



Munich Personal RePEc Archive

(Re)estimating marginal changes after “truncreg” and “tobit” in Stata

Drichoutis, Andreas

Department of Economics, University of Ioannina

23 March 2011

Online at <https://mpra.ub.uni-muenchen.de/33252/>
MPRA Paper No. 33252, posted 09 Sep 2011 10:46 UTC

(Re)estimating marginal changes after “truncreg” and “tobit” in Stata¹

Andreas C. Drichoutis

Dept. of Economics, University of Ioannina, Greece

1. Introduction

In censored and truncated regression models, coefficient estimates represent the effect of an independent variable on a latent variable. This effect is generally not useful, because the latent variable is unobserved. Thus we are mostly interested in estimating marginal changes which represent the effect of an independent variable on the conditional mean function. A typo in Stata’s help file has contributed to the incorrect estimation of marginal changes in such cases. This has significant implications for empirical practice, and results from some published studies may need to be reconsidered.

Popular software programs, such as Stata and Limdep, include automatic routines that ease the computational burden of estimating marginal changes and corresponding standard errors. Unfortunately, a typo that appears in the help files of some of Stata’s commands could potentially lead some users in choosing the wrong prediction option when estimating marginal changes after the associated commands (i.e., the `truncreg` and `tobit` commands). I should note that this is not a problem with the code but a problem with the help files. Stata has already acknowledged this error in personal communication. The same problem is pervasive in other commands of the censored regression model like `xttobit` (the random effects variant of tobit) and `ivtobit` (tobit model with continuous endogenous regressors). Stata has corrected the help files for `ivtobit` and published an erratum for the reference manuals on their website while an erratum was still pending at time of writing for the `truncreg`, `tobit` and `xttobit` commands. However, this typo has been fixed in the final free update of Stata 11 (as of September 1, 2011) and does not appear in Stata 12.

¹ I would like to thank, without implicating, William Greene for providing a preview version of Limdep ver. 10 and for sorting out several of my queries. All errors and omissions are my own.

The purpose of this note is to alert researchers about the potential flaws of past empirical studies and to avoid similar flaws in future empirical research. Relying on a single program for doing empirical economic analysis is potentially dangerous. It is unsettling to think of how many papers may have misreported the real effects of their models, with significant implications for empirical practice and policy making. I first discuss the formulas for estimating marginal changes in censored and truncated regression models and then pinpoint the error in Stata's help files. Then I provide an empirical illustration in which I compare Stata's output with Limdep's output and my own calculations in Stata.

2. The censored regression model

In the censored regression model, the latent underlying regression is:

$$Y_i^* = \mathbf{b}'\mathbf{X}_i + u_i, \quad (1)$$

while the observed dependent variable is:

$$Y_i = L_i \quad \text{if} \quad Y_i^* \leq L_i \text{ (lower tail censoring),}$$

$$Y_i = U_i \quad \text{if} \quad Y_i^* \geq U_i \text{ (upper tail censoring), and}$$

$$Y_i = Y_i^* = \mathbf{b}'\mathbf{X}_i + u_i \quad \text{if} \quad L_i < Y_i^* < U_i.$$

For the latent variable, the conditional mean function is simply $E(Y_i^* | \mathbf{X}_i) = \mathbf{b}'\mathbf{X}_i$. For an observation randomly drawn from the population, which may or may not be censored, the conditional mean function is (Greene, 2003, 2007):

$$E(Y_i | \mathbf{X}_i) = L_i \Phi(\alpha_i) + U_i (1 - \Phi(\gamma_i)) + (\Phi(\gamma_i) - \Phi(\alpha_i)) \left(\mathbf{b}'\mathbf{X}_i + \sigma \frac{\phi(\alpha_i) - \phi(\gamma_i)}{\Phi(\gamma_i) - \Phi(\alpha_i)} \right), \quad (2)$$

where $\alpha_i = \frac{L_i - \mathbf{b}'\mathbf{X}_i}{\sigma}$ and $\gamma_i = \frac{U_i - \mathbf{b}'\mathbf{X}_i}{\sigma}$. The most familiar censored regression model is the

Tobit model, which arises by setting $U_i = +\infty$ and $L_i = 0$ in (2).

