets: $\{\pi(k), \pi(k + 1), \dots, \pi(k + h - 1)\}$, $k = 1, \dots, n - h + 1$. Then we apply Proposition 1 to the subset $H_{(\hat{\alpha}, \hat{\beta})}$ and the proof is completed. According to this result, we only need to check the sum of squares for the $n - h + 1$ candidates subject to the constraint. \square

Our next move is to divide the range of β into regions in such a way that within each region the order of residuals is unchanged and thus we can apply the above proposition. It is motivated by two technical observations. First, we consider straight lines connecting each pair $\{(x_i, y_i), (x_j, y_j)\}$ satisfying $(x_i, y_i) \neq (x_j, y_j)$. The total number of such pairs is at most $(n(n - 1))/2$. The regions defined by these dividing lines satisfy the above order condition. Second, to compute the least-squares estimates for the $n - h + 1$ size- h subsets within each region, we do not have to repeatedly apply formulas (2). Instead, from one region to the next, we only need to update estimates if a few residuals change their orders. To illustrate the basic idea, we first describe a simplified version of the algorithm for LTS.

Algorithm 1 (LTS algorithm—version 1). *Consider the LTS problem in Definition 2.*

1. For any pair (i, j) such that $x_i \neq x_j$, compute $b^{(i,j)} = (y_j - y_i)/(x_j - x_i)$. Sort $\{b^{(i,j)}\}$ in the range $[b_1, b_2]$ and denote them by $b_1 < b^{[1]} \leq b^{[2]} \leq \dots \leq b^{[L]} < b_2$,

where $L \leq (n(n-1))/2$. Save the pairs of indices corresponding to these slopes.

Remark. The storing of these indices is one subtle yet key difference between the algorithms in this paper and those in Hössjer (1995). Also we exclude all the slopes outside (b_1, b_2) when sorting slopes. We first restrict the scope of this algorithm to the case without slope ties, namely,

Assumption. $b_1 < b^{[1]} < b^{[2]} < \dots < b^{[L]} < b_2$.

2. If $b_1 = -\infty$, go to Step 3. If $b_1 > -\infty$, then we consider the residuals along the lower bound.

- (a) Compute $r(b_1)_i = y_i - b_1 x_i$, $i = 1, \dots, n$. Sort them in the ascending order and denote the ordered residuals by $\{r(b_1)_{(i)}, i = 1, \dots, n\}$.
- (b) Compute the following quantities, for $l = 1, \dots, n - h + 1$,

$$\begin{cases} \bar{r}(b_1)_{(l:l+h-1)} = \frac{1}{h} \sum_{k=l}^{l+h-1} r(b_1)_{(k)}, \\ \overline{rr}(b_1)_{(l:l+h-1)} = \frac{1}{h} \sum_{k=l}^{l+h-1} r(b_1)_{(k)}^2, \\ SS(b_1)_{(l:l+h-1)} = \overline{rr}(b_1)_{(l:l+h-1)} - \bar{r}(b_1)_{(l:l+h-1)}^2, \end{cases} \quad (4)$$

by the recursion

$$\begin{cases} \bar{r}(b_1)_{(l+1:l+h)} = \bar{r}(b_1)_{(l:l+h-1)} + \frac{1}{h} [r(b_1)_{(l+h)} - r(b_1)_{(l)}], \\ \overline{rr}(b_1)_{(l+1:l+h)} = \overline{rr}(b_1)_{(l:l+h-1)} + \frac{1}{h} [r(b_1)_{(l+h)}^2 - r(b_1)_{(l)}^2]. \end{cases} \quad (5)$$

- (c) Save the size- h subset(s) that achieves the minimum of sum of squares.

3. Take a value β such that $b_1 < \beta < b^{[1]}$.

- (a) Compute $r(\beta)_i = y_i - \beta x_i$, $i = 1, \dots, n$. Sort them in the ascending order. For ties, we arrange them in the ascending order of x -values. We denote the ordered residuals by $r(\beta)_{\pi(i)}$, $i = 1, \dots, n$, where $\{\pi(1), \pi(2), \dots, \pi(n)\}$ is a permutation of $\{1, 2, \dots, n\}$. We also denote the inverse of $\{\pi(1), \pi(2), \dots, \pi(n)\}$ by $\{\lambda(1), \lambda(2), \dots, \lambda(n)\}$.

