# WORKING
# PAPERS

THE APPLICATION OF TOBIT AND PROBIT

ESTIMATION TO AGGREGATE DATA


Frederick I. Johnson


WORKING PAPER NO. 85


June 1983

---

---

BUREAU OF ECONOMICS
FEDERAL TRADE COMMISSION
WASHINGTON, DC 20580

# TOBIT AND PROBIT WITH AGGREGATE DATA

Limited dependent variable models have been applied with great success to many economic analyses. Unfortunately, these models have one drawback which has limited their wider applicability: in general they can only be applied to disaggregated data. In this paper we explore the aggregation problem in two limited dependent variable models, Tobit and Probit. Our conclusion is remarkably sanguine: we find that, if the explanatory variables are themselves normally distributed, we can use mean aggregate data in Tobit and Probit estimation.

## I. Tobit

### A. With Disaggregated Data

In the Tobit model the distribution of the dependent variable is truncated. (The truncation may be envisioned as either an upper or lower limit. Here we treat it as the former.) Consequently, values of the dependent variable ($y_{it}$) in excess of the truncation point ($v_{it}$) are not observed. The model is characterized as follows:

$$(1) \quad y_{it} = \begin{cases} X_{it}\beta + \varepsilon_{it} & <=> \quad I_{it} = 1 \\ v_{it} & <=> \quad I_{it} = 0 \end{cases}$$

where $I_{it}$ is the switching condition (i.e. whether or not the upper constraint $v_{it}$ is binding) and is defined as:

$$(2) \quad I_{it} = \begin{cases} 1 & <=> \quad X_{it}\beta + \varepsilon_{it} < v_{it} \\ 0 & <=> \quad X_{it}\beta + \varepsilon_{it} \geq v_{it} \end{cases}$$

The subscripts i and t represent individual i at time t. $X_{it}$ is a vector of explanatory variables, with $\beta$ the vector of corresponding parameters, and $\varepsilon_{it}$ is a disturbance term, $\varepsilon_{it} \sim N(0,\sigma^2)$. Equations (1) and (2) can be combined to form a single expression for $y_{it}$, from which its expectation can then be derived.

$$(3) \quad y_{it} = I_{it}(X_{it}\beta + \varepsilon_{it}) + (1-I_{it})v_{it}$$

$$= v_{it} + I_{it}(X_{it}\beta - v_{it}) + I_{it}\varepsilon_{it}$$

$$(4) \quad E(y_{it}) = v_{it} + (X_{it}\beta - v_{it}) \, F(\frac{v_{it}-X_{it}\beta}{\sigma}) - \sigma \, f(\frac{v_{it}-X_{it}\beta}{\sigma})$$

where $f(\cdot)$ is the standard normal density function and $F(\cdot)$ is the normal cumulative distribution function. The parameter $\beta$ and the standard error $\sigma$ can be estimated from equation (5):

$$(5) \quad \min_{\beta \, \sigma} \sum_i \sum_t \{y_{it} - E(y_{it})\}^2$$

B.  With Mean Aggregate Data

Suppose, however, that only mean data are available. That is, we have observations on $y_t$, $X_t$, and $v_t$ (each of which is the mean, at time t, of individuals i). Under what circumstances can we estimate $\beta$ and $\sigma$ from the mean analogs of equations (4) and (5) (the same equations without the subscript i)?

(6)  $E(Y_t) = \frac{1}{n} \sum E(y_{it})$

$$= v_t + \frac{1}{n} \sum (X_{it}\beta - v_t) \ F(\frac{v_{it} - X_{it}\beta}{\sigma}) - \frac{\sigma}{n} \sum f(\frac{v_{it} - X_{it}\beta}{\sigma})$$

$$= v_t + (X_t\beta - v_t) \ F(\frac{v_t - X_t\beta}{\sigma}) - \sigma \ f(\frac{v_t - X_t}{\sigma})$$

if  $(X_{it}\beta - v_{it}) = (X_t\beta - v_t)$, $i = 1, \ldots, n$

In general, if Tobit analysis is to be applied to mean data to estimate $\beta$ and $\sigma$, then each of the individuals at time t must face:

(i)    identical values for the explanatory variables  $(X_{it} = X_t)$

(ii)   identical truncation points    $(v_{it} = v_t)$

Actually, despite the formidable conclusion we have just reached, it turns out that Tobit analysis can still be applied to aggregate data even when individuals are not essentially identical.  That is, if the $X_{it}$ and $v_{it}$ are normally distributed we can still estimate the parameter vector $\beta$, but not the standard error $\sigma$.

