# REGRESSION BY MINIMUM SUM OF ABSOLUTE ERRORS: A NOTE ON PERFECT MULTICOLLINEARITY

## By Rolando A. Danao*

### 1. Introduction

Consider the multiple linear regression model

$$(1) \qquad y = x\beta + \epsilon$$

where $y$ is the regressand, $x = [x_1, x_2, \ldots, x_k]$ the vector of regressors, $\beta = [\beta_1, \beta_2, \ldots, \beta_k]'$ the vector of unknown coefficients, and $\epsilon$ the stochastic disturbance term. The most widely used method of estimating $\beta$ is by least squares, i.e., by minimizing the sum of squared errors (MSSE). Another method is that of minimizing the sum of absolute errors (MSAE), i.e., the MSAE estimate of $\beta$ is obtained by minimizing $\Sigma_i \, |\epsilon_i|$. Although MSAE estimation was suggested as far back as 1888 by Edgeworth (Bowley, 1928), its use has been limited because of computational difficulties. It was only in the 1950s that articles appeared (Charnes et al., 1955; Wagner, 1959) showing that the MSAE estimator can be obtained as a solution to a linear programming problem.

The MSAE regression problem is stated as follows:

$MSAE$: Minimize $\displaystyle\sum_{i=1}^{n} |\epsilon_i|$

$$\text{s.t.} \quad \sum_{j=1}^{k} x_{ij}\beta_j + \epsilon_i = y_i, \qquad i = 1, 2, \ldots, n.$$

---

* Professor of Economics, University of the Philippines.

where $x_{ij}$ is the $i$th observation on the $j$th regressor. To transform this into the standard form of the linear programming problem, we introduce the positive and negative parts of a variable $v$ of arbitrary sign by letting

$$v^+ = \max\{0, v\}$$
$$v^- = \max\{0, \bar{v}\}$$

Then

$$v = v^+ - v^-$$
$$|v| = v^+ + v^-,$$

and

$$v^+ \geq 0, v^- \geq 0.$$

Thus, the MSAE problem can be restated as follows:

*MSAE-LP1:*

Minimize $\displaystyle\sum_{i=1}^{n} (\epsilon_i^+ + \epsilon_i^-)$

s.t. $\displaystyle\sum_{j=1}^{k} x_{ij} (\beta_j^+ - \beta_j^-) + \epsilon_i^+ - \epsilon_i^- = y_i, \qquad i = 1, 2, \ldots, n$

$$\beta_j^+, \beta_j^-, \epsilon_i^+, \epsilon_i^- \geq 0.$$

If we set

$$u = [1, 1, \ldots, 1]' \text{ (an } n\text{-vector)}$$
$$\beta^+ = [\beta_1^+, \beta_2^+, \ldots, \beta_k^+]'$$
$$\beta^- = [\beta_1^-, \beta_2^-, \ldots, \beta_k^-]'$$
$$\epsilon^+ = [\epsilon_1^+, \epsilon_2^+, \ldots, \epsilon_n^+]'$$
$$\epsilon^- = [\epsilon_1^-, \epsilon_2^-, \ldots, \epsilon_n^-]'$$
$$y = [y_1, y_2, \ldots, y_n]'$$

126

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & & \cdot \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix}$$

then the MSAE problem can be written in matrix form:

MSAE-LP2:                Minimize    $u'\epsilon^+ + u'\epsilon^-$

s.t. $X\beta^+ - X\beta^- + I\epsilon^+ - I\epsilon^- = y$

$$\beta^+, \beta^-, \epsilon^+, \epsilon^- \geq 0.$$

## 2. The Problem of Perfect Multicollinearity

Any linear programming subroutine can solve MSAE-LP2. Narula and Wellington (1977) developed a procedure that is based on an efficient dual simplex algorithm. The authors, however, remark that "unlike the MSSE regression line, the calculations for the MSAE regression line are not affected by linear dependencies among the regressor variables." This statement is inaccurate. The presence of linear dependencies among regressor variables results in multiple estimates of $\beta$ under MSAE estimation; in fact, there are an infinite number of estimates. Furthermore, a similar situation also holds under MSSE estimation when the regressor variables are linearly dependent.

The presence of linear dependencies among the regressor variables is called the problem of extreme or perfect multicollinearity. The normal equations under MSSE estimation of (1) is given by

(2)                $X'X\beta_{MSSE} = X'y$

Perfect multicollinearity makes $X'X$ a singular matrix which makes it impossible to obtain $\hat{\beta}_{MSSE}$ since the ordinary inverse of $X'X$ does not exist. However, the normal equation (2) is always consistent

(Graybill, 1969) and a solution (in fact, an infinite set of solution for $\hat{\beta}_{MSSE}$ can be obtained by using the generalized inverse of $X'$. The general solution is given by

(3) $$\hat{\beta}_{MSSE} = (X'X)^g X' y + [I - (X'X)^g (X'X)] z$$

where $(X'X)^g$ is the generalized inverse of $X'X$, $I$ is the identity matrix, and $z$ is an arbitrary vector (Graybill, 1969).

