

1 **Finite Computation of the  $\ell_1$  Estimator from Huber's  $M$ -Estimator**  
 2 **in Linear Regression**

3 **M. Ç. Pınar, Ankara**

4 Received September 17, 2003; revised September 25, 2003

5 Published online: ■ ■ ■

6 © Springer-Verlag 2003

7 **Abstract**

8 We review and extend previous work on the approximation of the linear  $\ell_1$  estimator by the Huber  
 9  $M$ -estimator based on the algorithms proposed by Clark and Osborne [7], and Madsen and Nielsen  
 10 [12]. Although the Madsen-Nielsen algorithm is a promising one, it is guaranteed to terminate finitely  
 11 under certain assumptions. We describe a variant of the Madsen-Nielsen algorithm to compute the  $\ell_1$   
 12 estimator from the Huber  $M$ -estimator in a finite number of steps without any restrictive steps nor  
 13 assumptions. Summary computational results are given.

14 *Keywords:* Multiple linear regression, the  $\ell_1$  estimator, huber's  $M$ -estimator, finite algorithms.  
 15

16 **1 Introduction**

17 Consider the linear model

$$r = A^T x - b \quad (1)$$

19 where  $x \in \mathbb{R}^n$ ,  $b \in \mathbb{R}^m$  is the vector of dependent observations,  $A \in \mathbb{R}^{n \times m}$  (with  
 20  $m > n$ ) is the matrix of independent observations and  $r \in \mathbb{R}^m$  is the vector of  
 21 residuals. The purpose of this work is to review and extend algorithms for  
 22 computing the linear  $\ell_1$  estimator using Huber's  $M$ -estimator in (1). The linear  $\ell_1$   
 23 estimation problem consists of finding a vector  $x^* \in \mathbb{R}^n$  to the following mini-  
 24 mization problem:

25 [L1]

$$\text{minimize } G(x) \equiv \|r\|_1 \equiv \|A^T x - b\|_1. \quad (2)$$

27 The notation  $\|z\|_1$  is used to denote the  $\ell_1$  norm which is the sum of absolute  
 28 values of the components of  $z$ , i.e.,  $\|z\|_1 = \sum_{i=1}^n |z_i|$ . The linear  $\ell_1$  estimation  
 29 problem is more difficult than the the least squares problem where the  $\ell_2$  norm is  
 30 used. The  $\ell_2$  estimation problem admits a closed form solution whereas the  $\ell_1$   
 31 estimation problem is of a combinatorial nature as it can be recast as a linear  
 32 programming problem.

1 The  $\ell_2$  estimator is known to be the maximum likelihood estimator under the  
 2 assumption that the residuals  $r = A^T x - b$  have independent identical normal  
 3 distributions. However, the  $\ell_2$  estimator is quite sensitive to deviations from this  
 4 assumption, and the presence of a few outliers among data points may have a  
 5 significant effect. The interest for the  $\ell_1$  estimator stems from its robustness in the  
 6 face of outliers as discussed in [9]. Hence, despite the increase in computational  
 7 difficulty compared to the  $\ell_2$  case, the  $\ell_1$  estimator was studied also extensively;  
 8 see e.g. [4, 21] for a review of developments until 1982. There has been renewed  
 9 interest in the  $\ell_1$  estimator as evidenced by the emergence of recent ideas in [8, 19,  
 10 20, 22].

11 An alternative robust estimator which does not involve nonsmooth optimization  
 12 was proposed by Huber [9] as the minimizer  $x^*$  of the function

$$\Phi(x) = \sum_{i=1}^m \phi(r_i, \gamma) \quad (3)$$

14 where

$$\phi(r_i, \gamma) = \begin{cases} \frac{1}{2\gamma} r_i^2 & \text{if } |r_i(x)| \leq \gamma \\ |r_i(x)| - \frac{\gamma}{2} & \text{if } |r_i(x)| \geq \gamma, \end{cases} \quad (4)$$

16 and  $\gamma$  is a positive scalar to be estimated from the data. This estimator was shown  
 17 by Huber to be a maximum likelihood estimator for a perturbed normal distri-  
 18 bution and became known as Huber's  $M$ -estimator. Interestingly, there exist  
 19 intimate relationships between the  $\ell_1$  estimator and the Huber  $M$ -estimator. This  
 20 is to be expected since

$$\lim_{\gamma \rightarrow 0} \phi(r_i, \gamma) = |r_i|. \quad (5)$$

22 This fact has been noticed and studied by essentially three research groups  
 23 resulting in the papers by Clark [6], Clark and Osborne [7], Madsen and Nielsen  
 24 [11, 12], Madsen, Nielsen and Pınar [13, 15] and Li and Swetits [10], with valuable  
 25 insights and algorithms for computing both the  $M$ -estimator and the  $\ell_1$  estimator.  
 26 The first finite algorithm for computing the  $\ell_1$  estimator from the Huber  
 27  $M$ -estimator was proposed by Clark and Osborne. Later, this algorithm was  
 28 extended by Madsen and Nielsen and Madsen, Nielsen and Pınar. The Madsen-  
 29 Nielsen algorithm was reported in [12] to be up competitive with the Barrodale-  
 30 Roberts implementation of the simplex algorithm for the  $\ell_1$  estimation problem  
 31 [3], a significant contribution considering that the Barrodale-Roberts algorithm is  
 32 regarded as one of the most efficient algorithms in this area. This was testimony to  
 33 the promise of the new approach. Later, Madsen, Nielsen and Pınar [15] used this  
 34 algorithm to solve linear programming problems, extending both the theory and  
 35 practice of the new algorithm. All these algorithms have guaranteed finite ter-  
 36 mination under some restrictive assumptions. Later, Li and Swetits proposed a  
 37 recursive variant of the Madsen and Nielsen algorithm and proved its finiteness

1 property under the full rank assumption on  $A$ . Under the light of the above  
 2 discussion the purpose of the present paper is to review the computational ties  
 3 between the Huber  $M$ -estimator and the linear  $\ell_1$  estimator and to give a new  
 4 finite algorithm for computing the  $\ell_1$  estimator from the Huber  $M$ -estimator. The  
 5 new algorithm consists of a simple modification of the Madsen-Nielsen algorithm,  
 6 and terminates finitely without any assumptions. It is inspired from the original  
 7 Clark-Osborne algorithm. Comparative computational results with the modified  
 8 Madsen-Nielsen algorithm show that it is competitive with the most successful  
 9 implementations of the simplex type algorithms.

10 The plan of this paper is as follows. In Section 2, we will give basic properties of  
 11 the  $\ell_1$  estimator, some key results on Huber's  $M$ -estimator and the connections  
 12 between these two, respectively. We will study the algorithmic contributions of  
 13 Clark and Osborne, and Madsen and Nielsen in Section 3 and 4, respectively. We  
 14 give further results on the connection between the two estimators in Section 5. We  
 15 propose an extension of the Madsen-Nielsen algorithm in Section 6 and prove its  
 16 finite termination property. We also summarize computational experience with  
 17 the modified algorithm in Section 6.

18

## 2 Properties

19 In this section we review some relevant properties of the  $\ell_1$  estimator, the Huber  
 20  $M$ -estimator and their connections, in this order, respectively.