The marginal effect of a continuous variable X_j is given by differentiating (2) with respect to it:

$$ME_{X_j} = \partial E(Y_i | \mathbf{X}_i) / \partial X_{ji} = b_j (\Phi(\gamma_i) - \Phi(\alpha_i)) = b_j \times \text{Prob}(\alpha_i < Y_i^* < \gamma_i) . \quad (3)$$

If X_j is a dummy, then marginal effects are only approximately correct, and discrete changes should be estimated. From (2) we get:

$$DC_{X_j} = E(Y_i | \bar{\mathbf{X}}, X_j = 1) - E(Y_i | \bar{\mathbf{X}}, X_j = 0) . \quad (4)$$

3. The truncated regression model

The truncated regression model can be seen as a special case of the censored regression model where only data from the third group are observed. The latent regression is still given by (1), and the observed dependent variable is:

$$Y_i = Y_i^* = \mathbf{b}'\mathbf{X}_i + u_i \quad \text{if} \quad L_i < Y_i^* < U_i .$$

The conditional mean function in this case is:

$$E(Y_i | \mathbf{X}_i, L_i < Y_i < U_i) = \mathbf{b}'\mathbf{X}_i + \sigma \frac{\phi(\alpha_i) - \phi(\gamma_i)}{\Phi(\gamma_i) - \Phi(\alpha_i)} , \quad (5)$$

where $\alpha_i = \frac{L_i - \mathbf{b}'\mathbf{X}_i}{\sigma}$ and $\gamma_i = \frac{U_i - \mathbf{b}'\mathbf{X}_i}{\sigma}$.

If X_j is continuous, then marginal effects can be derived by differentiating (5):

$$ME_{X_j} = b_j \left[1 + \left(\frac{\alpha_i \phi(\alpha_i) - \gamma_i \phi(\gamma_i)}{\Phi(\gamma_i) - \Phi(\alpha_i)} \right) - \left(\frac{\phi(\alpha_i) - \phi(\gamma_i)}{\Phi(\gamma_i) - \Phi(\alpha_i)} \right)^2 \right] = b_j h(b, \sigma) , \quad (6)$$

where $h(b, \sigma)$ is the scale factor.

If X_j is a dummy, then discrete changes should be estimated as:

$$DC_{X_j} = E(Y_i | \bar{\mathbf{X}}, X_j = 1, L_i < Y_i < U_i) - E(Y_i | \bar{\mathbf{X}}, X_j = 0, L_i < Y_i < U_i) . \quad (7)$$

4. Estimating marginal changes in Stata

To estimate marginal changes in Stata, one needs to use the `mfxx` module (the `mfxx` module has been superseded by `margins` in Stata 11, although the error in the help file remained at time of writing). The help file of Stata indicates that the specified options in `mfxx`, `predict(options)` can be, among others:

`e(a,b)` which calculates $E(xb+u | a < xb+u < b)$, the expected value of $y|x$ conditional on $y|x$ being in the interval (a,b) , meaning that $y|x$ is censored,

or

`ystar(a,b)` which calculates $E(y^*)$, where $y^* = a$ if $xb+u \leq a$, $y^* = b$ if $xb+u \geq b$, and $y^* = xb+u$ otherwise, meaning that y^* is truncated.

However, if one tries to use `e(a,b)` after a censored regression model or `ystar(a,b)` after a truncated regression model, one gets incorrect results. Stata's technical support has acknowledged in personal communication that "the prediction option for '`ystar(a,b)`' corresponds to the case when y^* is censored (not truncated), whereas the prediction option for `e(a,b)` corresponds to the case when $y|x$ is truncated (not censored)."

To illustrate these errors with an empirical application, I use one of Stata's example datasets and similar specifications that appear in the help files for `tobit` and `truncreg`. The dataset is a 1978 automobile dataset which contains information on different car models, their selling price, fuel efficiency (miles per gallon), weight, gear ratio as well as the origin of the car (foreign vs. domestic) etc.