(b) Compute the following quantities, for $l = 1, \dots, n - h + 1$:

$$\left\{ \begin{array}{l} \bar{x}_{(l:l+h-1)} = \frac{1}{h} \sum_{i=l}^{l+h-1} x_{\pi(i)}, \\ \bar{y}_{(l:l+h-1)} = \frac{1}{h} \sum_{i=l}^{l+h-1} y_{\pi(i)}, \\ \overline{xx}_{(l:l+h-1)} = \frac{1}{h} \sum_{i=l}^{l+h-1} x_{\pi(i)}^2, \\ \overline{yy}_{(l:l+h-1)} = \frac{1}{h} \sum_{i=l}^{l+h-1} y_{\pi(i)}^2, \\ \overline{xy}_{(l:l+h-1)} = \frac{1}{h} \sum_{i=l}^{l+h-1} x_{\pi(i)} y_{\pi(i)}, \\ \overline{cx\bar{y}}_{(l:l+h-1)} = \overline{xy}_{(l:l+h-1)} - \bar{x}_{(l:l+h-1)} \bar{y}_{(l:l+h-1)}, \\ \hat{\beta}_{(l:l+h-1)} = \overline{cx\bar{y}}_{(l:l+h-1)} / [\overline{xx}_{(l:l+h-1)} - \bar{x}_{(l:l+h-1)}^2], \\ \hat{\alpha}_{(l:l+h-1)} = \bar{y}_{(l:l+h-1)} - \hat{\beta}_{(l:l+h-1)} \bar{x}_{(l:l+h-1)}, \\ SS_{(l:l+h-1)} = [\overline{yy}_{(l:l+h-1)} - \bar{y}_{(l:l+h-1)}^2] - \hat{\beta}_{(l:l+h-1)} \overline{cx\bar{y}}_{(l:l+h-1)}, \end{array} \right. \quad (6)$$

by the recursion

$$\left\{ \begin{array}{l} \bar{x}_{(l+1:l+h)} = \bar{x}_{(l:l+h-1)} + \frac{1}{h} [x_{\pi(l+h)} - x_{\pi(l)}], \\ \bar{y}_{(l+1:l+h)} = \bar{y}_{(l:l+h-1)} + \frac{1}{h} [y_{\pi(l+h)} - y_{\pi(l)}], \\ \overline{xx}_{(l+1:l+h)} = \overline{xx}_{(l:l+h-1)} + \frac{1}{h} [x_{\pi(l+h)}^2 - x_{\pi(l)}^2], \\ \overline{yy}_{(l+1:l+h)} = \overline{yy}_{(l:l+h-1)} + \frac{1}{h} [y_{\pi(l+h)}^2 - y_{\pi(l)}^2], \\ \overline{xy}_{(l+1:l+h)} = \overline{xy}_{(l:l+h-1)} + \frac{1}{h} [x_{\pi(l+h)} y_{\pi(l+h)} - x_{\pi(l)} y_{\pi(l)}]. \end{array} \right. \quad (7)$$

(c) Update the LTS solution.

4. For $k = 1, 2, \dots, L$, do the following.

(a) Consider $b^{[k]}$ and the corresponding pair of indices (i, j) . Update the permutation π by letting $\pi(\lambda(i)) = j$ and $\pi(\lambda(j)) = i$, and update λ by swapping $\lambda(i)$ and $\lambda(j)$.

(b) Update the quantities in (6).

(c) Update the LTS solution.

5. If $b_2 < \infty$, go through Step 2, replace b_1 by b_2 , and update the LTS solution.

Proposition 4. *The output of Algorithm 1 is the LTS solution.*

Proof. First, subject to the constraint $b_1 < \beta < b^{[1]}$, $r(\beta)_{\pi(1)} \leq r(\beta)_{\pi(2)} \leq \dots \leq r(\beta)_{\pi(n)}$ is always true for the permutation $\{\pi(1), \pi(2), \dots, \pi(n)\}$. Otherwise, we can find two indices (i, j) such that $y_i - \beta x_i < y_j - \beta x_j$ for some β values and $y_i - \beta x_i > y_j - \beta x_j$ for some other β values. Then there exists at least one value $b \in (b_1, b^{[1]})$ such that $y_i - bx_i = y_j - bx_j$ and $x_i \neq x_j$. This leads to $b = (y_j - y_i)/(x_j - x_i)$, a contradiction to the assumption that no other $b^{(i,j)}$ exists between b_1 and $b^{[1]}$. Hence, we can apply Proposition 3 to the interval $[b_1, b^{[1]}]$. As a consequence, we only need to check the $n-h+1$ subsample $\{\pi(l), \pi(l+1), \dots, \pi(l+h-1)\}$, where $1 \leq l \leq n-h+1$. This is exactly Step 2. Next we consider the order structure as β reaches $b^{[1]}$. Remember that we have assumed $b^{[1]} < b^{[2]}$. As β passes $b^{[1]}$ into the interval $(b^{[1]}, b^{[2]})$, the order structure is preserved except for the pair of indices (i, j) such that $b^{(i,j)} = (y_j - y_i)/(x_j - x_i) = b^{[1]}$. Otherwise, a self-contradiction would occur by the same argument as in the interval $(b_1, b^{[1]})$. If we do the swap as in Step 4(a): $\pi(\lambda(i)) = j$ and $\pi(\lambda(j)) = i$, and exchange $\lambda(i)$ and $\lambda(j)$, then the residuals under the new permutation π are still in the ascending order subject to $b^{[1]} \leq \beta < b^{[2]}$. Next, we update the quantities in (6) and apply Proposition 3 to this region. Recursively, we rotate the slope around the origin and repeat this procedure for each region $b^{[k]} \leq \beta < b^{[k+1]}$. \square

Proposition 3 guarantees that we can find the size- h subset(s) for the optimal solution at the end of search. Next, we evaluate the complexity of the algorithm. This requires a more detailed picture of the order structure.