21

### 2.1 The $\ell_1$ Estimator

22 The  $\ell_1$  estimator is characterized by the following necessary and sufficient con-  
 23 dition for optimality [21]:  $x$  is an  $\ell_1$  estimator iff there exists  $\lambda_j \in [-1, 1]$  such that

$$\sum_{j \in \mathcal{A}_0(x)} \lambda_j a_j + \sum_{j \notin \mathcal{A}_0(x)} a_j s_j = 0 \quad (6)$$

25 where  $\mathcal{A}_0(x) = \{j : r_j(x) = 0\}$  and  $s_j$  is defined for all  $j$  as

$$s_j(x) = \begin{cases} -1 & \text{if } r_j(x) < 0 \\ 0 & \text{if } |r_j(x)| = 0 \\ 1 & \text{if } r_j(x) > 0. \end{cases} \quad (7)$$

27 An interesting duality result links the  $\ell_1$  estimator with linear programming. It can  
 28 be shown using Lagrange duality that the dual problem to [L1] is given as

29 [D1]

$$\begin{aligned} \max \quad & -b^T y \\ \text{s.t.} \quad & Ay = 0 \\ & -\mathbf{1} \leq y \leq \mathbf{1} \end{aligned}$$

1 where  $\mathbf{-1}$  and  $\mathbf{1}$  denote vectors with components  $-1$  and  $1$ , respectively. Fur-  
 2 thermore,  $x$  solves [L1] and  $y$  solves the dual problem if and only if  $y$  satisfies  
 3  $Ay = 0$  and the following conditions hold

$$r_j(x) < 0 \implies y_j = -1, \quad (8)$$

$$r_j(x) > 0 \implies y_j = 1, \quad (9)$$

5 and

$$-1 < y_j < 1 \implies r_j(x) = 0. \quad (10)$$

7 It is easy to notice that these conditions are fully equivalent to the optimality  
 8 condition (6).

## 9 *2.2 The Huber M-Estimator*

10 Define for a given threshold  $\gamma > 0$  the sign vector

$$s^\gamma(x) = [s_1^\gamma(x), \dots, s_m^\gamma(x)] \quad (11)$$

12 with

$$s_i^\gamma(x) = \begin{cases} -1 & \text{if } r_i(x) \leq -\gamma \\ 0 & \text{if } |r_i(x)| < \gamma \\ 1 & \text{if } r_i(x) \geq \gamma. \end{cases} \quad (12)$$

14 If  $s = s^\gamma(x)$  then we also denote  $W_s$  the  $m \times m$  diagonal matrix whose  $i$ th diagonal  
 15 entry is given by  $1 - s_i^2$ . Alternatively, we will also use  $W(x)$  to denote the diag-  
 16 onalmatrix associated with  $s^\gamma(x)$  directly. Now, the Huber  $M$ -estimation problem  
 17 can be recast as the following minimization problem:

18 [SL1]

$$\text{minimize } G_\gamma(x) \equiv \frac{1}{2\gamma} r^T W_s r + s^{\gamma T} \left[ r - \frac{1}{2} \gamma s^\gamma \right] \quad (13)$$

20 where the argument  $x$  is dropped for notational convenience. Clearly,  $G_\gamma$  measures  
 21 the “small” residuals ( $|r_i(x)| < \gamma$ ) by their squares while the “large” residuals are  
 22 measured by the  $\ell_1$  function. Thus,  $G_\gamma$  is a piecewise quadratic function, and it is  
 23 continuously differentiable in  $\mathfrak{R}^n$ .

24  $G_\gamma$  is composed of a finite number of quadratic functions. In each domain  $D \subseteq \mathfrak{R}^n$   
 25 where  $s^\gamma(x)$  is constant  $G_\gamma$  is equal to a specific quadratic function as seen from the  
 26 above definition. These domains are separated by the following union of hyper-  
 27 planes,

$$B_\gamma = \{x \in \mathfrak{R}^n \mid \exists i : |r_i(x)| = \gamma\}. \quad (14)$$

1 A sign vector  $s$  is  $\gamma$ -feasible at  $x$  if

$$\forall \varepsilon > 0 \exists z \in \mathfrak{R}^n \setminus B_\gamma : \|x - z\| < \varepsilon \wedge s = s^\gamma(z). \quad (15)$$

3 If  $s$  is a  $\gamma$ -feasible sign vector at some point  $x$  then  $Q_s$  is the quadratic function  
4 which equals  $G_\gamma$  on the subset

$$\mathcal{C}_s^\gamma = \text{cl}\{z \in \mathfrak{R}^n | s^\gamma(z) = s\}. \quad (16)$$

6  $\mathcal{C}_s^\gamma$  is called a  $Q$ -subset of  $\mathfrak{R}^n$ . Notice that any  $x \in \mathfrak{R}^n \setminus B_\gamma$  has exactly one corre-  
7 sponding  $Q$ -subset ( $s = s^\gamma(x)$ ), whereas a point  $x \in B_\gamma$  belongs to two or more  
8  $Q$ -subsets. Therefore, we must in general give a sign vector  $s$  in addition to  $x$  in  
9 order to specify which quadratic function we are currently considering as repre-  
10 sentative of  $G_\gamma$ .

11  $Q_s$  can be defined as follows:

$$Q_s(z) = \frac{1}{2\gamma}(z - x)^T (AW_s A^T)(z - x) + G_\gamma'^T(x)(z - x) + G_\gamma(x). \quad (17)$$

13 The gradient of the function  $G_\gamma$  is given by

$$G_\gamma'(x) = A \left[ \frac{1}{\gamma} W_s r + s \right] \quad (18)$$

15 where  $s$  is a  $\gamma$ -feasible sign vector at  $x$ . For  $x \in \mathfrak{R}^n \setminus B_\gamma$ , the Hessian of  $G_\gamma$  exists,  
16 and is given by

$$G_\gamma''(x) = \frac{1}{\gamma} A W_s A^T. \quad (19)$$

18 The set of indices corresponding to “small” residuals

$$\mathcal{A}_\gamma(z) = \{i | 1 \leq i \leq m \wedge |r_i(z)| \leq \gamma\} \quad (20)$$

20 is called the  $\gamma$ -active set at  $z$ . The set of minimizers of  $G_\gamma$  is denoted by  $M_\gamma$ .

21 Interestingly, there exists a simple duality link between the Huber  $M$ -estimation  
22 problem and quadratic programming. More precisely, it can be shown using  
23 Lagrange duality (see e.g., [17]) that the dual of the Huber  $M$ -estimation is the  
24 following quadratic program:

25 [D2]

$$\begin{aligned} \max \quad & -b^T y - \frac{\gamma}{2} y^T y \\ \text{s.t.} \quad & Ay = 0 \\ & -\mathbf{1} \leq y \leq \mathbf{1} \end{aligned}$$

1 Furthermore, the optimal solutions  $x^*$  of [SL1] and its dual are related by the  
 2 identity:

$$y^* = \frac{1}{\gamma} W_s r(x^*), \quad (21)$$

4 where  $s = s^\gamma(x^*)$ . As  $y^*$  is the unique solution to the dual problem (the dual  
 5 problem is a strictly concave maximization problem) we have the following simple  
 6 but important consequences of the duality result.

7 **Lemma 1**  $s^\gamma(x_\gamma)$  is constant for  $x_\gamma \in M_\gamma$ . Furthermore  $r_i(x_\gamma)$  is constant for  $x_\gamma \in M_\gamma$   
 8 if  $s_i^\gamma = 0$ .

9 Following the lemma we use the notation  $s^\gamma(M_\gamma) = s^\gamma(x_\gamma), x_\gamma \in M_\gamma$  as the sign  
 10 vector corresponding to the solution set.

11 Based on the work of Mangasarian and Meyer [16], it can be shown that the point  
 12  $y^*$  defined in (21) is a least norm solution of the linear program [D1] provided that  
 13  $\gamma > 0$  is sufficiently small. Li and Swetits [10] use this result to give a recursive  
 14 procedure to compute the  $\ell_1$  estimator from Huber's  $M$ -estimator.

### 15 2.3 Connections between the $\ell_1$ Estimator and the Huber $M$ -Estimator

16 The purpose of this section is to summarize some key relationships between the  
 17 linear  $\ell_1$  estimator and the Huber  $M$ -estimator. In particular, the solution set of  
 18 the  $M$ -estimation problem allows a description of the solution set of the  $\ell_1$  esti-  
 19 mation problem.

20 Assume  $x_\gamma \in M_\gamma$ , and let  $s = s^\gamma(M_\gamma)$  and  $W = W_s$ . Then  $x_\gamma$  is a solution to the  
 21 following system of linear equations:

$$AW A^T x_\gamma = AW b - \gamma As. \quad (22)$$

23 Now, assume that  $x_\gamma + \delta h$  is a minimizer of  $G_{\gamma-\delta}$  with  $s^\gamma(x_\gamma + \delta h) = s$ . Thus, we  
 24 can write

$$AW A^T (x_\gamma + \delta h) = AW b - (\gamma - \delta) As.$$

26 This implies that  $h$  solves the system

$$AW A^T h = As. \quad (23)$$

28 This system of linear equations is always consistent since it is equivalent to the  
 29 following system:

$$AW A^T h = -\frac{1}{\gamma} AW r(x_\gamma)$$

31 which corresponds to normal equations associated with  $W A^T h = -\frac{1}{\gamma} W r(x_\gamma)$ .

