

COMMENTS ON COMPUTING MINIMUM ABSOLUTE DEVIATIONS REGRESSIONS BY
ITERATIVE LEAST SQUARES REGRESSIONS AND BY LINEAR PROGRAMMING

by

A. RONALD GALLANT

and

THOMAS M. GERIG

Institute of Statistics
Mimeograph Series No. 911
Raleigh - February 1974



COMMENTS ON COMPUTING MINIMUM ABSOLUTE DEVIATIONS REGRESSIONS BY
ITERATIVE LEAST SQUARES REGRESSIONS AND BY LINEAR PROGRAMMING

by

A. Ronald Gallant

and

Thomas M. Gerig ^{1/}

^{1/} Assistant Professor of Statistics and Economics and Assistant
Professor of Statistics, respectively. Department of
Statistics, North Carolina State University, Raleigh, North
Carolina 27607

COMMENTS ON COMPUTING MINIMUM ABSOLUTE DEVIATIONS REGRESSIONS BY
ITERATIVE LEAST SQUARES REGRESSIONS AND BY LINEAR PROGRAMMING

ABSTRACT

This note considers some aspects of the computational problem of fitting the regression model $y_t = \sum_{i=1}^k x_{it}\beta_i + \mu_t$ ($t = 1, 2, \dots, n$) by minimizing the sum of absolute deviations $\sum_{t=1}^n |y_t - \sum_{i=1}^k x_{it}\beta_i|$. The iterative method recently proposed by Schlossmacher (1973) is shown to have undesirable features under certain conditions. The linear programming approach using the simplex method as suggested by Fisher (1961) requires that the simplex tableau contain a submatrix of order n by $2k + n$ which restricts the method to relatively small problems. We show that an n by k matrix and a $2k + n$ vector are sufficient to represent this submatrix. This representation improves the efficiency of the simplex method and allows its use in a large proportion of the problems which occur in applications.

1. INTRODUCTION

Consider the linear model

$$y = X\beta + \mu,$$

where y is the $(n \times 1)$ vector of responses, X the $(n \times k)$ matrix of inputs, β the $(k \times 1)$ vector of unknown parameters, and μ the $(n \times 1)$ vector of unobservable errors. The parameter β is to be estimated by minimizing

$$\delta = \sum_{t=1}^n |y_t - x_t' \beta|.$$

An iterative technique was recently proposed by Schlossmacher (1973) for solving this problem. The method consists of doing weighted regressions using diagonal weights equal to the reciprocal of the residual for each observation from the previous iterative step. When a residual becomes small relative to the others, the weight for the corresponding observation is set to zero. The procedure terminates when residuals from successive steps differ by only a small amount.

We have found this method is unsatisfactory in some cases. An example is given where there exists a serious lack of stability at the answer and another where the procedure converges to the wrong answer.

The objection to the use of linear programming and the simplex method of Dantzig as suggested by Fisher (1961) is the excessive size of the simplex tableau which limits the method to small problems. We show that a submatrix of the simplex tableau may be represented by a

considerably smaller matrix which both reduces storage requirements and improves the efficiency of the simplex method.

We call attention to the paper by Wagner (1959) which gives an alternative method of achieving storage compression by recasting the primal linear program as a bounded variables dual linear program. We think that the reader will find the simplex method with storage compression as suggested in this note to be the simpler approach.

2. COMMENTS ON ITERATIVELY REWEIGHTED LEAST SQUARES

The following example illustrates the lack of stability of Schlossmacher's (1973) iteratively reweighted regression procedure. Suppose we wish to estimate the location, β , of a population from the following data: -1, -1, 0, 0, 2. The starting estimate (the sample mean) is $\hat{\beta}(1) = 0$ yielding associated residuals -1, -1, 0, 0, 2 and weights 1, 1, 0, 0, 1/2. The second step yields an estimate $\hat{\beta}(2) = \frac{2}{5}[(-1)(1) + (-1)(1) + (0)(0) + (0)(0) + (2)(\frac{1}{2})] = -\frac{2}{5}$. Thus, the first step estimate is the desired answer (the median) while the second step is a significant step away from the answer.