We now show that in the presence of perfect multicollinearity the estimate of $\beta$ under MSAE estimation is not unique, i.e., MSAE-LP2 has an infinite number of optimal solutions. It is clear that MSAE-LP2 has a feasible solution given by

$$\beta_j^+ = \beta_j^- = 0, \qquad j = 1, 2, \ldots, k;$$

$$\left. \begin{array}{l} \epsilon_i^+ = y_i \\[2mm] \epsilon_i^- = 0 \end{array} \right\} \text{ if } y_i > 0;$$

$$\left. \begin{array}{l} \epsilon_i^+ = 0 \\[2mm] \epsilon_i^- = - y_i \end{array} \right\} \text{ if } y_i < 0.$$

Moreover, the objective function is bounded below by zero; hence, MSAE-LP2 has an optimal solution.

For expositional convenience, consider the case where a column of $X$ is scalar multiple of another column. Let $\beta_1$ and $\beta_2$ be the coefficients of the linearly dependent regressors $x_1$ and $x_2$, respectively. Then in MSAE-LP2, the variables $\beta_1^+, \beta_1^-, \beta_2^+, \beta_2^-$ have zero coefficients in the objective function while the columns associated with them are pairwise linearly dependent. Suppose that $\beta_1^+$ is in the basis of the optimal solution obtained by the simplex algorithm. Then $\beta_1^-, \beta_2^+, \beta_2^-$ cannot be in the basis since this would violate linear independence of the basis vectors. The portion of the optimal tableau (in canonical form) corresponding to these variables would look like the following:

| Basic Variables | ... | $\beta_1^+$ | $\bar{\beta}_1$ | $\beta_2^+$ | $\bar{\beta}_2$ | ... | Right Hand Side |
|---|---|---|---|---|---|---|---|
| $\left(\begin{array}{l}\text{Objective}\\\text{Function Row}\end{array}\right) \rightarrow$ | ... | 0 | 0 | 0 | 0 | ... | $b_0$ |
| | ... | 0 | 0 | 0 | 0 | ... | . |
| | | : | : | : | : | | . |
| $\beta_1^+$ | ... | 1 | 1 | $\alpha$ | $-\alpha$ | ... | $\hat{\beta}_1^+$ |
| | | 0 | 0 | 0 | 0 | | . |
| | | : | : | : | : | | |
| | | . | . | . | . | | |

where $-b_0$ is the optimal value of the objective function. The column vector associated with $\beta_1^+$ is the unit vector since $\beta_1^+$ is in the basis.

The other column vectors follow from the fact that they are scalar multiples of the vector associated with $\beta_1^+$. Note that $\beta_2^+$ is a nonbasic variable whose objective function coefficient is zero in the optimal tableau. This implies that the optimal solution is not unique since a necessary and sufficient condition for the uniqueness of an optimal solution is that the objective function coefficients of the nonbasic variables in the optimal tableau are positive (Simmonard, 1966). Another optimal solution can be obtained by pivoting on $\alpha$ (if $\alpha > 0$) or on $-\alpha$ (if $\alpha < 0$). This would put $\beta_2^+$ or $\bar{\beta}_2$ in the basis, replacing $\beta_1^+$ which now becomes zero. The new optimal solution has $\hat{\beta}_1 = 0$, $\hat{\beta}_2 = \dfrac{\beta_1^+}{\alpha}$ (if $\alpha > 0$) or $\hat{\beta}_2 = \dfrac{\beta_1^+}{\alpha}$ (if $\alpha < 0$). Since the set of optimal solutions is convex, it follows that there are an infinite number of optimal solutions.

*Remarks:* (1) The general case is proved in a similar manner. If rank $X < k$, then in the optimal solution, not all of the $\beta$'s will appear as basic variables since this would violate the linear independence of the basis vectors. Suppose that $\beta_1^*, \beta_2^*, \ldots \beta_{k_1}^*$ (where $\beta_j^* = \beta_j^+$ or $\bar{\beta}_j$) are in the optimal basis. This implies that the associated regressors $x_1, x_2, \ldots x_{k_1}$ are linearly independent. Moreover, there is a nonbasic variable $\beta_\ell^*$ whose associated regressor $x_\ell$ is a linear combination of $x_1, x_2, \ldots k_1$. One can show that

129

the column associated with $\beta_\ell^*$ in the optimal tableau is a linear combination of the columns associated with $\beta_1^*, \ldots, \beta_k^*$ with a zero objective function coefficient. This shows the existence of nonbasic variable with a zero objective function coefficient in the optimal tableau which implies nonuniqueness of the optimal solution.