1 Next, we state an important result without proof from [13]. Let  $\mathcal{S}$  denote the set  
 2 of minimizers of [L1] and  $\mathcal{D}_s^0 = \{x | r_i(x) \leq 0, i \in \sigma_-(s) \wedge r_i(x) \geq 0, i \in \sigma_+(s)\}$   
 3 where  $\sigma_+(s) = \{i | s_i = 1\}$  and  $\sigma_-(s) = \{i | s_i = -1\}$ . Let  $\sigma(s)$  denote the comple-  
 4 ment of  $\sigma_-(s) \cup \sigma_+(s)$  with respect to  $\{1, \dots, m\}$ .

5 **Theorem 1** (a) *There exists  $\gamma_0 > 0$  such that  $s^\gamma(M_\gamma)$  is constant for  $0 < \gamma \leq \gamma_0$ .*

6 (b) *For  $0 < \gamma \leq \gamma_0$ , where  $\gamma_0$  is given in (a), let  $s = s^\gamma(M_\gamma)$ , and let  $\mathcal{N}_s$  denote the*  
 7 *orthogonal complement to  $\text{span}\{a_i | s_i = 0\}$ . If  $x_\gamma \in M_\gamma$ , and  $d$  solves (23) then*

$$M_0 \equiv \mathcal{S}$$

9 *where*

$$M_0 = (x_\gamma + \gamma d + \mathcal{N}_s) \cap \mathcal{D}_s^0, \quad (24)$$

11 *and*

$$y^* = \frac{1}{\gamma} W_s r(x_\gamma) + s \quad (25)$$

13 *solves [D1].*

14 The above theorem gives a description of the set of  $\ell_1$  estimators from the set of  
 15  $M$ -estimators for small enough values of  $\gamma$ . We will use this result in our  
 16 description of the Madsen-Nielsen algorithm and the variant of it we will propose.

### 17 3 The Clark-Osborne Continuation Algorithm

18 The Clark-Osborne algorithm is a continuation algorithm which was not origi-  
 19 nally intended as a device for solving the linear  $\ell_1$  estimation problem. Its pre-  
 20 scribed use was to compute the Huber  $M$ -estimator for suitable values of  $\gamma$   
 21 starting from a large enough value so that the  $\gamma$ -active set includes all the indices.  
 22 In otherwords, the Clark-Osborne algorithm begins with a large value of  $\gamma$  to  
 23 mimic the  $\ell_2$  estimator and decreases  $\gamma$  until its desired value by following the  
 24 piecewise linear path of Huber  $M$ -estimators. In this section we give a slightly  
 25 modified version of this algorithm, tailored to compute the  $\ell_1$  estimator.

26 To carry on with a preliminary description of this algorithm we give a new sign  
 27 vector definition:

$$s_{\gamma i}(x) = \begin{cases} -1 & \text{if } r_i(x) < -\gamma \\ 0 & \text{if } |r_i(x)| \leq \gamma \\ 1 & \text{if } r_i(x) > \gamma. \end{cases} \quad (26)$$

29 We will refer to  $s_\gamma$  as an ‘‘extended’’ sign vector. Notice that  $s^\gamma$  and  $s_\gamma$  differ only  
 30 for those residuals that are on the boundary  $B$ . The Clark-Osborne algorithm

1 works with the above definition of a sign vector rather than (12). In this section we  
 2 will refer to the sign vector  $s_\gamma$  associated with the unique Huber  $M$ -estimator as  
 3 the “optimal extended sign vector”. We assume in this section that  $A$  has full rank.

4 The key idea that motivates the Clark-Osborne algorithm is the linear system (23).  
 5 Assume the Huber  $M$ -estimator  $x_\gamma$  is unique and that it is non-degenerate, i. e. , at  
 6 any value of  $\gamma$  the set  $\{i|r_i(x_\gamma) = \gamma\}$  is a singleton. Clark shows that if the  
 7  $M$ -estimator is unique the matrix  $AW A^T$  has full rank, c. f. Lemma 6 of [6]. Since  
 8 there exists a continuum of values of  $\gamma$  for a finite set of possible  $\gamma$ -feasible sign  
 9 vectors  $s_\gamma$ , one can immediately deduce from our analysis of the previous section  
 10 that that  $s_\gamma$  remains constant by intervals. The intervals corresponding to sign  
 11 vectors constitute “segments” of the piecewise linear path of  $M$ -estimators. We  
 12 refer to these as the “sign intervals”.

13 The algorithm consists of following the unique path of  $M$ -estimators using the  
 14 linear system of equations (23). Under the assumption of uniqueness, and the  
 15 nondegeneracy of  $M$ -estimators, the Clark-Osborne algorithm traces the piecewise  
 16 linear segments of this path. They use the nondegeneracy assumption to show that  
 17 when moving from one segment to another, at the change of segment the adjacent  
 18 sign vectors differ by a single entry. Furthermore, the sign vector obtained from  
 19 this single change is the optimal extended sign vector of the next segment.

20 In the rest of this section we will make the Clark-Osborne algorithm mathemat-  
 21 ically precise.

22 The basic algorithm can be formulated as follows:

```

find the  $\ell_2$  estimator
choose initial  $\gamma$ 
repeat
    Compute  $h$  from (23)
    Decrease  $\gamma$  along  $h$ 
until  $\gamma = 0$ .
  
```

24 The  $\ell_2$  estimator is found as the solution  $x_{ls}$  of the linear system:

$$AA^T x = Ab.$$

26 The parameter  $\gamma$  is initialized to  $\max_{j=1,\dots,m} |r_j(x_{ls})|$ . The next step in the algorithm  
 27 is to trace the path of  $M$ -estimators. To do this, one computes the unique solution  
 28  $h$  to the system

$$AW A^T h = As$$

30 where  $s = s_\gamma(x_\gamma)$  and  $W = W(x_\gamma)$  with  $x_\gamma = x_{ls}$  for initialization. Let  
 31  $x_{\gamma-\delta} \equiv x_\gamma + \delta h$  and  $r(\gamma - \delta) \equiv r(x_\gamma) + \delta Ah$ . The algorithm finds the smallest of  
 32  $\delta > 0$  where one of the components of  $r(\gamma - \delta)$  changes status, i.e., where



- 1  $|r_j(\gamma - \delta)| = \gamma - \delta$  for some  $j$ ,  $1 \leq j \leq m$ . More precisely, let  $\{\delta_i\}$   $i = 1, \dots, K$ ,
- 2 with  $\delta_1 < \delta_2 < \dots < \delta_K$ , be the set of points in  $(0, \gamma)$  where  $|r_j(\gamma - \delta)| = \gamma - \delta$
- 3 for some  $j$ . Then  $\gamma$  is replaced with  $\gamma - \delta_1$ ,  $x_\gamma$  is replaced with  $x_\gamma + \delta_1 h$ ,  $s$  is
- 4 updated as  $s_\gamma(x_\gamma)$ , and the loop is repeated.
- 5 We summarize the steps of this algorithm below.

find  $x_{ls}$  from  $AA^T x = Ab$

choose  $\gamma = \max_{j=1, \dots, m} \|r_j(x_{ls})\|$

find  $s = s_\gamma(x_\gamma)$  and  $W = W_s$

**repeat**

compute  $h$  from  $AW A^T h = As$

compute  $d = Ah$

compute  $\delta_i^+ = \frac{\gamma - r_i(x_\gamma)}{1 + d_i}$  for  $i \in \sigma(s) \cup \sigma_+(s)$

compute  $\delta_i^- = \frac{-\gamma - r_i(x_\gamma)}{-1 + d_i}$  for  $i \in \sigma(s) \cup \sigma_-(s)$

find  $\delta = \min_i \{\delta_i^+, \delta_i^-\}$

$x_\gamma \leftarrow x_\gamma + \delta h$

$\gamma \leftarrow \gamma - \delta$

find  $s = s_\gamma(x_\gamma)$  and  $W = W_s$

**until**  $\gamma = 0$

- 7 Notice that the algorithm stops with an  $\ell_1$  estimator (the unique  $\ell_1$  estimator) if
- 8  $\delta = \gamma$  since the uniqueness of the  $M$ -estimator for sufficiently small  $\gamma > 0$  implies
- 9 the uniqueness of the  $\ell_1$  estimator; see [10].