Proposition 5. *If we assume that*

1. *Data contain no identical observations.*
2. *No tie exists in the slopes $\{b^{[l]}\}$, namely, $b_1 < b^{[1]} < b^{[2]} < \dots < b^{[L]} < b_2$.*

Then

- *as β goes from an interval $[b^{[k-1]}, b^{[k]})$ to the next $[b^{[k]}, b^{[k+1]})$, the pair of indices (i_k, j_k) involved in the swap, cf. Step 4(a) of Algorithm 1, must be adjacent to one another in the permutations π in the two intervals $[b^{[k-1]}, b^{[k]})$ and $[b^{[k]}, b^{[k+1]})$;*
- *consequently, we only need to adjust two partial regressions given in Eq. (6) and thus the complexity of Algorithm 1 is quadratic in time except for the part of sorting $\{b^{(i,j)}\}$.*

Proof. When $\beta \in [b^{[k-1]}, b^{[k]})$, we have either $y_{i_k} - \beta x_{i_k} < y_{j_k} - \beta x_{j_k}$ or $y_{i_k} - \beta x_{i_k} > y_{j_k} - \beta x_{j_k}$. Without loss of generality, we assume the former is true, i.e.,

$$\begin{cases} y_{i_k} - \beta x_{i_k} < y_{j_k} - \beta x_{j_k}, & b^{[k-1]} < \beta < b^{[k]}, \\ y_{i_k} - \beta x_{i_k} = y_{j_k} - \beta x_{j_k}, & \beta = b^{[k]}, \\ y_{i_k} - \beta x_{i_k} > y_{j_k} - \beta x_{j_k}, & b^{[k]} < \beta < b^{[k+1]}. \end{cases}$$

Suppose, we have another term (x_u, y_u) whose residual is between these two terms for $\beta \in [b^{[k-1]}, b^{[k]}]$. That is, $y_{i_k} - \beta x_{i_k} < y_u - \beta x_u < y_{j_k} - \beta x_{j_k}$. This implies

$$0 < (y_u - y_{i_k}) - \beta(x_u - x_{i_k}) < (y_{j_k} - y_{i_k}) - \beta(x_{j_k} - x_{i_k}).$$

When $\beta \rightarrow b^{[k]}$, we have

$$0 \leq (y_u - y_{i_k}) - \beta(x_u - x_{i_k}) \leq (y_{j_k} - y_{i_k}) - \beta(x_{j_k} - x_{i_k}) = 0.$$

Hence, we have $y_u - y_{i_k} = \beta(x_u - x_{i_k})$ and $b^{[k]} = (y_u - y_{i_k})/(x_u - x_{i_k})$. This conflicts the assumption. Similarly, the two terms indexed by i_k, j_k are still next to each other in the interval $[b^{[k]}, b^{[k+1]}]$. The only change in π is the swap of their positions. This exchange involves only two partial regressions given in (6). \square

Now, we consider the general case in which ties may occur in $b^{(i,j)}$'s. First, we define a relation between data pairs.

Definition 5. For any pair of data (x_{i_1}, y_{i_1}) , (x_{i_2}, y_{i_2}) , we denote their relation by: $i_1 \stackrel{b}{\sim} i_2$ if

- either $x_{i_1} \neq x_{i_2}$, $b = (y_{i_2} - y_{i_1})/(x_{i_2} - x_{i_1})$,
- or $(x_{i_1}, y_{i_1}) = (x_{i_2}, y_{i_2})$.

In other words, $i_1 \stackrel{b}{\sim} i_2$ if and only if $r(b)_{i_1} = r(b)_{i_2}$.

Proposition 6. The relation $\stackrel{b}{\sim}$ is an equivalent relation because it is

- reflexive, namely, $i \stackrel{b}{\sim} i$;
- symmetric, namely, if $i_1 \stackrel{b}{\sim} i_2$, then $i_2 \stackrel{b}{\sim} i_1$;
- transitive, namely if $i_1 \stackrel{b}{\sim} i_2, i_2 \stackrel{b}{\sim} i_3$, then $i_1 \stackrel{b}{\sim} i_3$.

For each value of b , we collect all the indices $\{i_1^b, i_2^b, \dots, i_{b_m}^b\}$ that are involved in the relation $\stackrel{b}{\sim}$. Due to the equivalence, we can partition $\{i_1^b, i_2^b, \dots, i_{b_m}^b\}$ into non-overlapping transitive sets in such a way that any pair $\{j_1, j_2\}$ within one set satisfy $j_1 \stackrel{b}{\sim} j_2$.

