# The Uses of Tobit Analysis

John F. McDonald, Robert A. Moffitt

*The Review of Economics and Statistics*, Volume 62, Issue 2 (May, 1980), 318-321.

## THE USES OF TOBIT ANALYSIS

### John F. McDonald and Robert A. Moffitt*

As econometric models with truncated or censored error terms come into increasing, almost routine, use, it is important that the information they provide be used fully and correctly. One of the models that is seeing increasing use is Tobit analysis, a model devised by Tobin (1958) in which it is assumed that the dependent variable has a number of its values clustered at a limiting value, usually zero. For example, data on demand for consumption goods often have values clustered at zero; data on hours of work often have the same clustering. The Tobit technique uses all observations, both those at the limit and those above it, to estimate a regression line, and it is to be preferred, in general, over alternative techniques that estimate a line only with the observations above the limit.

In this paper we point out that the coefficients obtained from using Tobit—here called "beta" coefficients—provide more information than is commonly realized. In particular, we show that Tobit can be used to determine both changes in the probability of being above the limit and changes in the value of the dependent variable if it is already above the limit; and we show that this decomposition can be quantified in rather useful and insightful ways. In addition, we apply the decomposition to several recent journal articles that have used Tobit analysis, and we show the additional information that could have been obtained in these articles—and the errors that could have been avoided—if the decomposition had been used. Thus the paper illustrates an important use of Tobit analysis which could be usefully employed as the model is more widely used.

The decomposition also has important substantive economic and policy implications. For example, we were first led to the problem by the following question: How will the labor-supply reduction induced by a negative income tax be spread between marginal decreases in hours worked and decreases in the probability of working any hours? As it turns out, policymakers in the executive and legislative branches are *not* indifferent as to the composition of the total reduction. It was not obvious how to use our Tobit beta coefficients to answer this question, and a review of the literature revealed that no one else had; hence, this paper.

The first section below briefly states the mathematical relationships involved in the decomposition. The second section applies it to several recent pieces of research. The third section concludes the paper and discusses the generalizability of the results.

### I. The Tobit Model

The stochastic model underlying Tobit may be expressed by the following relationship:

$$y_t = X_t\beta + u_t \quad \text{if } X_t\beta + u_t > 0$$
$$= 0 \quad \text{if } X_t\beta + u_t \leq 0,$$
$$t = 1, 2, \ldots, N, \quad (1)$$

where $N$ is the number of observations, $y_t$ is the dependent variable, $X_t$ is a vector of independent variables, $\beta$ is a vector of unknown coefficients, and $u_t$ is an independently distributed error term assumed to be normal with zero mean and constant variance $\sigma^2$. Thus the model assumes that there is an underlying, stochastic index equal to $(X_t\beta + u_t)$ which is observed only when it is positive, and hence qualifies as an unobserved, latent variable.

As Tobin shows, the expected value of $y$ in the model is

$$Ey = X\beta F(z) + \sigma f(z), \quad (2)$$

where $z = X\beta/\sigma$, $f(z)$ is the unit normal density, and $F(z)$ is the cumulative normal distribution function (individual subscripts are omitted for notational convenience). Furthermore, the expected value of $y$ for observations above the limit, here called $y^*$, is simply $X\beta$ plus the expected value of the truncated normal error term (see, e.g., Amemiya, 1973):

$$Ey^* = E(y|y > 0)$$
$$= E(y|u > -X\beta)$$
$$= X\beta + \sigma f(z)/F(z). \quad (3)$$

Consequently, the basic relationship between the expected value of all observations, $Ey$, the expected value conditional upon being above the limit, $Ey^*$, and the probability of being above the limit, $F(z)$, is

$$Ey = F(z)Ey^*. \quad (4)$$

The decomposition that we have found useful is obtained by considering the effect of a change in the $i^{th}$ variable of $X$ on $y$:

$$\partial Ey/\partial X_i = F(z)(\partial Ey^*/\partial X_i) + Ey^*(\partial F(z)/\partial X_i). \quad (5)$$

Thus the total change in $y$ can be disaggregated into two, very intuitive parts: (1) the change in $y$ of those

above the limit, weighted by the probability of being above the limit; and (2) the change in the probability of being above the limit, weighted by the expected value of $y$ if above. The relative magnitudes of these two quantities is an important indicator with substantive economic implications, as shown in the next section.