- 10 **Example 1** Consider the example problem with  $r(x_1, x_2) = (3x_1 + 2x_2,$
- 11  $4x_1 - 4, 3x_2 - 3, 2x_1 + 3x_2 - 5, 7.5x_1 + 7x_2 - 20)^T$  from [6]. The  $\ell_1$  estimator cor-
- 12 responding to this problem is  $(1, 1)^T$ . The least squares solution is
- 13  $x_{ls} = (1.0135, 13892)^T$  with  $r(x_{ls}) = (5.8188, 0.0539, 1.1675, 1.1345, 0.2674)$ . We
- 14 choose  $\gamma = 5.8188$  and initialize  $s = (1, 0, 0, 0, 0)^T$ . We solve the system (23) which
- 15 in this case gives

$$\begin{pmatrix} 76.25 & 58.5 \\ 58.5 & 67 \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} = \begin{pmatrix} 3 \\ 2 \end{pmatrix}. \quad (27)$$

- 17 The unique solution is  $h = (0.0498, -0.0136)^T$ . We find  $\delta = 4.3551$ . Therefore,
- 18  $\gamma \leftarrow \gamma - \delta = 1.4637$  with  $x = x_{ls} + \delta h = (1.2304, 1.3298)^T$  and the corresponding
- 19 residual vector  $r = r(x_{ls}) + \delta d = (6.3507, 0.9216, 0.9893, -1.4501, -1.4637)^T$ .
- 20 Notice that the optimal extended sign vector is  $(1, 0, 0, 0, 0)^T$  in the sign interval

1 [1.4637, 5.8188]. Now, we update  $s$  to become  $s = (1, 0, 0, 0, -1)^T$ . We solve the  
 2 linear system (23) again:

$$\begin{pmatrix} 20 & 6 \\ 6 & 18 \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} = \begin{pmatrix} -4.5 \\ -5 \end{pmatrix}. \quad (28)$$

4 The solution is  $h = (-0.1574, -0.2253)^T$ . We compute  $\delta = 1.4637$ . Since  $\delta = \gamma$ , the  
 5 algorithm stops with  $x \leftarrow x + \gamma h = (1, 1)^T$  as the  $\ell_1$  estimator.

6 The finiteness of this algorithm depends on the following property proved by  
 7 Clark and Osborne:

8 **Theorem 2** If  $s_\gamma$  is the optimal extended sign vector for  $\gamma > \bar{\gamma}$  but fails to be optimal  
 9 for  $\gamma < \bar{\gamma}$ , the difference being caused by the size of a single residual  $r_k$ , then the sign  
 10 vector  $s'_\gamma$  with  $\sigma(s') = \sigma(s) \setminus \{k\}$  or  $\sigma(s') = \sigma(s) \cup \{k\}$  is optimal for some  $\gamma < \bar{\gamma}$ .

11 This gives the algorithm a look-ahead ability in that at the change of intervals the  
 12 algorithm knows what the optimal sign vector will be in the next interval. Now,  
 13 since the algorithm moves from one optimal sign vector to another (the adjacent  
 14 one) while decreasing  $\gamma$  (c.f. Theorem 2), and since  $x_\gamma$  is a piecewise linear function  
 15 of  $\gamma$  under the absence of degeneracy (c.f. Theorem 2. 6 of [7]), the algorithm  
 16 never repeats an optimal extended sign vector. As the number of distinct sign  
 17 vectors is finite, the algorithm terminates finitely.

18 However, when the difference alluded to in the theorem above is caused by more  
 19 than one residual we are no longer sure of the optimal extended sign vector in the  
 20 next interval. To overcome this difficulty, Clark and Osborne propose a finite  
 21 partitioning algorithm to find the  $M$ -estimator for a slightly smaller  $\gamma$  value than  
 22 the current one and continue the algorithm from this point. However, the  
 23 expression “slightly smaller value” is numerically ill-defined, and Clark and  
 24 Osborne do not incorporate this finite partitioning algorithm into their imple-  
 25 mentations.

26 An important feature of the Clark-Osborne algorithm is the update of a suitable  
 27 factorization of the symmetric, positive definite matrix  $AWA^T$  at the change of  
 28 sign intervals. Since there is only one entry that changes in the matrix  $W$  at a  
 29 change of interval and that the matrix always retains its positive definiteness, the  
 30 factorization can be updated in a stable and efficient way by means of orthogonal  
 31 transformations; see [7] for details.

32

#### 4 The Madsen-Nielsen Algorithm

33 The Madsen-Nielsen algorithm is essentially an extension of the Clark-Osborne  
 34 algorithm. The main difference between the two is that the Madsen-Nielsen  
 35 algorithm does not require a unique path of  $M$ -estimators and does not stay on  
 36 the path(s) of  $M$ -estimators. Although no analytical result is available to support  
 37 the superiority of the Madsen-Nielsen algorithm over the Clark-Osborne algo-

1 rithm, the former was shown experimentally to be significantly faster than the  
 2 well-known Barrodale-Roberts simplex  $\ell_1$  algorithm. No such experimental result  
 3 is available for the Clark-Osborne algorithm to date.

4 It can be easily shown using the results of Section 2.3 that the  $M$ -estimators form  
 5 a family of piecewise linear paths. The algorithm then consists of the following  
 6 steps. First, an  $M$ -estimator for some initial value of  $\gamma$  is computed. This is done  
 7 using a finite, modified Newton algorithm earlier proposed by Madsen and  
 8 Nielsen [11]. Then, using a solution to (23) the paths of  $M$ -estimators for  
 9 decreasing values of  $\gamma$  are explored. However, unlike the Clark-Osborne algorithm  
 10 the Madsen-Nielsen algorithm never moves to a point where there is a change of  
 11 sign vectors. Instead, the algorithm allows a larger reduction in  $\gamma$  than the nearest  
 12 end point of a sign interval. With the new value of  $\gamma$ , the modified Newton  
 13 algorithm is invoked using a projected initial guess at the  $M$ -estimator for the new  
 14 value of  $\gamma$ . This is repeated until suitable termination criteria are satisfied.

15 Notice that the most critical departure from the Clark-Osborne continuation  
 16 scheme is that the Madsen-Nielsen algorithm leaves the paths of  $M$ -estimators to  
 17 return to them later.

18 The basic algorithm can be formulated as follows:

```

    choose initial  $\gamma$ 
    repeat
      compute an  $M$ -estimator  $x_\gamma$ 
      decrease  $\gamma$ 
    until  $\gamma = 0$ 
  
```

20 The algorithm has three main components: (1) stopping criterion, (2) computa-  
 21 tion of an  $M$ -estimator, (3) decreasing  $\gamma$ . We study these components in the above  
 22 order.

#### 23 *4.1 Stopping Criteria*

24 The original Madsen-Nielsen algorithm in [12] used the same stopping criteria as  
 25 the Clark-Osborne algorithm. Later Madsen, Nielsen and Pinar in [15] use dif-  
 26 ferent stopping criteria which consist of checking the duality gap and comple-  
 27 mentarity as follows. Let  $x_\gamma \in M_\gamma$  for some  $\gamma > 0$  with  $s = s_\gamma(x_\gamma)$  and  $y_\gamma = \frac{1}{\gamma} W_s r(x_\gamma)$ .  
 28 Let  $h$  be a solution to the system  $AW_s A^T h = As$ . Let  $x_0 = x_\gamma + \gamma h$ . The algorithm  
 29 stops with output  $x_0$  if

$$G(x_0) + b^T y_\gamma = 0, \quad (29)$$

31 and

$$s_i r_i(x_0) \geq 0, \quad \forall i \in \sigma_+(s) \cup \sigma_-(s). \quad (30)$$

1 Clearly,  $x_0$  and  $y_\gamma$  that satisfy these criteria are optimal solutions to [L1] and [D1],  
 2 respectively as these criteria are equivalent to the optimality condition (6) in the  $\ell_1$   
 3 estimation problem.