I regress miles per gallon (*mpg*) on a continuous variable (*weight*) and a dummy variable indicating origin of the car (*foreign*). This specification is similar to the one that appears in the example of the viewer window of Stata if one types: `help tobit postestimation##predict`. The only addition is the *foreign* variable. The lower limit is 17 and the upper limit is 24. I used the `nlcom` module in Stata to calculate standard errors for marginal changes of own calculations.

It is obvious from Table 1 that, in the censored regression model, Limdep's output, the `ystar(a,b)` option, and my own calculations coincide. According to Stata's official help file, the `e(a,b)` option should have been used, but it provides a very different set of results. Note that Limdep results were obtained from the preview version of Limdep ver. 10. The effect of the dummy variable *foreign* in Limdep, is slightly different than the corrected output of Stata and my own calculations in Stata. This is because the routine that requests estimation of marginal changes at the means, uses the scaled coefficient, i.e., it treats the dummy as a continuous variable. In essence, the program uses formula (3) instead of (4) and formula (6) instead of (7) which, however, produces only approximately correct results.

For the truncated regression model I use the same dataset and a specification similar to the one that appears in the example of the viewer window of Stata if one types: `help truncreg postestimation##predict`. I regress price of cars on miles per gallon (*mpg*) and origin of the car (*foreign*). The lower limit is 4,000 and the upper limit is 9,500. Similarly, in the truncated regression model, the `e(a,b)` option, Limdep's output², and own calculations coincide. However, according to Stata's help file, the `ystar(a,b)` option rather than the `e(a,b)` option should have been used³. The former evidently produces a different set of results.

5. Conclusions

Estimating marginal changes in Stata after a censored or a truncated regression model has been a black box for many researchers. Unfortunately, a simple typo in the help file of Stata may have caused researchers to use the wrong formulas for calculating marginal changes. This could have had significant impact on policy recommendations if Stata's users ignored the mathematics (which appear correct) but not the labels in the help file. If this is the case, many studies may

² Users of Limdep ver. 9 may also notice that standard errors in the case of the truncated regression are different than the ones reported in Table 1 that were obtained with ver. 10. This is due to a simplification of an extremely complicated derivative that was used by the program in version 9. Limdep ver. 10 now uses the full correct derivative and standard errors will coincide with output of Stata, as evident in Table 1.

³ In the case of a truncated regression the results obtained with the wrong option i.e., the `ystar(a,b)` option, are meaningless. In contrast, in a censored regression model the results obtained with the wrong option, i.e., the `e(a,b)` option, correspond to the marginal effects on the expected value of the dependent variable conditional on being uncensored.

need to publish errata of their empirical findings after re-estimating marginal changes for truncated or censored regression models.

6. References

Greene, W.H. 2003. *Econometric Analysis (5th ed.)*. Prentice Hall: New Jersey.

Greene, W.H. 2007. *Limdep version 9.0, Econometric modeling guide (Vol. 1)*. Econometric Software Inc.: Plainview, New York.

Table 1. Estimation results using Stata and Limdep

		Stata's output				Limdep's output		Own calculations in Stata	
		<i>e(a,b)</i>		<i>ystar(a,b)</i>					
Variables		Coef.	Std. Error	Coef.	Std. Error	Coef.	Std. Error	Coef.	Std. Error
Censored regression (lower limit=17, upper limit=24, Dependent variable is <i>mpg</i>)	<i>weight</i>	-0.328	0.059	-0.582	0.066	-0.582	0.066	-0.582	0.066
	<i>foreign</i>	-1.541	0.492	-2.706	0.780	-2.846	0.914	-2.706	0.780
Truncated regression (lower limit=4000, upper limit=9500, Dependent variable is <i>price</i>)	<i>mpg</i>	-78.901	31.836	-89.264	60.430	-78.905	31.836	-78.901	31.836
	<i>foreign</i>	1119.988	439.570	1541.508	676.680	876.649	283.876	1119.978	439.570