Proof. The reflexive and symmetric parts are obvious. From $b = (y_{i_2} - y_{i_1})/(x_{i_2} - x_{i_1}) = (y_{i_3} - y_{i_1})/(x_{i_3} - x_{i_1})$, we have $b = (y_{i_3} - y_{i_2})/(x_{i_3} - x_{i_2})$ by straightforward calculation. This proves the transitive part. The partition follows the general practice of defining equivalent classes based on an equivalent relation. We note that it is possible to have more than one non-overlapping transitive sets for one relation $\stackrel{b}{\sim}$. \square

Consider all the index pairs $\{(i_b, j_b)\}$ associated with a slope value b . We can represent them by an undirected graph $G = (V, E)$, where V is the set of all indices involved and E includes all the edges $\{(i_b, j_b)\}$. Then each transitive set is equivalent to a connected component of G . The depth-first search algorithm can find the connected components of a graph in $O(|V| + |E|)$ running time. We can also use data structure for disjoint sets to solve the problem, see Cormen et al. (1990).

Proposition 7. Suppose $\{j_1^b, j_2^b, \dots, j_t^b\}$ is one transitive set associated with $\overset{b}{\sim}$. Then

1. In the two regions right before and right after b , the indices $\{j_1^b, j_2^b, \dots, j_t^b\}$ form a consecutive block in the permutation $(\pi(1), \pi(2), \dots, \pi(n))$. Remember that π is defined in Algorithm 1 so that $r(\beta)_{\pi(1)} \leq r(\beta)_{\pi(2)} \leq \dots \leq r(\beta)_{\pi(n)}$.
2. If $r(\beta)_{j_1^b} \leq r(\beta)_{j_2^b} \leq \dots \leq r(\beta)_{j_t^b}$ in the region right before b ($\beta < b$), then the order structure is valid for the reversed indices: $r(\beta)_{j_1^b} \geq r(\beta)_{j_2^b} \geq \dots \geq r(\beta)_{j_t^b}$ in the region right after b ($\beta > b$). In fact, $x_{j_1^b} \leq x_{j_2^b} \leq \dots \leq x_{j_t^b}$ must be true.

Proof. The argument to prove the first part is similar to that of Proposition 5. That is, if another term were between a pair of $\{j_1^b, j_2^b, \dots, j_t^b\}$, then a contradiction to the definition of transitive sets would occur. To prove the second part, we consider $r(\beta) - r(b)$. Due to the definition of $\overset{b}{\sim}$ and transitive sets, we have

$$r(b)_{j_1^b} = r(b)_{j_2^b} = \dots = r(b)_{j_t^b},$$

then in the region right before b , $\beta < b$, and

$$r(\beta)_{j_1^b} - r(b)_{j_1^b} \leq r(\beta)_{j_2^b} - r(b)_{j_2^b} \leq \dots \leq r(\beta)_{j_t^b} - r(b)_{j_t^b},$$

namely

$$-(\beta - b)x_{j_1^b} \leq -(\beta - b)x_{j_2^b} \leq \dots \leq -(\beta - b)x_{j_t^b}.$$

Hence

$$x_{j_1^b} \leq x_{j_2^b} \leq \dots \leq x_{j_t^b}.$$

As β passes b , $\beta - b > 0$, the above inequalities reverse, namely

$$r(\beta)_{j_1^b} - r(b)_{j_1^b} \geq r(\beta)_{j_2^b} - r(b)_{j_2^b} \geq \dots \geq r(\beta)_{j_t^b} - r(b)_{j_t^b}.$$

This completes the proof. \square

Next, we describe the complete version of the LTS algorithm. It deals with the complication of slope ties.

Algorithm 2 (LTS algorithm—version 2).

- In Step 1 of Algorithm 1, we compute $b^{(i,j)} = (y_j - y_i)/(x_j - x_i)$ for each pair satisfying $x_i \neq x_j$. Then for each value $b \in \{b^{(i,j)}\}$, we find all the transitive sets $\{j_1^b, j_2^b, \dots, j_t^b\}$. We keep only one copy of b for each transitive subset in the sorted slopes, $-\infty < b^{[1]} \leq b^{[2]} \leq \dots \leq b^{[L]} < \infty$. Please note that it is still possible to have ties in $\{b^{[k]}\}$ because two different transitive subsets may associate with the same b .
- In Step 1 of Algorithm 1, we also sort $\{x_i\}$.
- In Step 4 of Algorithm 1, we rotate the slope around the origin through intervals divided by $\{b^{[k]}\}$. For each specific value $b \in \{b^{[k]}\}$, it may associate with more than one transitive sets and we deal with them one at a time. Suppose $\{j_1^b, j_2^b, \dots, j_t^b\}$ is a transitive set such that $x_{j_1^b} \leq x_{j_2^b} \leq \dots \leq x_{j_t^b}$. Then as β moves from the interval to the next divided by b , we reverse the positions of these indices in the permutation π .

This results in $\{j_t^b, j_{t-1}^b, \dots, j_1^b\}$. Next, we update the quantities in (6). We iterate this operation for each of the transitive set under $\stackrel{b}{\sim}$.

According to Proposition 7, the residuals indexed by $\{j_1^b, j_2^b, \dots, j_t^b\}$ is indeed in the ascending order in the region just before b . This justifies Step 4. In addition, we do not have to do sorting again in Step 4 because $\{x_i\}$ has been sorted in Step 1.