4 The problem with both termination criteria is that there is nothing that guarantees  
 5 that an arbitrary solution  $h$  to (23) satisfies these conditions. Theorem 1 guar-  
 6 antees the existence of such a solution  $h$  for sufficiently small  $\gamma > 0$  under the  
 7 condition that we use the sign vector definition (4) to compute  $s$ . However, no  
 8 information is conveyed in this theorem as to which solution to (23) to compute.  
 9 In the special case where the  $M$ -estimator is unique and  $AW_s A^T$  has full rank ( $A$   
 10 needs to have full rank for this to hold) then the above stopping criteria lead to a  
 11 finite termination argument. For implementation, one usually computes a basic  
 12 solution or a least-norm solution of (23). But, there is no analytical result to  
 13 justify such choices.

#### 14 4.2 Computing an $M$ -estimator

15 The Newton method of Madsen and Nielsen [11] is a modified Newton method  
 16 with a line search procedure. We will refer to this algorithm as the MN algorithm  
 17 for convenience.

18 The MN algorithm consists of inspecting the domains  $\mathcal{C}_s^\gamma$  to find the quadratic  
 19 representation of  $G_\gamma$  where the global minimizer is located. A search direction  $h$  is  
 20 computed by minimizing the quadratic  $Q_s(x)$  where  $s$  is the sign vector of the  
 21 current iterate. More precisely, let  $x$  be the current iterate and  $s = s(x)$  and  
 22  $W = W(x)$ , we consider the system of equations

$$Q_s'' h = -Q_s'(x). \quad (31)$$

24 This system is expressed as

$$(AW A^T)h = -A[Wr(x) + \gamma s]. \quad (32)$$

26 Clearly,  $x + h$  minimizes the quadratic  $Q_s$  for any  $h$  that solves (32). For ease of  
 27 notation let  $C \equiv AW A^T$ . Furthermore, let  $\mathcal{N}(C)$  denote the null space of  $C$ . If  $C$   
 28 has full rank, then  $h$  is the unique solution to (32). The algorithm checks whether  
 29  $x + h \in \mathcal{C}_s^\gamma$ . If the answer is affirmative, the algorithm stops with  $x + h$  as the  
 30 minimizer of  $G_\gamma$ . Otherwise, it proceeds with a piecewise linear one-dimensional  
 31 search along  $h$ . If the system of equations (32) is consistent, a minimum norm  
 32 solution is computed. The algorithm checks whether  $x + h \in \mathcal{C}_s^\gamma$  and stops with  
 33  $x + h$  as the minimizer if the answer is affirmative. Otherwise, the next iterate is  
 34 found by moving to the first kinkpoint  $\alpha_1$  along  $h$ , i.e., the smallest value of  $\alpha$   
 35 where  $s_\gamma(x + \alpha) \neq s_\gamma(x)$ . Notice that if  $h$  is the least norm solution of (32) the  
 36 point  $x + h$  is the orthogonal projection of  $x$  onto the set of minimizers of the  
 37 quadratic  $Q_s$ .

38 If the system is inconsistent a suitable descent direction  $h$  is computed and a  
 39 piecewise linear one-dimensional search along  $h$  is performed. Madsen and

1 Nielsen showed that under the full rank assumption on  $A$  the iteration is finite,  
 2 i.e., after a finite number of iterations we have  $x + h \in \mathcal{C}_s^\gamma$  and therefore,  $x + h$  is a  
 3 minimizer of  $G_\gamma$ .

4 Recently, Chen and Pinar [5] proposed a modification of this algorithm and  
 5 proved finite termination without the full rank assumption on  $A$ . The modified  
 6 algorithm allows any solution of the system (32) to be used as a descent direction  
 7 as long as its norm is bounded by a constant times the norm of the minimum  
 8 norm solution  $h_m$  while the original algorithm is restricted to the use of a least  
 9 norm solution in the consistent case. Furthermore, in this case, the original  
 10 algorithm moved to the first kink point along the search direction whereas the  
 11 modified algorithm prescribes a line search along this direction. With these  
 12 computational enhancements Chen and Pinar proved that the modified MN  
 13 algorithm stops at an  $M$ -estimator after a finite number of iterations. The proof of  
 14 this result is quite involved. Therefore, the interested reader is referred to [5] for  
 15 details.

16

#### 4.3 Reduction of $\gamma$

17 Assume  $\gamma \notin (0, \gamma_0]$  as defined in Theorem 1. Let  $x_\gamma$  be an  $M$ -estimator corre-  
 18 sponding to the present value of  $\gamma$ . Let  $x_{\gamma-\delta} \equiv x_\gamma + \delta h$  and  $r(\gamma - \delta) \equiv r(x_\gamma) + \delta Ah$ .  
 19 The algorithm finds the smallest of  $\delta > 0$  where one of the components of  $r(\gamma - \delta)$   
 20 changes status, i.e., where  $|r_j(\gamma - \delta)| = \gamma - \delta$  for some  $j$ ,  $1 \leq j \leq m$ . More pre-  
 21 cisely, let  $\{\delta_i\}$   $i = 1, \dots, K$ , with  $\delta_1 < \delta_2 < \dots < \delta_K$ , be the set of points in  $(0, \gamma)$   
 22 where  $|r_j(\gamma - \delta)| = \gamma - \delta$  for some  $j$ . Then  $\gamma$  is replaced with  $\gamma - \delta$  where  $\delta > \delta_1$ ,  $x$   
 23 is replaced with  $x_\gamma + \delta h$ ,  $s$  is updated as  $s_\gamma(x)$ , and the modified Newton algorithm  
 24 is invoked with  $x$  as the starting point.

25 Note that there is some flexibility involved in the choice of  $\delta$  in the reduction  
 26 strategy as long as a change of interval is assured. Madsen and Nielsen[12]  
 27 describe a strategy based on inspecting the points of interval change  $\{\delta_i\}$  as in the  
 28 Clark-Osborne and picking  $\delta$  according to some heuristic criteria. Another heu-  
 29 ristic method is described in Madsen, Nielsen and Pinar [15]. The important point  
 30 here is to find a good heuristic that decreases  $\gamma$  neither too fast nor too slowly.  
 31 This is usually problem dependent, but the two heuristics mentioned above seem  
 32 to give good average performances.

33 As in the Clark-Osborne algorithm, the efficiency of the Madsen-Nielsen algo-  
 34 rithm strongly depends on the efficient solution of linear systems (23) and (32).  
 35 Both these systems involve the same symmetric, positive (semi) definite matrix  
 36  $AW A^T$ . However, the modified Newton algorithm may allow more than one  
 37 index to change its sign unlike the Clark-Osborne case. Nielsen [18] describes a  
 38 software package for updating  $LDL^T$  factors of  $AW A^T$  in a stable and efficient  
 39 way within the modified Newton (MN) algorithm. When the  $M$ -estimator has  
 40 been computed using the MN algorithm, the system (23) is solved to check  
 41 optimality and reduce  $\gamma$ . Since the factors of  $AW A^T$  from the last MN iteration  
 42 are available, no update or refactorization is needed at that stage.

1

**5 Further Results**

2 In this section, we give some further results that are useful in the analysis of the  
 3 extension of the Madsen-Nielsen algorithm. We use  $S(M_\gamma)$  to denote the set of all  
 4 distinct extended sign vectors corresponding to the elements of  $M_\gamma$ . That is, for  
 5 any  $x_\gamma \in M_\gamma$   $s_\gamma(x_\gamma) \in S(M_\gamma)$ .

6 The following result is a consequence of the linearity of the problem.

7 **Lemma 2** *If  $S(M_{\gamma_1}) = S(M_{\gamma_2})$  where  $0 < \gamma_2 < \gamma_1$  then  $S(M_\gamma) = S(M_{\gamma_1}) = S(M_{\gamma_2})$  for*  
 8  $\gamma_2 \leq \gamma \leq \gamma_1$ .

9 **Theorem 3** *There exists  $\bar{\gamma}$  such that  $S(M_\gamma)$  are constant for  $\gamma \in (0, \bar{\gamma})$  where*  
 10  $0 < \bar{\gamma} \leq \gamma_0$ .

11 *Proof:* Since  $s^\gamma(M_\gamma)$  remains constant in  $(0, \gamma_0]$  following Theorem 1 and the  
 12 number of different extended sign vectors is finite, the result is a consequence of  
 13 the previous lemma.  $\square$

14 The above result indicates that when  $\gamma$  is sufficiently small, the boundaries of the  
 15 set of  $M$ -estimators also remain constant. In other words, the set of extended sign  
 16 vectors corresponding to  $M$ -estimators remain constant. This property allows us  
 17 to prove the following important result.