To see why this happened let $\hat{\beta}(j)$ be the $(k \times 1)$ estimate from the j^{th} step, $\mathbf{r}(j) = \mathbf{y} - \mathbf{X}\hat{\beta}(j)$, $\mathbf{W} = \text{diag}(w_1, w_2, \dots, w_n)$, $w_i = |r_i(j)|^{-1}$ if $|r_i(j)| > \delta$, $w_i = 0$ if $|r_i(j)| \leq \delta$, and $\mathbf{s}' = (s_1, s_2, \dots, s_n)$, and $s_i = \text{sgn}(r_i(j))$. Then the condition for stopping is

$$\mathbf{r}(j) = \mathbf{r}(j+1)$$

$$\Rightarrow (\mathbf{y} - \mathbf{X}\hat{\beta}(j+1)) - (\mathbf{y} - \mathbf{X}\hat{\beta}(j)) = \mathbf{0}$$

$$\Rightarrow \mathbf{X}[(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}\mathbf{y} - \hat{\beta}(j)] = \mathbf{0}$$

$$\Rightarrow \mathbf{X}(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}[\mathbf{X}'\mathbf{W}\mathbf{y} - \mathbf{X}'\mathbf{W}\mathbf{X}\hat{\beta}(j)] = \mathbf{0}$$

$$\Rightarrow \mathbf{X}(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}[\mathbf{y} - \mathbf{X}\hat{\beta}(j)] = \mathbf{0}$$

$$\Rightarrow \mathbf{X}(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}\mathbf{r}(j) = \mathbf{0}$$

$$\Rightarrow \mathbf{X}'(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{s} = \mathbf{0}$$

$$\Rightarrow \mathbf{X}'\mathbf{s} = \mathbf{0}$$

Thus, at the solution, $\underline{X}'\underline{S} = \underline{0}$ is a necessary condition for the residuals from the "next step" to be identical to those at the answer. No such \underline{S} vector exists for the example given above and can be shown not to exist in many situations. In view of this, it is not surprising that the method exhibits the observed instability.

The following example illustrates an instance when the process converges to the wrong answer. Suppose we wish to estimate the location, β , of a population from the following sample: $-1, 3t, 1$, where $0 < t < \frac{1}{3}$ and $2t$ is less than the test value for setting weights to zero. Then the first step estimate is $\hat{\beta}(1) = \frac{1}{3}[-1 + 3t + 1] = t$ with residuals $(-1-t), 2t, (1-t)$ and associated weights $(1+t)^{-1}, 0, (1-t)^{-1}$. The second step estimate is $\hat{\beta}(2) = [(1+t)^{-1} + (1-t)^{-1}]^{-1}[-(1+t)^{-1} + 0 + (1-t)^{-1}] = t$ with residuals $(-1-t), 2t, (1-t)$. That is, the residuals are unchanged, the procedure has converged but to t not the correct answer $3t$.

3. THE PROBLEM AS A LINEAR PROGRAM

The problem of fitting the regression model $y = X\beta + \mu$ by minimizing the sum of absolute deviations may be put in the form of a linear programming problem as follows. Decompose the vector β into its positive and negative parts β^+ and β^- so that $\beta = \beta^+ - \beta^-$ where $\beta^+ \geq 0$ and $\beta^- \geq 0$. Similarly decompose the residuals $e = y - X\beta$ into $e = e^+ - e^-$ where $e^+ \geq 0$ and $e^- \geq 0$. Then the problem

$$\begin{aligned} \min \quad & \delta = \sum_{t=1}^n |e_t| \\ \text{s.t.} \quad & X\beta + e - y = 0 \end{aligned}$$