Assuming that one has estimates of $\beta$ and $\sigma$ (see below), each of the terms in equation (5) can be evaluated at some value of $X\beta$, usually at the mean of the $X$'s, $\bar{X}$. The value of $Ey^*$ can be calculated from equation (3), and the value of $F(z)$ can be obtained directly from statistical tables. The two partial derivatives are also calculable:

$$\partial F(z)/\partial X_i = f(z)\beta_i/\sigma \qquad (6)$$

and, from equation (3),

$$\begin{aligned}
\partial Ey^*/\partial X_i &= \beta_i + (\sigma/F(z))\partial f(z)/\partial X_i \\
&\quad - (\sigma f(z)/F(z)^2)\partial F(z)/\partial X_i \\
&= \beta_i[1 - zf(z)/F(z) - f(z)^2/F(z)^2], \qquad (7)
\end{aligned}$$

using $F'(z) = f(z)$ and $f'(z) = -zf(z)$ for a unit normal density.

It should be noted from equation (7) that the effect of a change in $X_i$ on $y^*$ is not equal to $\beta_i$. It is a common error in the literature to assume that the Tobit beta coefficients measure the correct regression coefficients for observations above the limit. As can be seen from equation (7), this is true only when $X$ equals infinity, in which case $F(z) = 1$ and $f(z) = 0$. This will of course not hold at the mean of the sample or for any individual observation.

It should also be noted that when equations (6) and (7) are substituted into equation (5), the total effect $\partial Ey/\partial X_i$ can be seen to equal simply $F(z)\beta_i$. Furthermore, by dividing both sides of equation (5) by $F(z)\beta_i$,

it easily can be seen that the fraction of the total effect due to the effect above the limit, $\partial Ey^*/\partial X_i$, is just $[1 - zf(z)/F(z) - f(z)^2/F(z)^2]$. Thus, the information we wish to seek in the decomposition can be obtained by calculating this fraction. In addition, as mentioned in the previous paragraph, this is also the fraction by which the $\beta_i$ coefficients must be adjusted to obtain correct regression effects for observations above the limit.

The discussion thus far has consisted only of a further elaboration of the implications of Tobin's original model. In order to make an empirical application of the model, we must assume that estimates of $\beta$ and $\sigma$ have been obtained. Tobin (1958) and Amemiya (1973) have shown that consistent estimates of these parameters—call them $b$ and $s$—can be obtained with maximum likelihood techniques, or plim$(b) = \beta$ and plim$(s) = \sigma$. Furthermore, we shall assume that large samples have been used to obtain estimates of these parameters in each of the studies examined below. The small-sample properties of expressions such as those in equations (2) to (7) are unknown. Investigation of small-sample properties of models such as the Tobit model would be a fruitful area for future research.

## II. Examples from the Literature

The application of the decomposition is illustrated in table 1 for several articles in the literature, dating from Tobin's original article up to the present time (Keeley et al., 1978). For each study we show the fraction of the sample above the (zero) limit and the fraction of the mean total response (of a change in an independent variable) that is due to response above the limit value, evaluated at the $z$ corresponding to the $F(z)$ of the first column.

TABLE 1.—DECOMPOSITION OF TOBIT EFFECTS IN SEVERAL RECENT STUDIES

| Study | Dependent Variable | Fraction of Sample Above Limit $[F(z)]$ | Fraction of Mean Total Response Due to Response Above Limit $[1 - zf(z)/F(z) - f(z)^2/F(z)^2]$ |
|---|---|---|---|
| Tobin (1958) | Durable goods expenditure divided by disposable income | .75 | .54 |
| Dagenais (1975) | Value of auto purchase | .21 | .23 |
| Keeley et al. (1978)[a] | Hours worked per year: | | |
| | Husbands | .94 | .78 |
| | Wives | .36 | .29 |
| | Female Heads | .56 | .40 |
| Rosen (1976) | Hours worked per year, married women | .44 | .33 |
| | Hours worked per week, married women | .44 | .33 |
| Shishko-Rostker (1976) | Moonlighting hours worked per week, male heads | .15 | .20 |

Notes: For all studies, effects are calculated at the $F(z)$ shown, equal to the fraction of the sample above the limit. This fraction is *not* equal to $F(\bar{z})$, i.e., the cumulative distribution function evaluated at the mean of the $x$'s, since $F$ is a nonlinear function of $z$. Although the latter is to be preferred, most of the studies did not present $\bar{z}$.

[a] The fraction of the sample above the limit was calculated by dividing the number of observations with wages (presumably equal to the number of workers) by the number of observations in the total hours equation.