18 **Theorem 4** *Let  $\gamma \in (0, \bar{\gamma})$  and  $x_\gamma \in M_\gamma$  with  $s = s_\gamma(x_\gamma)$  and  $W = W_s$ . Then*

$$Wr(x_\gamma + \gamma h) = 0 \quad (33)$$

20 *for any solution  $h$  to (35). Furthermore, if*

$$s_i r_i(x_\gamma + \gamma h) \geq 0, \quad \forall i \in \sigma_+(s) \cup \sigma_-(s) \quad (34)$$

22 *then  $x_\gamma + \gamma h$  solves [L1].*

23 *Proof:* Let  $\gamma \in (0, \bar{\gamma})$  and  $x_\gamma \in M_\gamma$  with  $s = s_\gamma(x_\gamma)$  and  $W = W_s$ . Consider the system

$$(AW \ A^T)h = As. \quad (35)$$

25 This is a consistent system of linear equations as we have shown in Section 2.3. By  
 26 Theorem 3 there exists  $x_\gamma \in M_\gamma$  such that  $s_\gamma(x_\gamma) = s$  for all  $\gamma \in (0, \bar{\gamma})$ . This implies  
 27 that there exists  $h$  that solves (35) such that  $x_\gamma + \delta h \in M_{\gamma-\delta}$  for all  $\delta \in (0, \gamma]$ . A  
 28 consequence of this using the continuity of  $r$  and (5) is that  $x_\gamma + \gamma h$  solves [L1], and  
 29  $Wr(x_\gamma + \gamma h) = 0$ . Since  $h$  can be replaced by  $h + \eta$  in the above identity where  
 30  $\eta \in \mathcal{N}(AW \ A^T)$ , it follows that

$$Wr(x_\gamma + \gamma h) = 0. \quad (36)$$

32 Now, define  $y_\gamma = \frac{1}{\gamma} Wr(x_\gamma) + s$ . It is easy to verify that if (34) holds,  $x_\gamma + \gamma h$  and  $y_\gamma$   
 33 satisfy the complementarity condition. Since  $y_\gamma$  is feasible for [D1], this implies  
 34 that  $x_\gamma + \gamma h$  is an  $\ell_1$  estimator.  $\square$

1 The theorem says that the “small” residuals in the sense of definition (26) are  
 2 approaching zero as  $\gamma$  approaches zero using any solution to (35) provided  $\gamma$  is  
 3 sufficiently small.

4 Under a certain regularity assumption on the  $\ell_1$  problem it is possible to relate the  
 5 magnitude of  $\bar{\gamma}$  to the nonzero optimal residuals magnitudes of the  $\ell_1$  solution.

6 **Theorem 5** Let  $x$  be an  $\ell_1$  estimator with  $s = s(x)$ . If for some solution  $h$  to the  
 7 system

$$AW_s A^T h = As \quad (37)$$

9 we have

$$\|W_s A^T h\|_\infty \leq 1, \quad (38)$$

11 then, there exists  $x_\gamma \in M_\gamma$  with  $s_\gamma(x_\gamma) = s$  for all  $\gamma \in (0, \xi)$  where  $\xi \leq$   
 12  $\min\{|r_i(x)| : i \in \sigma_-(s) \cup \sigma_+(s)\}$ .

13 *Proof:* Let  $s = s(x)$  and  $\delta = \min\{|r_i(x)| : i \in \sigma_-(s) \cup \sigma_+(s)\}$ . The linear system  
 14 (37) is consistent following (6). By the regularity assumption we have  $\|WA^T h\| \leq 1$   
 15 for any solution  $h$  to the system. Choose  $0 < \xi \leq \delta$  so that for all  $0 < \gamma \leq \xi$ ,

$$r_i(x) - \gamma(A^T h)_i > \xi, \quad i \in \sigma_+(s), \quad (39)$$

$$r_i(x) - \gamma(A^T h)_i < -\xi, \quad i \in \sigma_-(s). \quad (40)$$

18 Now using (37) and the fact that  $W_s(A^T x - b) = 0$  we have:

$$\begin{aligned} 0 &= AW_s A^T(-\gamma h) + \gamma As \\ &= AW_s(A^T(x - \gamma h) - b) + \gamma As. \end{aligned}$$

20 Since  $\|W_s A^T h\|_\infty \leq 1$ , using (39) and (40) we have  $s_\gamma(x - \gamma h) = s$ . Hence,  
 21  $x - \gamma h \in M_\gamma$ . □

22 Following this theorem, we can expect to decrease  $\gamma$  to the level of the smallest  
 23 nonzero optimal residual(s) to enter the final sign interval of  $M$ -estimators.

## 24 **6 An Extension of the Madsen-Nielsen Algorithm**

25 Before going into the details of the algorithm to be proposed below, it is  
 26 instructive to examine how the theory developed in Section 5 motivates the  
 27 algorithm.

28 Notice that for  $\gamma$  sufficiently small ( $\gamma \in (0, \bar{\gamma})$ ) the point  $x_\gamma + \gamma h$  gives  
 29  $W(r(x_\gamma + \gamma h)) = 0$  regardless of the choice of  $h$ . Hence, if  $x_\gamma + \gamma h$  is complementary  
 30 to  $y_\gamma \equiv \frac{1}{\gamma} W r(x_\gamma) + s$ , then  $(y_\gamma, x_\gamma + \gamma h)$  is clearly a primal-dual optimal pair. If

- 1  $W(r(x_\gamma + \gamma h) = 0$  but  $x_\gamma + \gamma h$  and  $y_\gamma$  are not complementary, we move to the  
 2 smallest positive point along  $h$  where a change of sign occurs. If  $\gamma \in (0, \bar{\gamma})$  this  
 3 leads to an expansion of the active set. Continuing this way, the algorithm stops in  
 4 a finite calculation. If  $W(r(x_\gamma + \gamma h) \neq 0$ , we reduce  $\gamma$  exactly as in Madsen and  
 5 Nielsen [12] or, as in [15]. As far as the finite termination arguments are concerned  
 6 it suffices that  $\gamma$  is replaced by  $\beta\gamma$  where  $\beta \in (0, 1)$  in this case.
- 7 More precisely, we propose the following algorithm:

Choose  $\gamma$  and compute a minimizer  $x_\gamma$  of  $G_\gamma$ (call MN)  
**while** not STOP  
 find  $s = s_\gamma(x_\gamma)$  and  $W = W_s$   
 compute  $h$  from  $AW A^T h = As$   
**if**  $Wr(x_\gamma + \gamma h) = 0$  **then**  
 compute  $d = Ah$   
 compute  $\delta_i^+ = \frac{\gamma - r_i(x_\gamma)}{1 + d_i}$  for  $i \in \sigma_+(s)$   
 compute  $\delta_i^- = \frac{-\gamma - r_i(x_\gamma)}{d_i - 1}$  for  $i \in \sigma_-(s)$   
 find  $\delta = \min_i \{\delta_i^+, \delta_i^-\}$   
 $x_\gamma \leftarrow x_\gamma + \delta h$   
 $\gamma \leftarrow \gamma - \delta$   
**else**  
 reduce  $\gamma(\gamma \leftarrow \beta\gamma)$   
 $x \leftarrow x_\gamma + (1 - \beta)\gamma h$   
 compute a minimizer  $x_\gamma$  of  $G_\gamma$  starting from  $x$  (call MN)  
**endif**  
**end while.**

- 9 In the above iteration STOP is a function that returns TRUE if

$$s_i r_i(x_\gamma + \gamma h) \geq 0, \forall i \in \sigma_+(s) \cup \sigma_-(s). \quad (41)$$

- 11 Notice that when the condition  $Wr(x_\gamma + \gamma h) = 0$  holds the new algorithm uses a  
 12 strategy similar to the Clark-Osborne algorithm. However, no restriction about  
 13 uniqueness of  $x_\gamma$  nor non-singularity of the matrix  $AW A^T$  is required. Further-  
 14 more, we do not impose any requirements on the solution  $h$  of  $AW A^T h = As$  as  
 15 far as the proof of finite termination is concerned.