Let

$$X_{it}\beta = X_t\beta + \delta_{it}$$

$$v_{it} = v_t + \eta_{it}$$

where $\delta_{it}$ and $\eta_{it}$ are normal with zero mean.  The model of equations (1) and (2) can be rewritten as:

(7)  $Y_{it} = \begin{cases} X_t\beta + \varepsilon_{it} + \delta_{it} & \Leftrightarrow \quad I_{it} = 1 \\ v_t + \eta_{it} & \Leftrightarrow \quad I_{it} = 0 \end{cases}$

$$
(8) \quad I_{it} = \begin{cases} 1 & <=> \quad X_t\beta + \varepsilon_{it} + \delta_{it} < v_t + \eta_{it} \\ \\ 0 & <=> \quad X_t\beta + \varepsilon_{it} + \delta_{it} \geq v_t + \eta_{it} \end{cases}
$$

As before, these equations can be combined to solve for $y_{it}$:

$$
(9) \quad y_{it} = I_{it}(X_t\beta + \varepsilon_{it} + \delta_{it}) + (1-I_{it})(v_t + \eta_{it})
$$

$$
= v_t + I_{it}(X_t\beta - v_t) + I_{it}\gamma_{it} + \eta_{it}
$$

where

$$
\gamma_{it} = \varepsilon_{it} + \delta_{it} - \eta_{it}
$$

and therefore

$$
\gamma_{it} \sim N(0, \tau^2)
$$

As in equation (4) we can take the expectation of $y_{it}$ (noting that $E(\eta_{it})=0$) and obtain:

$$
(10) \quad E(y_{it}) = v_t + (X_t\beta - v_t) \, F\left(\frac{v_t - X_t\beta}{\tau}\right) - \tau \, f\left(\frac{v_t - X_t\beta}{\tau}\right)
$$

Finally, since $E(y_{it})$ has the same value for all i:

$$
(11) \quad E(y_t) = E(y_{it})
$$

Consequently, we can still use Tobit to estimate $\beta$ even if we only have observations on the means $y_t$, $X_t$ and $v_t$. However, the variance estimated, $\tau^2$, is that of the composite term $\gamma_{it}$:

$$
(12) \quad \tau^2 = var(\gamma_{it})
$$

$$
= var(\varepsilon_{it} + \delta_{it} - \eta_{it})
$$

$$
= var \, \varepsilon_{it} + var \, \delta_{it} + var \, \eta_{it} + 2cov(\varepsilon_{it}, \delta_{it})
$$

$$
- 2cov(\varepsilon_{it}, \eta_{it}) - 2cov(\delta_{it}, \eta_{it})
$$

Presumably, $\varepsilon_{it}$ will be independent of either $\delta_{it}$ or $\eta_{it}$, in which case:

(13)     $\tau^2 = \text{var } \varepsilon_{it} + \text{var } \delta_{it} + \text{var } \eta_{it} - 2\text{cov}(\delta_{it}, \eta_{it})$

Furthermore, it is quite possible that either $\delta_{it}$ or $\eta_{it}$ will have zero variance. For example, each individual faces the same truncation point. Then

(14)     $\tau^2 = \text{var } \varepsilon_{it} + \text{var } \delta_{it}$

And, of course, if both $\delta_{it}$ and $\eta_{it}$ have zero variance, then

(15)     $\tau^2 = \sigma^2 \ (\equiv \text{var } \varepsilon_{it})$

One should note that in order to estimate $\beta$ it is necessary that $\tau$ be constant over time (or at least that $\tau$ be a known function of time). This in turn implies that the various constituents of $\tau$, itemized in equation (12), themselves be constant (or well known functions of time).

C.   Applications

Maddala and Nelson (1975) investigate the problem of Tobit analysis with aggregate data in the context of bank interest rates. The upper bound on permissable interest rates is set by Regulation Q. The pattern of interest rates in the absence of this regulation could be inferred by Tobit analysis if disaggregated data on individual banks were available. Maddala and Nelson show that if only data on mean interest rates (unweighted by the varying size of deposits at different banks) are available, then Tobit can still be used to estimate $\beta$ and $\sigma$. In essence, the restrictions they

impose are that all banks face identical values of the explanatory variables ($X_{it}=X_t$) and truncation points ($v_{it}=v_t$).