Proposition 8. *The complexity of Algorithm 2 is $\max\{O(L^2 \log L), O(n^2)\}$ in time and $O(L^2)$ in memory, where L is the number of $\{b^{(i,j)}\}$ in the range $[b_1, b_2]$.*

Proof. In Step 1, we need to compute slopes $\{b^{(i,j)}\}$, the number of which is at most $(n(n-1))/2$. However, we only need to save and sort the slope values in the range $[b_1, b_2]$. The heap-sort algorithm allows us to sort these slopes with $O(L^2 \log L)$ operations in place; see Cormen et al. (1990). Similarly, we can sort $\{x_i\}$ in place in $O(n \log n)$ time. In Step 3, it takes $9(n-h+1)$ memory units to save those quantities in (6). In Step 4, if one transitive set $\{j_1^b, j_2^b, \dots, j_t^b\}$ has t terms, then we need to update $2t$ quantities in (6). By doing so we have dealt with $(t(t-1))/2$ terms in the set $\{b^{(i,j)}\}$. Thus throughout the iteration, the total number of terms being updated is $O(L^2)$. \square

Stability: The algorithm described in Hössjer (1995) requires only $O(n)$ memory. But it is not stable as the author noted. Instead of storing and sorting the values of $\{b^{(i,j)}\}$, the refined algorithm computes $b^{[k]}$ at each dividing line recursively. This hinges on the correct ranking of current residuals. Once a mistake is caused by roundoff errors, it propagates in the recursion and the whole algorithm sidetracks. In comparison, we compute $\{b^{(i,j)}\}$ from raw data and sort them systematically. In the meantime, we record the corresponding pairs of indices. It is the sorted slopes and their associated indices that totally determine which indices should be updated as the searching process moves from one region to the next. In fact, the ranks of residuals in each region have already been determined even before the search starts. In other words, the ranks of residuals in each region are independent of the quantities in (6) and thus are free from their roundoff errors. Especially, in cases of ties, which correspond to intersection points in dual plots, transitive sets are constructed solely from $\{b^{(i,j)}\}$ and associated indices and are not subject to roundoff errors of other intermediate quantities. Thus, the stability of Algorithms 1 and 2 depends primarily on the correct sorting of $\{b^{(i,j)}\}$. Consequently, the well-developed sorting algorithms and softwares provide a solid support to Algorithms 1 and 2.

3. Examples

3.1. DNA sequencing and color-correction

In DNA sequencing, we want to measure four kinds of dye concentrations; see Adams et al. (1994). However, the spectra of the four dyes used in fluorescence-based

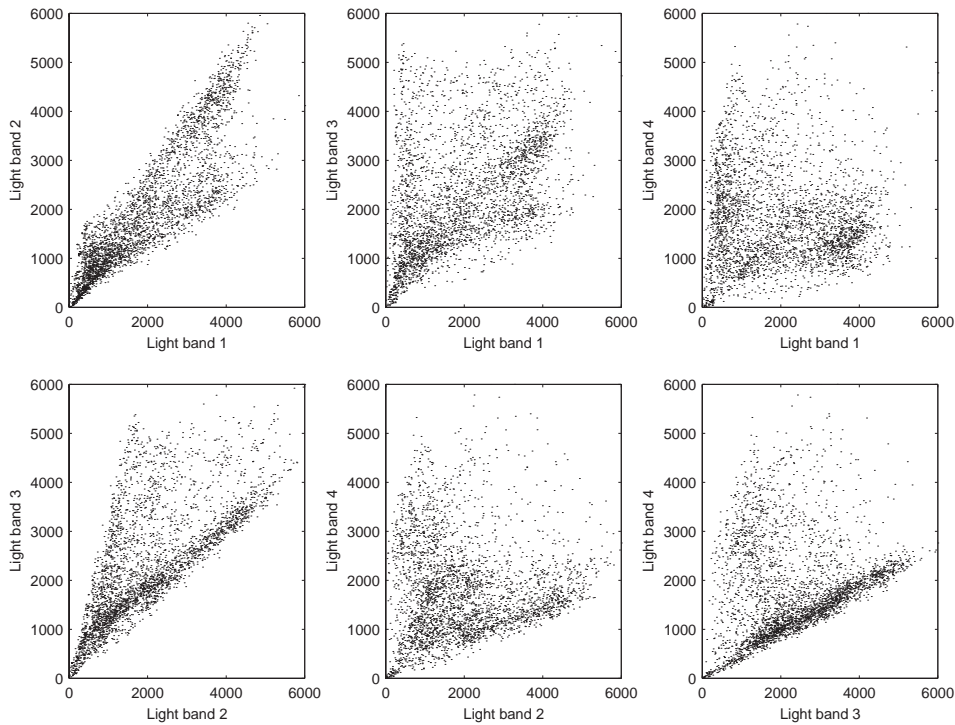


Fig. 1. Pairwise scatter plots of the four light intensities of a slab gel data set.