- 16 **Example 2** Consider the following example problem from [10]. We have  
 17  $r(x_1, x_2) = (x_1 + 8x_2, x_1 - 8x_2, 2x_2, 17x_2 - 1)^T$ . The  $\ell_1$  estimator corresponding to  
 18 this problem is unique,  $x = (0, 0)^T$ . Interestingly, for  $\gamma \in (0, 4/21]$ , the Huber  
 19  $M$ -estimator is not unique. Let  $\gamma = 3/21$ . It can be verified easily that



1  $x_\gamma = (1/10, 3/84)^T$  is an M-estimator for  $\gamma = 3/21$ , with  $r(x_\gamma) = (0.3857,$   
 2  $-0.1857, 0.0714, -0.3929)^T$ . This corresponds to  $s_\gamma(x_\gamma) = (1, -1, 0, -1)^T$ . We solve  
 3 the system (23) which is in this case:

$$\begin{pmatrix} 0 & 0 \\ 0 & 4 \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} = \begin{pmatrix} 0 \\ -1 \end{pmatrix}. \quad (42)$$

5 The least norm solution of this system is  $h = (0, -1/4)^T$ . This gives  
 6  $d = A^T h = (-2, 2, -0.5, -4.25)^T$ . The point  $r(x_\gamma) + \gamma d$  gives  $W(r(x_\gamma) + \gamma d) = 0$  but  
 7 does not satisfy complementarity criterion (41). Evaluating the kink points  $\{\delta_i^+\}$  and  
 8  $\{\delta_i^-\}$  we find  $\delta = 0.0429$ . Hence,  $\gamma \leftarrow \gamma - \delta = 0.1$ , the algorithm moves to  
 9  $x = x_\gamma + \delta h = (0.1, 0.025)^T$  with  $r = r(x_\gamma) + \delta d = (0.3, -0.1, 0.05, -0.575)^T$ . The  
 10 extended sign vector associated with this point ( $x$  is the Huber M-estimator for  
 11  $\gamma = 0.1$ ) is  $(1, 0, 0, -1)^T$ . We form (23) again:

$$\begin{pmatrix} 1 & -8 \\ -8 & 68 \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} = \begin{pmatrix} 1 \\ -9 \end{pmatrix} \quad (43)$$

13 with the unique solution  $h = (-1, -0.25)$  and  $d = A^T h = (-3, 1, -0.5, -4.25)^T$ .  
 14 This time, the point  $x \leftarrow x + \gamma h = (0, 0)^T$  yields  $r \leftarrow r + \gamma d = (0, 0, 0, -1)$ , which  
 15 satisfies the termination criteria. Thus, the algorithm stops with the unique  $\ell_1$  esti-  
 16 mator in two iterations starting from the Huber M-estimator at  $\gamma = 3/21$ .

## 17 6.1 Finite Convergence

18 In this section we show that the algorithm of Section 6 converges finitely.

19 **Lemma 3** Assume  $\gamma \in (0, \bar{\gamma})$ . Let  $x \in M_\gamma$  with  $s = s_\gamma(x)$ . Let  $h$  solve (35), and  $x_{\text{next}}$  be  
 20 generated by one iteration of the algorithm. Then either

$$x_{\text{next}} \equiv x + \gamma h \in \mathcal{S}$$

22 and the algorithm stops, or

$$x_{\text{next}} \equiv x + \delta h \in M_{\gamma_{\text{next}}},$$

$$\gamma_{\text{next}} = \gamma - \delta$$

25 where  $\delta$  is as defined in the algorithm, and  $\mathcal{A}_{\gamma_{\text{next}}}(x_{\text{next}})$  is an extension of  $\mathcal{A}_\gamma(x)$ .

26 *Proof:* Let  $y = \frac{1}{\gamma} W r(x) + s$ . Clearly  $W r(x + \gamma h) = 0$  from Theorem 4. If  $x + \gamma h$  and  
 27  $y$  are complementary then  $x_{\text{next}} \equiv x + \gamma h$  is a solution to [L1] by Theorem 4 and  
 28 the algorithm stops. Otherwise, Theorem 4 implies that  $\mathcal{A}_\gamma(x) \subseteq \mathcal{A}_0(x + \gamma h)$ .  
 29 Hence, using the definition of  $\delta$ ,

$$\mathcal{A}_{\gamma-\alpha}(x + \alpha h) = \mathcal{A}_\gamma(x)$$

1 for  $\alpha \in [0, \delta)$ . Since there exists  $j \in \{1, \dots, m\} \setminus \mathcal{A}_\gamma(x)$  such that  $|r_j(x + \delta h)| =$   
 2  $\gamma - \delta$ ,  $\mathcal{A}_{\gamma-\alpha}(x + \delta h)$  is an extension of  $\mathcal{A}_\gamma(x)$ . Furthermore  $x + \delta d \in \mathcal{C}_s^{\gamma-\delta}$ .  
 3 Therefore, using the continuity of the gradient  $G'_\gamma$ , (18) and the definition of  $h$ , we  
 4 have

$$G'_\gamma(x) = G'_{\gamma-\delta}(x + \delta h) = 0.$$

6 Thus,  $x_{next}$  minimizes  $G_{\gamma-\delta}$ . □

7 **Theorem 6** *The algorithm defined in Section 6 terminates in a finite number of*  
 8 *iterations with an  $\ell_1$  estimator.*

9 *Proof:* Let  $x \in M_\gamma$  for some  $\gamma > 0$ . Unless the stopping criteria are met and the  
 10 algorithm stops with a primal-dual optimal pair,  $\gamma$  is reduced by a nonzero factor.  
 11 Since the modified Newton iteration of Section 4.2 is a finite process,  $\gamma$  enters the  
 12 range  $(0, \bar{\gamma})$  where  $\bar{\gamma}$  is as defined in Theorem 3 in a finite number of iterations  
 13 unless the algorithm stops. Now assume  $\gamma \in (0, \bar{\gamma})$ . From Lemma 3 either the  
 14 algorithm terminates or the  $\gamma$ -active set  $\mathcal{A}_\gamma$  is expanded. Repeating this argument,  
 15 the algorithm should stop with an  $\ell_1$  estimator since the  $\gamma$ -active set has finite  
 16 cardinality. □

## 17 6.2 Computational Behavior

18 A software system that implements the original algorithm of Madsen, Nielsen and  
 19 Pınar, called LPASL1, was developed in [14], and later modified by the present  
 20 author to include the changes proposed above. In preliminary tests, it was found  
 21 that the additional precautions proposed above for finite convergence did not cause  
 22 a discernible slowdown of the algorithm. Recently, while the present paper was  
 23 under review, Shi and Lukas [20] introduced a new reduced gradient type algorithm  
 24 for the  $\ell_1$  estimation problem, and reported extensive comparative computational  
 25 results with the most important  $\ell_1$  codes available in the public domain and our  
 26 modification of LPASL1. These include the algorithm ACM551 of [1], ACM552 of  
 27 [3], the algorithm AFK of [2] which are considered to be the fastest  $\ell_1$  codes  
 28 available. While Shi and Lukas' new reduced gradient algorithm turns out to be the  
 29 fastest in a wide range of computational tests with randomly generated over-  
 30 determined linear systems with up to 2430 equations and 1215 unknowns, modified  
 31 LPASL1 is quite competitive with the aforementioned well-established codes. We  
 32 give a brief summary of the cpu time results in Table 1 where 10 instances were  
 33 solved for each size, and average cpu seconds reported. Shi and Lukas [20] indicate  
 34 that the problems are degenerate. The sign – indicates that AFK was not able to  
 35 solve all 10 problems. However, modified LPASL1 was not able to solve success-  
 36 fully all 10 of the largest  $2430 \times 1215$  instances, due to numerical difficulties. We  
 37 believe that the reasons for this behavior are related to the degenerate nature of  
 38 some test problems because for non-degenerate  $2430 \times 1215$  instances, LPASL1  
 39 successfully obtained an optimal solution in competitive CPU times; see [20]. This  
 40 point is to be examined in more detail in further research.