Johnson (1982) encounters the aggregation problem in estimating crop yields. Sales of the crop are limited by a quota, which suggests the application of Tobit. However, the observations are only on aggregate data: sales (s), area cultivated (a), and quotas (q). Is Tobit appropriate?

If we denote mean yields by $\mu$, we have the following model of observed yields:

$$
(19) \quad y_{it} =
\begin{cases}
\mu_{it} + \eta_{it} & \Longleftrightarrow \quad \mu_{it} + \eta_{it} < v_{it} \\
v_{it} & \Longleftrightarrow \quad \mu_{it} + \eta_{it} \geq v_{it}
\end{cases}
$$

where

$$
y_{it} = s_{it}/a_{it}
$$
$$
v_{it} = q_{it}/a_{it}
$$

While mean yields $\mu_{it}$ may vary from farm to farm, it is reasonable to assume that they are normally distributed about the group mean $\mu_t$. As for the truncation point $v_{it}$ it is not immediately apparent that all growers should face the same constraint. However, it turns out that if all the quotas are assigned efficiently (by which it is meant that all growers produce the last unit at the same marginal cost) then the $v_{it}$ will be identical, at time t. Allowing for errors in calculating and assigning quotas efficiently implies

variance in the $v_{it}$, but again it is reasonable to assume that such error is distributed normally. Consequently, the Tobit model described in equation (7) is appropriate.

## II. Probit

The same conclusions extend also to the Probit model. This model is similar to Tobit, except that now we are concerned only with the expectation of $I_{it}$. As before

$$(16) \quad I_{it} = \begin{cases} 1 & \Longleftrightarrow & X_{it}\beta + \varepsilon_{it} < v_{it} \\ 0 & \Longleftrightarrow & X_{it}\beta + \varepsilon_{it} \geq v_{it} \end{cases}$$

$$(17) \quad E(I_{it}) = \Pr\{\varepsilon_{it} < v_{it} - X_{it}\beta\}$$

$$= F\left(\frac{v_{it} - X_{it}\beta}{\sigma}\right)$$

If observations are only available on the means $X_t$ and $v_t$, and

$$X_{it}\beta = X_t\beta + \delta_{it}$$

$$v_{it} = v_t + \eta_{it}$$

then $I_{it}$ can be written as in equation (8) and

$$(18) \quad E(I_{it}) = \Pr\{\varepsilon_{it} + \delta_{it} - \eta_{it} < v_t - X_t\beta\}$$

$$= F\left(\frac{v_t - X_t\beta}{\tau}\right)$$

and

$$(19) \quad E(I_t) = E(I_{it})$$

Note that while $I_{it}$ can only take on values of 0 or 1, $I_t$ can fall anywhere on the interval $(0,1)$. Furthermore, by the Central Limit Theorem the means $I_t$ are asymplotically normal if the $I_{it}$ are independent and identically distributed (i.e. $X_{it}\beta - v_{it} = X_t\beta - v_t$, for all i at time t). This attribute has led to the use of mean data when the underlying distribution is not known. A prime example is the analysis of pesticide dosages on insect mortality, when all the insects are assumed to be drawn from the same population. The implication of equation (18) is that it is not necessary that the insects all be drawn from the same population: rather, it is sufficient that the distribution of insects from each population be normal, and that this distribution be stable over time.

III. Summary

Tobit and Probit models are formulated for observations on individuals. If the data consist only of observations on the mean, then in general neither Tobit nor Probit is appropriate except in the unusual event that each of the individuals has identical characteristics. However, if these characteristics differ from individual to individual, but the pattern of these characteristics is itself normal, then Tobit and Probit methods can be employed to estimate the mean ($X_t\beta$) but not the variance ($\sigma^2$).

# REFERENCES

Johnson, Frederick I. (1982) "Estimating Supply Despite a
Quota: Brazil's Supply of Sugar," Economic Inquiry
XX(1) pp. 54-71.

Lee, Lung-Fei (1975) "Limited Information Estimation of Some
Switching Regression Models," University of Rochester,
Department of Economics Discussion Paper 75-20.

Maddala, G.S. and Forest D. Nelson (1975) "Specification
Errors in Limited Dependent Variable Models," National
Bureau of Economic Research Working Paper (Preliminary:
not numbered).

Tobin, James (1958) "Estimation of Relationships for Limited
Dependent Variables," Econometrica XXVI, pp. 24-36.