is equivalent to the linear program

$$\begin{aligned} \min \quad & \delta = \mathbf{1}'e^+ + \mathbf{1}'e^- \\ \text{s.t.} \quad & X\beta^+ - X\beta^- + e^+ - e^- - y = 0 \\ & \beta^+ \geq 0, \quad \beta^- \geq 0, \quad e^+ \geq 0, \quad e^- \geq 0 \end{aligned}$$

where $\mathbf{1}'$ is a $(1 \times n)$ row vector whose elements are ones. The most convenient form of this linear program for our purposes is

$$\begin{aligned} \max \quad & -\delta = \mathbf{1}'X\beta^+ - \mathbf{1}'X\beta^- - 2\mathbf{1}'e^- - \mathbf{1}'y \\ \text{s.t.} \quad & X\beta^+ - X\beta^- - e^- \leq y \\ & \beta^+ \geq 0, \quad \beta^- \geq 0, \quad e^- \geq 0. \end{aligned}$$

Set $\underline{A} = [\underline{X}'_1 - \underline{X}'_2 - \underline{I}]$ ($n \times 2k + n$), $\underline{b} = \underline{y}$ ($n \times 1$),
 $\underline{c}' = (\underline{1}'\underline{X}_1, -\underline{1}'\underline{X}_2, -2\underline{1}'\underline{1})$ ($1 \times 2k + n$), and $\underline{r}' = (\underline{\beta}^+, \underline{\beta}^-, \underline{e}^-)$
 ($2k + n \times 1$). The matrix formulation of the linear program is

$$\max \quad -\delta = \underline{c}'\underline{r} - \underline{1}'\underline{y}$$

$$\text{s.t.} \quad \underline{A}\underline{r} \leq \underline{b}$$

$$\underline{r} \geq \underline{0}.$$

When the program has been solved to obtain $\hat{\underline{r}}$ and the associated minimum - max($-\delta$) the estimate of the parameter $\underline{\beta}$ is recovered from $\hat{\underline{r}} = (\hat{\underline{\beta}}^+, \hat{\underline{\beta}}^-, \hat{\underline{e}}^-)$ by setting $\hat{\underline{\beta}} = \hat{\underline{\beta}}^+ - \hat{\underline{\beta}}^-$.

4. THE SIMPLEX ALGORITHM AND THE SUBMATRIX \underline{A}

Our discussion of the simplex algorithm is heavily dependent on the tabular representation of the problem and description of the method contained in Chapter III of Owen (1968). It is recommended that the reader have this reference in hand while reading this section.

The problem is represented initially by the tableau

$$\underline{T} = \begin{bmatrix} \underline{A} & -\underline{b} \\ \underline{c}' & d \end{bmatrix}$$

where the entry $d = -\underline{1}'\underline{y}$. The elementary step of the simplex algorithm is a pivot on some non-zero element $a_{\alpha\beta}$ of the submatrix \underline{A} to obtain the tableau

$$\underline{T} = \begin{bmatrix} \bar{\underline{A}} & -\bar{\underline{b}} \\ \bar{\underline{c}}' & \bar{d} \end{bmatrix} .$$

The submatrix $\bar{\underline{A}}$ is obtained from \underline{A} by pivoting on $a_{\alpha\beta}$ as follows:

- i) $\bar{a}_{\alpha'\beta'} = a_{\alpha'\beta'} - a_{\alpha\beta} a_{\alpha\beta}^{-1} a_{\alpha'\beta}$, $\alpha' \neq \alpha, \beta' \neq \beta$,
- ii) $\bar{a}_{\alpha\beta'} = a_{\alpha\beta}^{-1} a_{\alpha\beta'}$, $\beta' \neq \beta$,
- iii) $\bar{a}_{\alpha'\beta} = -a_{\alpha\beta}^{-1} a_{\alpha'\beta}$, $\alpha' \neq \alpha$,
- iv) $\bar{a}_{\alpha\beta} = a_{\alpha\beta}^{-1}$.

The simplex algorithm is a set of rules for choosing these pivots $a_{\alpha\beta}$. These rules are explained in complete detail by Owen (1968). Our concern is with their effect on the submatrix \underline{A} .