The first two articles, by Tobin (1958) and Dagenais (1975), are concerned with the demand for consumer goods. In Tobin's study, 75% of the observations had nonzero expenditures on durable goods. Evaluating his data at this point, we can thus say that 54% of the total change in durable-goods expenditure resulting from a change in the independent variables would be generated by marginal changes in the value of (positive) expenditures, whereas 46% would be generated by changes in the probability of spending anything at all. In Dagenais' study, on the other hand, there were fewer nonzero purchases—only 21% of the sample purchased an auto in the given year. Evaluating at this point, we find a correspondingly lower percentage (23%) of any total change that would be due to marginal consumption changes. Most of the response would be due to changes in the probability of purchasing in the first place.

The three following studies concern labor supply. The study of Keeley et al. (1978) measured the effect of an experimental negative income tax, and found a significant labor supply disincentive. At the mean of their sample, we can see from the table that most of the work disincentives for wives (71%) would take the form of reductions in the probability of working, whereas most of the work disincentives for husbands (78%) would take the form of a smaller, still-positive hours level.[1] Rosen (1976) examined married women, and his results also show that, for wives, most of any change (67%) would take the form of changes in the probability of working at all. Surprisingly, the fraction above the limit was the same for hours worked per year and per week. Finally, Shishko and Rostker's (1976) study of hours worked on secondary jobs indicates that, at the mean of their sample, a change in total moonlighting hours would be generated more by a change in the probability of moonlighting at all than by a change in the number of hours worked by those who moonlight. Shishko and Rostker themselves briefly examine this question (p. 307) and incorrectly reason their way to the opposite conclusion, primarily because they were not aware of the decomposition technique presented here. In addition, they err quite severely (p. 303) by stating that the Tobit beta coefficients measure the change in expected moonlighting hours of those who are above the zero limit. As we have shown, this is true only when $X$ equals infinity; at the mean of their sample, the beta

coefficients must be multiplied by 0.20 to obtain coefficients for conditional moonlighting hours.

## III.  Generalizability

The decomposition outlined here disaggregates Tobit effects into (1) effects on the probability of being above zero, and (2) effects conditional upon being above zero. It is readily generalizable to several Tobit extensions, such as an upper limit rather than a lower, a nonzero limit rather than a zero, and a different limit for each observation (as Tobit originally assumed, in fact). It is also applicable to recent, more sophisticated variants, such as models with an unobserved, stochastic limit (Nelson, 1977); those with two limits, an upper *and* a lower (Rosett and Nelson, 1975); and simultaneous equations models (Amemiya, 1974). In addition, the decomposition can be applied in somewhat different form to models of market disequilibrium (Goldfeld and Quandt, 1975) to obtain separate effects on the probability of supply being greater than demand, and so on with obvious extensions.

On the other hand, there are some applications in which the decomposition is not of much interest. For example, the estimation of earnings equations on artificially-created, income-truncated samples (such as those containing only poor families) requires Tobit-type maximum-likelihood estimation (Hausman and Wise, 1976, 1977), but a decomposition is not of much interest because here one *is* interested in the underlying Tobit index, whose expected value is that of the underlying population of interest. Thus the decomposition is not relevant whenever one is indeed interested in the *un*truncated population; in this case the beta coefficients are directly usable.

A model in which both the problems of a limit value for the dependent variable and sample truncation arise has been developed by Heckman (1974, 1976). In this model the labor supply of a woman is limited at zero hours, and a wage rate is observed only for those who work. In the case of the wage rate, we are interested in Tobit beta coefficients if we wish to obtain unbiased estimates of the function describing the wages women can earn in the market. For the labor-supply function, however, the technique pointed out in this paper is relevant. Indeed, Heckman (1976) has used the relationship in equation (3) for labor supply to develop an estimation technique for the model.

---

[1] The decomposition here is erroneous if institutional constraints prevent marginal reductions in hours worked. But the problem is not with the decomposition in this case, but with the underlying Tobit model, which assumes no constraints on hours other than that at zero. The second author is presently adapting the Tobit model to introduce institutional constraints.

## REFERENCES

Amemiya, Takeshi, "Regression Analysis When the Dependent Variable Is Truncated Normal," *Econometrica* 41 (Nov. 1973), 997–1016.

———, "Multivariate Regression and Simultaneous Equations Models When the Dependent Variables Are

Truncated Normal," *Econometrica* 42 (Nov. 1974), 999–1011.