DNA sequencing overlap, and thus the cross-talk phenomenon arises. Approximately, we describe the relationship between the unknowns—four dye concentrations $C(t)$, $G(t)$, $A(t)$ and $T(t)$ —and the observations—four fluorescence intensities— $I_1(t)$, $I_2(t)$, $I_3(t)$, $I_4(t)$ —at an electrophoretic time t by the following linear system

$$\begin{bmatrix} I_1(t) \\ I_2(t) \\ I_3(t) \\ I_4(t) \end{bmatrix} = \begin{bmatrix} 1 & w_{12} & w_{13} & w_{14} \\ w_{21} & 1 & w_{23} & w_{24} \\ w_{31} & w_{32} & 1 & w_{34} \\ w_{41} & w_{42} & w_{43} & 1 \end{bmatrix} \begin{bmatrix} C(t) \\ G(t) \\ A(t) \\ T(t) \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}, \quad (8)$$

where (b_1, b_2, b_3, b_4) is the baseline due to background light.

The goal of color-correction is to reconstruct the dye concentrations using data of fluorescence intensities. The effective cross-talk matrix $[w_{ij}]$, according to our data analysis, varies under different DNA sequencing conditions and needs to be estimated adaptively using data to be color-corrected. Because both the input and the system are unknown, this is a “blind inversion problem”. In Li (2003), we propose a general framework, referred to as “blind inversion needs distribution (BIND)”, to solve this kind of problems. The basis idea is to use the distributional information of the light

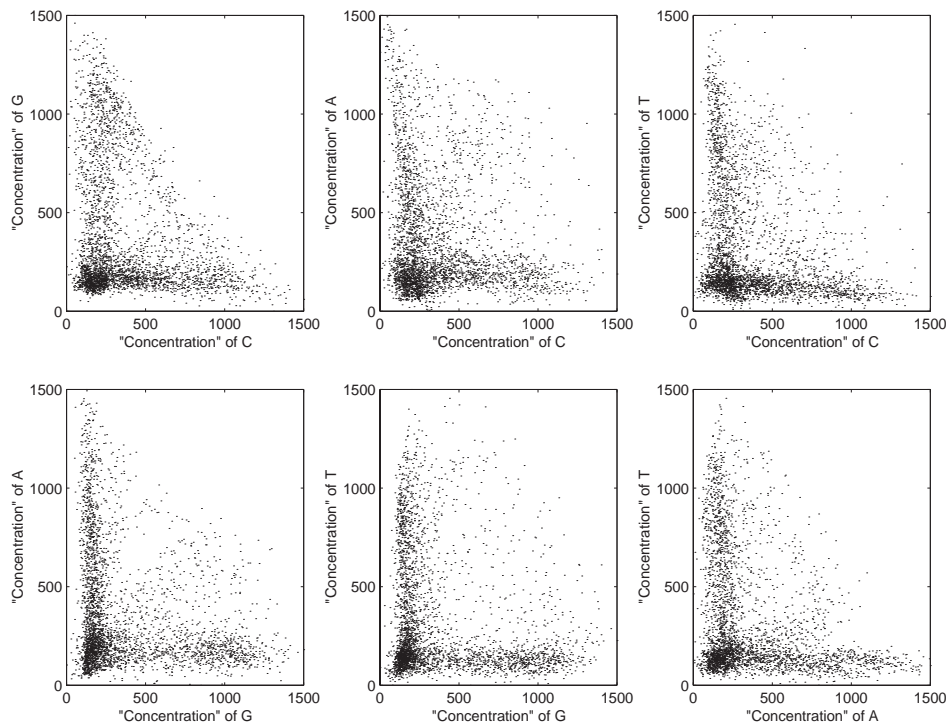


Fig. 2. Pairwise scatter plots of the reconstructed dye concentrations.

intensities, which is known from the sequencing design. We display the distribution of light-intensities and dye concentrations by pairwise scatter plots; see Figs. 1 and 2. Roughly speaking, the data on the two boundary arms of the six pairwise scatter plots in Fig. 1 are the “sufficient” statistics for the 12 parameters in Eq. (8).

One step of our color-correction procedure relies on the estimation of slopes of the boundary arms. We use a binning strategy to sample data from these boundaries. But, occasionally data from other regions are included too. Thus we need a robust regression algorithm. We note that the four dyes used in sequencing are so designed that the constraints $\{0 \leq w_{ij} \leq 1, i \neq j\}$ are valid. Details on color-correction can be found in Li and Speed (1999), and Li (2003). This is one example that motivates the work in this paper.