(2)  In effect, MSAE estimation in the presence of perfect multicollinearity will choose a maximal set of linearly independent regressors (whose number equals the rank of $X$) and drops the other regressors from the equation by setting their coefficient equal to zero. This is also one of the remedies resorted to by researchers when confronted with perfect multicollinearity under MSSE estimation.

*Example*    Consider the following data set:

| y | $x_1$ | $x_2$ |
|---|-------|-------|
| 1 | 1 | 2 |
| 3 | 2 | 4 |
| 2 | 3 | 6 |
| 3 | 4 | 8 |
| 4 | 5 | 10 |

and the regression model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon .$$

Note that $x_2 = 2x_1$. Using the standard simplex algorithm on the MSAE-LP2 of this model, we obtain the optimal tableau shown in Figure 1.

Figure 1

| Basic Variables | $\beta_0^+$ | $\beta_0^-$ | $\beta_1^+$ | $\beta_1^-$ | $\beta_2^+$ | $\beta_2^-$ | $\epsilon_1^+$ | $\epsilon_1^-$ | $\epsilon_2^+$ | $\epsilon_2^-$ | $\epsilon_3^+$ | $\epsilon_3^-$ | $\epsilon_4^+$ | $\epsilon_4^-$ | $\epsilon_5^+$ | $\epsilon_5^-$ | Right Hand Side |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Objective Function Row → | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 2 | 2 | 0 | 2 | 0 | 0 | 2 | −2 |
| $\beta_0^+$ | 1 | −1 | 0 | 0 | 0 | 0 | 4/3 | −4/3 | 0 | 0 | 0 | 0 | −1/3 | 1/3 | 0 | 0 | 1/3 |
| $\epsilon_2^+$ | 0 | 0 | 0 | 0 | 0 | 0 | −2/3 | 2/3 | 1 | −1 | 0 | 0 | −1/3 | 1/3 | 0 | 0 | 4/3 |
| $\beta_2^+$ | 0 | 0 | $\boxed{1/2}$ | −1/2 | 1 | −1 | −1/6 | 1/6 | 0 | 0 | 0 | 0 | 1/6 | −1/6 | 0 | 0 | 1/3 |
| $\epsilon_3^-$ | 0 | 0 | 0 | 0 | 0 | 0 | 1/3 | −1/3 | 0 | 0 | −1 | 1 | 2/3 | −2/3 | 0 | 0 | 1/3 |
| $\epsilon_5^+$ | 0 | 0 | 0 | 0 | 0 | 0 | 1/3 | −1/3 | 0 | 0 | 0 | 0 | −4/3 | 4/3 | 1 | −1 | 1/3 |

131

The optimal solution corresponding to this optimal tableau is given by

$$\hat{\beta}_0 = \frac{1}{3}$$

$$\hat{\beta}_1 = 0$$

$$\hat{\beta}_2 = \frac{1}{3}$$

$$\hat{\epsilon}_1 = 0$$

$$\hat{\epsilon}_2 = \frac{4}{3}$$

$$\hat{\epsilon}_3 = -\frac{1}{3}$$

$$\hat{\epsilon}_4 = 0$$

$$\hat{\epsilon}_5 = \frac{1}{3}$$

where the $\beta_j$'s are the MSAE regression coefficients and the $\hat{\epsilon}_i'$s are the residuals. Another optimal solution can be obtained by pivoting on the element $\frac{1}{2}$ (enclosed in a rectangle) thus putting $\beta_1^+$ into the basis and removing $\beta_2^+$ from the basis. This other optimal solution is given by

$$\hat{\beta}_0 = \frac{1}{3}$$

$$\hat{\beta}_1 = \frac{2}{3}$$

$$\hat{\beta}_2 = 0$$

$$\hat{\epsilon}_1 = 0$$

$$\hat{\epsilon}_2 = \frac{4}{3}$$

$$\hat{\epsilon}_3 = -\frac{1}{3}$$

$$\hat{\epsilon}_4 = 0$$

$$\hat{\epsilon}_5 = \frac{1}{3}$$

resulting in another set of MSAE estimates. Convex combinations of these two optimal solutions are also optimal solutions.

## REFERENCES

Bowley, A. L. (1928), *F. Y. Edgeworth's Contributions to Mathematical Statistics,* London: Royal Statistical Society.

Charnes, A., Cooper, W.W. and Ferguson, R. O. (1955), "Optimal Estimation of Executive Compensation by Linear Programming," *Management Science,* 1: 138-151.

Graybill, F.A. (1969), *Introduction to Matrices with Applications in Statistics,* Belmont, California: Wadsworth Publishing Company.

Narula, S.C. and Wellington, J.F. (1977), "Multiple Linear Regression with Minimum Sum of Absolute Errors," *Applied Statistics,* XXVI: 1: 106-111.

Simmonard, M. (1966), *Linear Programming,* Prentice-Hall, N.J.

Wagner, H.M. (1959), "Linear Programming Techniques for Regression," *Journal of the American Statistical Association.*