Consider the Figure. Initially the submatrix \underline{A} can be represented by a copy of \underline{X} , a permutation vector \underline{j} initialized at $j_\beta = \beta$ ($\beta = 1, 2, \dots, 2k + n$), and these rules:

- i) if $j_\beta \leq k$ then $a_{\alpha\beta} = x_{\alpha j_\beta}$,
- ii) if $k < j_\beta \leq 2k$ then $a_{\alpha\beta} = -x_{\alpha(j_\beta - k)}$,
- iii) if $2k < j_\beta$ then $a_{\alpha\beta} = -1$ if $j_\beta = 2k + \alpha$ and $= 0$ otherwise.

We call this set of rules the $(\underline{X}, \underline{j})$ representation of \underline{A} .

Again consider the Figure. It is quite clear from the three examples given that any number of pivots does not destroy the structure of \underline{A} . There will always be k columns with arbitrary entries, k columns whose entries are their negatives, and n negative elementary vectors. A pivot on $a_{\alpha\beta}$ amounts to no more than an alteration of the entries in k columns of \underline{A} and permuting two columns. Thus, it is possible to obtain a $(\bar{\underline{X}}, \bar{\underline{j}})$ representation of $\bar{\underline{A}}$ via the rules of the previous paragraph by suitably altering $(\underline{X}, \underline{j})$.

Three possibilities must be taken into account when modifying $(\underline{X}, \underline{j})$ to reflect the effect of a pivot on $a_{\alpha\beta}$. These are:

- i) a pivot occurs on $a_{\alpha\beta}$ such that $j_\beta \leq k$,
- ii) a pivot occurs on $a_{\alpha\beta}$ such that $k < j_\beta \leq 2k$,
- iii) a pivot occurs on $a_{\alpha\beta}$ such that $2k < j_\beta$ and $j_\beta = 2k + \alpha$.

In the first case we have $a_{\alpha\beta} = x_{\alpha j_\beta}$. Alter $(\underline{X}, \underline{j})$ as follows:

- i) Pivot on $x_{\alpha j_\beta}$ to obtain \bar{X} using the pivoting rules given above with \underline{X} , replacing \underline{A} conceptually.
- ii) If $j_{\beta'} = j_\beta + k$ set $\bar{j}_{\beta'} = 2k + \alpha$. If $j_{\beta'} = 2k + \alpha$ set $\bar{j}_{\beta'} = j_\beta + k$. Otherwise set $\bar{j}_{\beta'} = j_{\beta'}$.

In the second case we have $a_{\alpha\beta} = -x_{\alpha(j_\beta - k)}$. Alter $(\underline{X}, \underline{j})$ as follows:

- i) Pivot on $x_{\alpha(j_\beta - k)}$ to obtain \bar{X} .
- ii) Change the sign on each element in row α of \bar{X} except $\bar{x}_{\alpha(j_\beta - k)}$. Change the sign on each element in column j_β of \bar{X} except $\bar{x}_{\alpha(j_\beta - k)}$.
- iii) If $j_{\beta'} = j_\beta - k$ set $\bar{j}_{\beta'} = 2k + \alpha$. If $j_{\beta'} = 2k + \alpha$ set $\bar{j}_{\beta'} = j_\beta - k$. Otherwise set $\bar{j}_{\beta'} = j_{\beta'}$.

In the third case we have $a_{\alpha\beta} = -1$ and $j_\beta - 2k = \alpha$. Alter $(\underline{X}, \underline{j})$ as follows:

- i) Do not pivot.
- ii) Change the sign on each element in row α of \underline{X} .
- iii) Do not alter \underline{j} .

Clearly the above alterations of $(\underline{X}, \underline{j})$ to obtain (\bar{X}, \bar{j}) may be performed in place and no additional storage is needed. One word of caution, modify the vectors \underline{b} , \underline{c} and the scalar d of the tableau \underline{T} before altering $(\underline{X}, \underline{j})$.