3.2. A simulated example

We simulated a data set of size 3000 from a mixture model with three components. The description of the model is as follows:

1. 2000 independent samples from: $Y = X + \varepsilon$, where $X \sim \text{Uniform}(0, 100)$, $\varepsilon \sim N(0, 30^2)$;

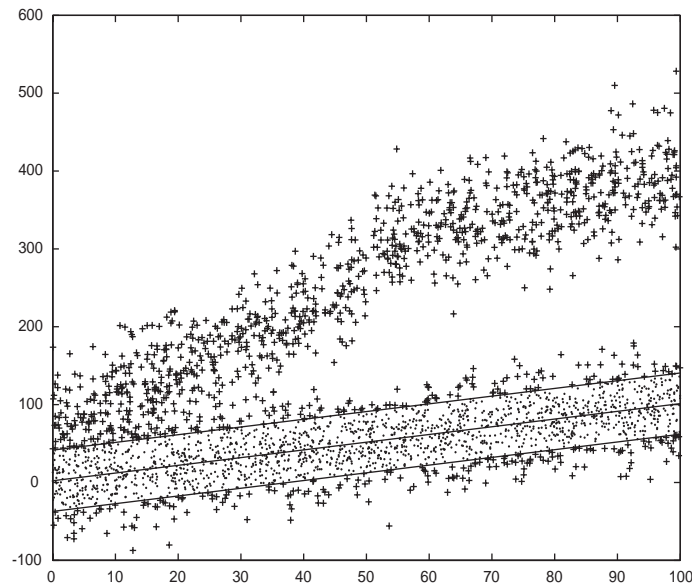


Fig. 3. The scatter plot of a simulated data set of size 3000 from a mixture model with three components. The straight line in the middle corresponds to the LTS solution with trimming fraction 0.45. We include an upper and lower band to show the status of each observation. Those observations falling inside the bands, marked by dots, are in the size- h subset for the LTS. Those observations falling outside the bands, marked by +’s, are trimmed out, or are not in the final size- h subset.

2. 500 independent samples from: $Y = 60 + 4.0 * X + \varepsilon$, where $X \sim \text{Uniform}(0, 50)$, $\varepsilon \sim N(0, 40^2)$;
3. 500 independent samples from: $Y = 200 + 2.0 * X + \varepsilon$, where $X \sim \text{Uniform}(51, 100)$, $\varepsilon \sim N(0, 40^2)$.

The first component contributes two-thirds of the observations and is the main component. The other two contribute one-sixth each. The scatter plot is shown in Fig. 3. We want to estimate the slope and intercept of the first component and minimize the influence from the other two components. We applied our LTS algorithm to the data set with various trimming fractions. Table 1 shows the numerical results. The intercept and slope of the main component are 0.0 and 1.0, respectively. The most accurate estimate of the slope is obtained when the trimming fraction is 0.45. In Table 1, we also show the distribution of the three components in the optimal size- h subsets. As the trimming fraction goes under one-third, the influence from the other two components becomes serious. In the extreme case, we keep all the residuals and this results in the normal least-squares estimate. In all cases except for the normal least squares, the solutions given by “ltsreg()” function in Splus did not match the exact solutions. If we restrict the slope to the range $[0, 2]$, then the total amount of calculation reduces to 25% of that without constraints. The computation takes 2 seconds on a desktop PC equipped with a Pentium 4 processor (2.26G Hz).

Table 1

The LTS estimates with different trimming fractions. The distribution of the three components in the optimal size- h subsets are also shown.

Trimming fraction	0.5	0.45	0.35	0.25	0.15	0.05	0
Intercept	−1.1989	2.1508	6.3484	28.1769	50.6181	35.2475	26.7885
Slope	1.0318	0.9912	0.9190	0.6311	0.5714	1.4285	1.8546
Data distribution							
Component 1	1483	1623	1897	1997	2000	2000	2000
Component 2	17	27	53	253	487	499	500
Component 3	0	0	0	0	63	351	500

In Fig. 3, the straight line in the middle corresponds to the LTS solution with trimming fraction 0.45. We include an upper and lower band to show the status of each observation. Those observations falling inside the bands, marked by dots, are in the size- h subset for the LTS. Those observations falling outside the bands, marked by +’s, are trimmed out, or are not in the final size- h subset.

4. Discussion

4.1. Single-case and high-breakdown diagnostics

Like the concept of robustness, outliers have also been considered from different views. One treatment is to consider the influence of one observation at a time. Cook’s distance (Cook, 1977) and jackknife residuals, see Velleman and Welsch (1981), are two standard techniques along this line. As Rousseeuw and Leroy (1987, Chapter 6) pointed out, single-case diagnostics has its limitations such as masking effect in the presence of multiple outliers. The telephone data in Rousseeuw and Leroy (1987, pp. 222 and 229) is one case that single-case diagnostics fails. On the other hand, high-breakdown diagnostics is able to cope with simultaneous influence of multiple outliers. LMS has been used to serve such a purpose in the literature. Likewise, LTS has the same desirable high-breakdown property. The algorithm offered in the article can help researchers to conduct diagnostics using LTS.

4.2. Multiple solutions

In general, the LTS solution is not unique. Consider a sample of $\{1, 2, 3, 4, 5, 6, 700\}$ and we want to estimate its location parameter. This is a special case of simple linear regression by letting $x_i = 0$. Now the sample size is seven and we take h to be five. One solution is three obtained by discarding 6 and 700. The other solution is 4 obtained by discarding 1 and 700. However, if data are real numbers with random perturbations, the chance of multiple solutions is small.

4.3. Complexity bound

It is useful to evaluate the lower bound on the complexity of computing LTS for simple linear regression. The overall complexity of our algorithm is $O(n^2 \log n)$ in time.