**Table 1** Summary Computational Results

Size	ACM551	ACM552	AFK	modified LPASL1
$480 \times 240$	13.4	4.79	209	0.69
$720 \times 360$	60.4	18.4	1193	2.39
$1080 \times 540$	287	87.1	–	9.06
$1620 \times 810$	1299	340	–	30.2

1

## 7 Conclusions

2 We studied the finite computation of the  $\ell_1$  estimator from Huber's  $M$ -estimator.  
3 We have reviewed and extended the contributions of Clark and Osborne, and  
4 Madsen and Nielsen to give a new finite algorithm to compute the  $\ell_1$  estimator  
5 from the Huber  $M$ -estimator. The new method has guaranteed finite termination  
6 property without any restrictive assumptions. In particular, we removed the  
7 assumption of full rank on the matrix  $A$ , an assumption which is also present in  
8 the recent paper by Li and Swetits [10]. This reference also gives a recursive  
9 algorithm of a somewhat different nature than the algorithms of the present paper  
10 to compute the  $\ell_1$  estimator from the  $M$ -estimator based on the least norm  
11 solution of the dual linear program [D1] although no computational experience or  
12 any evidence to the efficiency of this algorithm is reported. We also summarized  
13 promising results of comparative computational tests obtained with the modified  
14 Madsen-Nielsen algorithm.

15

## Acknowledgement

16 The author is indebted to K. Madsen, H.B., Nielsen and B., Chen for collabora-  
17 tion and many useful discussions.

18

## References

- 19 [1] Abdelmalek, N. N.: Algorithm 551, A Fortran subroutine for the  $L_1$  solution of overdetermined  
20 systems of linear equations. *ACM Trans. Math. Software* 6, 229–230 (1980).  
21 [2] Armstrong, R. D., Frome, E. L., Kung, D. S.: A revised simplex algorithm for the absolute  
22 deviation curve fitting problem. *Comm. Statist. B* 8, 175–190 (1979).  
23 [3] Barrodale, I., Roberts, F.: An improved algorithm for discrete  $L_1$  approximation. *SIAM*  
24 *J. Numer. Anal.* 10, 839–848 (1972).  
25 [4] Bloomfield, P., Steiger, W. L.: *Least Absolute Deviations: Theory, Applications and Algorithms*,  
26 Boston: Birkäuser (1983).  
27 [5] Chen, B., Pınar, M.Ç.: On Newton's method for Huber  $M$ -estimation problems in robust linear  
28 regression. *BIT* 38, 674–684 (1998).  
29 [6] Clark, D.: The mathematical structure of Huber's  $M$ -estimator. *SIAM J. on Scientific and*  
30 *Statistical Computing* 6, 209–219 (1985).  
31 [7] Clark, D.I., Osborne, M.R.: Finite algorithms for Huber's  $M$ -estimator. *SIAM J. on Scientific*  
32 *and Statistical Computing* 7, 72–85 (1986).  
33 [8] Coleman, T., Li, Y.: A globally and quadratically convergent algorithm for the linear  $\ell_1$  problem.  
34 *Math. Prog.* 56, 189–222 (1992).  
35 [9] Huber, P.J.: *Robust Statistics*. New York, Wiley (1981).  
36 [10] Li, W., Swetits, J. J.: Linear  $\ell_1$  estimator and Huber  $M$ -estimator. *SIAM J. on Optimization* 8,  
37 457–475 1997.  
38 [11] Madsen, K., Nielsen, H. B.: Finite algorithms for robust linear regression. *BIT* 30, 682–699  
39 (1990).

- 1 [12] Madsen, K., Nielsen, H. B.: A finite smoothing algorithm for linear  $\ell_1$  estimation. *SIAM J. on*  
2 *Optimization* 3, 223–235 (1993).
- 3 [13] Madsen, K., Nielsen, H. B., Pinar, M. Ç.: New characterizations of  $\ell_1$  solutions of  
4 overdetermined systems of linear equations. *Oper. Res. Letters* 16, 159–166 (1994).
- 1 [14] Madsen, K., Nielsen, H. B., Pinar, M. Ç.: User's guide to LPASL1: A Fortran 77 package, based  
2 on a continuation algorithm for dense linear programming. Lyngby, DK, Technical University of  
3 Denmark 1994.
- 4 [15] Madsen, K., Nielsen, H. B., Pinar, M. Ç.: A new finite continuation algorithm for linear  
5 programming. *SIAM J. on Optim.* 6, 600–616 (1996).
- 6 [16] Mangasarian, O.L., Meyer, R. R.: Nonlinear perturbations of linear programs. *SIAM J. on*  
7 *Control Optim.* 25, 583–595 (1979).
- 8 [17] Michelot, C., Bougeard, M. L.: Duality results and proximal solutions of the Huber  $M$ -estimator  
9 problem. *Appl. Math. Optim.* 30, 203–221 (1994).
- 10 [18] Nielsen, H. B.: AAFAC: A Fortran 77 package for solving  $AA^T x - c$ , Technical Report NI-90-11.  
11 Lyngby, DK Technical University of Denmark, 1990.
- 12 [19] Ruzinski, S. A., Olsen, E. T.:  $L_1$  and  $L_\infty$  minimization via a variant of Karmarkar's algorithm,  
13 *IEEE Trans. on Acous. Speech and Sig. Proc.* 37, 245–253 (1989).
- 14 [20] Shi, M., Lukas, M. A.: An  $L_1$  estimation algorithm with degeneracy and linear constraints.  
15 *Computational Statistics and Data Analysis* 39, 35–55 (2002).
- 16 [21] Watson., G. A.: *Approximation Theory and Numerical Methods*. John Wiley & Sons, New York  
17 1980.
- 18 [22] Zhang, Y.: Primal-dual interior point approach to computing  $L_1$  and  $L_\infty$  solutions of  
19 overdetermined linear systems. *J. Optim. Theory Appl.* 77, 323–341 (1993).

Mustafa Ç. Pinar  
Department of Industrial Engineering  
Bilkent University  
06533 Ankara, Turkey  
e-mail: mustafap@bilkent.edu.tr



Instruction to printer	Mark	Examples	
		In the text	In the margin
Character to be corrected	/	Litter to be corrected	e /
Group of characters to be corrected	H	Letters to be corrected	ed H
Several identical characters to be corrected	/	Council for Commission	o ///
Differentiation of several errors in the same paragraph	1 F L J	There are many faults in this line	r / L m / i / a F
Character or word to be deleted	o	Commission and Parliament	o y o H
Character or word to be added	h	A word missing	is h
Superior character required	^	The Court's judgment.	(^) /
Omitted text to be added (see copy)	h	1. January 12. December	h (Out see copy)
Inferior character required	v	H <sub>2</sub> SO <sub>4</sub>	4 /
Change to italic		Ad infinitum	(ital.)
Change italic characters to roman	o	status quo	(rom.)
Change capitals to lower case	o	UNESCO	(l.c.)
Change to capitals or small capitals	= =	Robert Burns, AD 1759-96	(Caps.) (S.C.)
Change to bold face	~~~~	This word needs emphasising!	(bold)
To be letter-spaced		<del>This line is crooked</del>	/
Correct horizontal alignment		<del>This line is crooked</del>	/
Text to be raised or lowered	∩ ∪	This line is uneven	∩ / ∪ /
Text to be aligned (to the left)	⌋	This text is to be aligned	⌋ /
Text to be aligned (to the right)	⌈	This text is to be aligned	⌈ /
Text to be centred	[ ]	This text is to be centred	[ ] /
Take back to previous line	] ]	This hyphen is unnecessary	] /
Text to run on (no new paragraph)	~	... line. No new paragraph here	~ /
Take forward to next line	[ ]	This hyphen is badly placed	[ /
Create new paragraph	⌋ ⌈	... line. A new paragraph should begin here	⌋ / ⌈ /
Close up	o o	A space is wrong here	o /
Equalise space	/	This spacing is very uneven	∩ /
Add space between words	z	A space is missing here	z # /
Reduce space between words	∩	These spaces are too big!	∩ /
Add space between lines	Y #	These lines are too close together	Y # /
Reduce space between lines	↑	These lines are too far apart.	↑ /
Stet (let original text stand)	⋮	This text was corrected in error	⊙
Transpose characters	S	These letters are transposed	S /
Transpose words	∩	These words are transposed	∩ /
Transpose lines	∩ ∪	These lines are transposed	∩ / ∪ /

NB: A correction made in the text must always have a corresponding mark in the margin, otherwise it may be overlooked when the corrections are made. The same marks should be used, where appropriate, by copy-editors marking up copy. Where instructional words are used in marginal marks, e.g. 'ital.', 'bold', etc., they must always be encircled to show that they are not to be printed.