Dagenais, Marcel, "Application of a Threshold Regression Model to Household Purchases of Automobiles," this REVIEW 57 (Aug. 1975), 275–285.

Goldfeld, Stephen, and Richard Quandt, "Estimation in a Disequilibrium Model and the Value of Information," *Journal of Econometrics* 3 (1975), 325–348.

Hausman, Jerry, and David Wise, "The Evaluation of Results from Truncated Samples: The New Jersey Income Maintenance Experiment," *Annals of Economic and Social Measurement* 5 (1976), 421–445.

———, "Social Experimentation, Truncated Distributions, and Efficient Estimation," *Econometrica* 45 (May 1977), 919–938.

Heckman, James, "Shadow Prices, Market Wages, and Labor Supply," *Econometrica* 42 (July 1974), 679–694.

———, "The Common Structure of Statistical Models of Truncation, Sample Selection, and Limited Dependent Variables and a Simple Estimator for Such Models,"

*Annals of Economic and Social Measurement* 5 (1976), 475–492.

Keeley, Michael, Philip Robins, Robert Spiegelman, and Richard West, "The Labor Supply Effects and Costs of Alternative Negative Income Tax Programs," *Journal of Human Resources* 13 (Winter 1978), 3–36.

Nelson, Forrest, "Censored Regression Models with Unobserved, Stochastic Censoring Thresholds," *Journal of Econometrics* 6 (1977), 309–327.

Rosen, Harvey, "Taxes in a Labor Supply Model with Joint Wage-Hours Determination," *Econometrica* 44 (May 1976), 485–507.

Rosett, Richard, and Forrest Nelson, "Estimation of the Two-Limit Probit Regression Model," *Econometrica* 43 (Jan. 1975), 141–146.

Shishko, Robert, and Bernard Rostker, "The Economics of Multiple Job Holding," *American Economic Review* 66 (June 1976), 298–308.

Tobin, James, "Estimation of Relationships for Limited Dependent Variables," *Econometrica* 26 (Jan. 1958), 24–36.

# ACCOUNTING FOR SEASONALITY WITH SPLINE FUNCTIONS

## A. Leslie Robb[*]

Suppose that one wishes to estimate a monthly model from $k$ years of data of the form

$$y(t) = A(s) + B(s)x(t) + e(t)$$
$$t = 1, \ldots, 12*k$$
$$s = 1, \ldots, 12 \qquad (1)$$

where $s$ indexes the month, $t$ indexes the observation, and the parameters $A$ and $B$ are allowed to vary with the month of the year. One could proceed by estimating a separate relation for each month of the year, or, equivalently, by using both slope and intercept monthly dummy variables. This procedure may, however, be unsatisfactory in situations with few years of data or in situations where many slope parameters should be allowed to vary monthly. In these situations, the use of spline functions can save on degrees of freedom and still allow a fair degree of flexibility in the shape of the seasonal pattern. In addition, spline functions allow for testing whether or not the reduction in parameters is warranted.[1]

[1] The formulation presented here owes much to the paper by Suits, Mason and Chan (1978).

The use of spline functions is often best illustrated by an example and that is the approach taken in this note. In the example, the parameter $B$ is allowed to vary with the seasons and $A$ is treated as fixed. As will be clear from the example, however, it is a trivial matter to extend the seasonal variation to $A$ as well. At the end of the paper, a specific example using Canadian unemployment data is given.

Suppose then that $B(s)$ has the pattern given in diagram 1. The points marked $+$ in diagram 1 can be thought of as the estimates of the 12 separate slopes that might be estimated in one of the ways mentioned earlier.[2] What I propose to do is to fit a polynomial to each of the four subdivisions of the year shown in the diagram, and to test for continuity and smoothness at the join points $B$, $C$, $D$ and $E$ ($=A$). In particular, a second degree (quadratic) polynomial is employed. Define $B(s)$ as follows:

$$\begin{aligned}
B(s) = {} & a0 + a1*s + a2*s^2 \\
& + [b0 + b1*(s-3) + b2*(s-3)^2]*D3 \\
& + [c0 + c1*(s-6) + c2*(s-6)^2]*D6 \\
& + [d0 + d1*(s-9) + d2*(s-9)^2]*D9 \qquad (2)
\end{aligned}$$

[2] The $X$'s should be thought of as end-of-month estimates. Thus, $s = 0$ corresponds to the end of the previous year or the beginning of the current one. If data points were to refer to the middle of the month, the subdivisions would be moved to the right by 0.5.