The core iteration of Algorithm 2 takes $O(n^2)$ operations. The extra term $\log n$ comes from sorting $\{b^{(i,j)} = (y_j - y_i)/(x_j - x_i)\}$. We are curious whether the logarithm factor can be removed. The sorting of slopes belongs to the general class of problems: sorting with partial information. For a simpler problem of sorting $\{x_i + y_j, i, j = 1, \dots, n\}$, to date, no $O(n^2)$ algorithm is known; see Fredman (1976) and Steiger and Streinu (1995). This suggests that we do not have much room for improvement along the line of this article. In comparison, the existing algorithms of computing LMS for simple linear regression achieve complexity of $O(n^2)$ in time and $O(n^2)$ in memory or $O(n^2 \log n)$ in time and $O(n)$ in memory, cf. Souvaine and Steele (1987). Thus our algorithm for LTS is competitive in consideration of its advantage of efficiency over LMS.

4.4. Final remarks

In this article, we offer a fast and stable algorithm to compute LTS estimate for simple linear regression. It has been known for quite some time that LTS enjoys features such as reasonable efficiency and high-breakdown value. Unfortunately, it has not been widely used due to the lack of fast implementation. We hope our work will help improve the situation for the case of simple linear regression. The generalization to the multiple linear regression is an interesting research topic.

Acknowledgements

This research is supported by the NSF Grant DMS-9971698 and a NIH CEGS grant. I would also like to thank Prof. Michael Waterman and Simon Tavaré at University of Southern California for helpful discussions, and Chao Cheng for help in programming.

References

- Adams, M.D., Fields, C., Venter, J.C. (Eds.), 1994. Automated DNA Sequencing and Analysis. Academic Press, London, San Diego.
- Cook, D.R., 1977. Detection of influential observation in linear regression. *Technometrics* 19, 15–18.
- Cormen, T.H., Leiserson, C.E., Rivest, R.L., 1990. Introduction to Algorithms. The MIT Press, Cambridge, MA.
- Donoho, D.L., Huber, P.J., 1983. The notion of the breakdown point. In: Bickel, P., Doksum, K., Hodges Jr., J.L. (Eds.), *A Festschrift for Erich Lehmann*, Wadsworth, Belmont, CA.
- Fredman, M.L., 1976. How good is the information theory bound in sorting? *Theoret. Comput. Sci.* 1, 355–361.
- Hampel, F.R., 1971. A general qualitative definition of robustness. *Ann. Math. Statist.* 42, 1887–1896.
- Hampel, F.R., 1974. The influence curve and its role in robust estimation. *J. Amer. Statist. Assoc.* 69, 383–393.
- Hawkins, D.M., 1993. The feasible set algorithm for least median of squares regression. *Comput. Statist. Data Anal.* 16, 81–101.
- Hössjer, O., 1995. Exact computation of the least trimmed squares estimate in simple linear regression. *Comput. Statist. Data Anal.* 19, 265–282.
- Huber, P.J., 1964. Robust estimation of a location parameter. *Ann. Math. Statist.* 35, 73–101.
- Lawson, C.L., Hanson, R.J., 1974. Solving Least Squares Problems. Prentice-Hall, Englewood Cliffs, NJ.

- Li, L., 2003. Blind inversion needs distribution (BIND): the general notion and case studies. In: Goldstein, D. (Ed.), *Science and Statistics: A Festschrift for Terry Speed*, IMS Lecture Note Series, Vol. 40, pp. 273–293.
- Li, L., Speed, T.P., 1999. An estimate of the color separation matrix in four-dye fluorescence-based DNA sequencing. *Electrophoresis* 20, 1433–1442.
- Luenberger, D.G., 1984. *Linear and Nonlinear Programming*. Addison-Wesley Publishing Company, Reading, MA.
- Rousseeuw, P.J., 1984. Least median of squares regression. *J. Amer. Statist. Assoc.* 79, 871–880.
- Rousseeuw, P.J., Leroy, A.M., 1987. *Robust Regression and Outlier Detection*. Wiley, New York.
- Rousseeuw, P.J., Van Driessen, K., 1999. Computing LTS regression for large data sets. Technical Report, University of Antwerp.
- Souvaine, D.L., Steele, J.M., 1987. Time- and space-efficient algorithms for least median of squares regression. *J. Amer. Statist. Assoc.* 82, 794–801.
- Stark, P.B., Parker, R.L., 1995. Bounded-variable least-squares: an algorithm and applications. *Comput. Statist.* 10, 129–141.
- Steiger, W., Streinu, I., 1995. A pseudo-algorithmic separation of lines from pseudo-lines. *Inform. Process. Lett.* 53, 295–299.
- Velleman, P.F., Welsch, R.E., 1981. Efficient computing of regression diagnostics. *Amer. Statist.* 35, 234–242.
- Víšek, J.Á., 1996. Sensitivity analysis of M-estimates. *Ann. Inst. Statist. Math.* 48, 469–495.
- Víšek, J.Á., 2000. On the diversity of estimates. *Comput. Statist. Data Anal.* 34, 67–89.