5. STORAGE AND EFFICIENCY

Assuming double word storage for floating point variables and single word storage for integers, the $(\underline{X}, \underline{j})$ representation of \underline{A} requires $2nk + 2k + n$ words rather than the $8nk + 8k^2 + 2n^2$ words required to store \underline{A} itself. The elimination of the n^2 term extends the simplex method of solving of the minimum absolute deviations estimation problem to a considerable proportion of the regression problems occurring in applications.

A discussion of the relative efficiency of a pivot via the $(\underline{X}, \underline{j})$ representation as opposed to operating directly on \underline{A} is probably irrelevant--for small problems it doesn't matter and for moderate sized problems storage of \underline{A} is infeasible. Nevertheless, in cases (i) and (ii) a pivot on \underline{A} using $(\underline{X}, \underline{j})$ requires nk arithmetic and indexing operations as opposed to $n^2 + 2nk$ using \underline{A} . For case (iii) using $(\underline{X}, \underline{j})$ requires k operations as opposed to $n^2 + 2nk$. The increased overhead required to index $a_{\alpha\beta}$ indirectly via \underline{j} in other portions of the program will be more than offset by these savings.

CONFIGURATIONS OF A SUBMATRIX OF THE SIMPLEX TABLEAU

Initial

$$[\bar{X} \quad -\bar{X} \quad -\bar{I}]$$

Possibilities after Pivoting

$$\begin{bmatrix} x_{11}^* & x'_{12} & -x_{11} & -x'_{12} & -1 & 0' \\ x_{21} & X_{22} & -x_{21} & -X_{22} & 0 & -I \end{bmatrix} \rightarrow \begin{bmatrix} \bar{x}_{11} & \bar{x}'_{12} & -1 & -\bar{x}'_{12} & -\bar{x}_{11} & 0' \\ \bar{x}_{21} & \bar{X}_{22} & 0 & -\bar{X}_{22} & -\bar{x}_{21} & -I \end{bmatrix}$$

$$\begin{bmatrix} x_{11} & x'_{12} & -x_{11}^* & -x'_{12} & -1 & 0' \\ x_{21} & X_{22} & -x_{21} & -X_{22} & 0 & -I \end{bmatrix} \rightarrow \begin{bmatrix} -1 & -\bar{x}'_{12} & -\bar{x}_{11} & \bar{x}'_{12} & \bar{x}_{11} & 0' \\ 0 & \bar{X}_{22} & -\bar{x}_{21} & -\bar{X}_{22} & -\bar{x}_{21} & -I \end{bmatrix}$$

$$\begin{bmatrix} x_{11} & x'_{12} & -x_{11} & -x'_{12} & -1^* & 0' \\ x_{21} & X_{22} & -x_{21} & -X_{22} & 0 & -I \end{bmatrix} \rightarrow \begin{bmatrix} -x_{11} & -x'_{12} & x_{11} & x'_{12} & -1 & 0' \\ x_{21} & X_{22} & -x_{21} & -X_{22} & 0 & -I \end{bmatrix}$$

$$\bar{x}_{11} = x_{11}^{-1} \quad \bar{x}_{21} = -x_{11}^{-1} x_{21} \quad \bar{x}'_{12} = x_{11}^{-1} x'_{12}$$

$$\bar{X}_{22} = X_{22} - x_{21} x_{11}^{-1} x'_{12}$$

REFERENCES

- [1] Fisher, W. D. (1961) "A note on curve fitting with minimum deviations by linear programming," Journal of the American Statistical Association, 56, pp. 359-362.
- [2] Owen, G. (1968) Game Theory. Philadelphia: W. B. Saunders Company.
- [3] Schlossmacher, E. J. (1973) "An iterative technique for absolute deviation curve fitting," Journal of the American Statistical Association, 68, pp. 857-859.
- [4] Wagner, H. M. (1959) "Linear programming techniques for regression analysis," Journal of the American Statistical Association, pp. 206-212.