

Censored variables and censored regression

Karl G Jöreskog

(3 December 2002, updated August 2021)

A censored variable has a large fraction of observations at the minimum or maximum. Because the censored variable is not observed over its entire range ordinary estimates of the mean and variance of a censored variable will be biased. Ordinary least squares (OLS) estimates of its regression on a set of explanatory variables will also be biased. These estimates are not consistent, i.e., the bias does not become smaller when the sample size increases. This note explains how maximum likelihood estimates can be obtained using PRELIS. The maximum likelihood estimates are consistent, i.e., the bias is small in large samples.

Examples of censored variables are

Econometrics The first example of censored regression appears to be that of Tobin (1958). This is a study of the demand for capital goods such as automobiles or major household appliances. Households are asked whether they purchased such a capital good in the last 12 months. Many households report zero expenditures. However, among those households that made such an expenditure, there will be a wide variation in the amount of money spent.

Greene (2000) p. 205 lists several other examples of censored variables:

1. The number of extramarital affairs
2. The number of hours worked by a woman in the labor force (Quester & Greene, 1982)
3. The number of arrests after release from prison (Witte, 1980)
4. Vacation expenditures (Melenberg & van Soest, 1996)

Biomedicine or Epidemiology Censored variables are common in biomedical, epidemiological, survival and duration studies. For example, in a five-year follow-up study, time to death or time to recovery after surgery, medical treatment or diagnosis, are censored variables if, after five years, many patients are still alive or not yet recovered.

Education Testing If a test is too easy or too difficult there will be a large number of examinees with all items or no items correctly answered.

In econometric dependent censored variables are often called limited dependent variables and censored regression is sometimes called the tobit model (This model was first discussed by Tobin (1958). Goldberger (1964, p. 253) gave it this nickname in analogy with the probit model).

1. Censored normal variables

A censored variable can be defined as follows. Let y^* be normally distributed with mean μ and variance σ^2 . An observed variable y is censored below if

$$y = c \text{ if } y^* \leq c$$

$$= y^* \text{ otherwise,}$$

where c is a constant. This is illustrated in Figure 1.

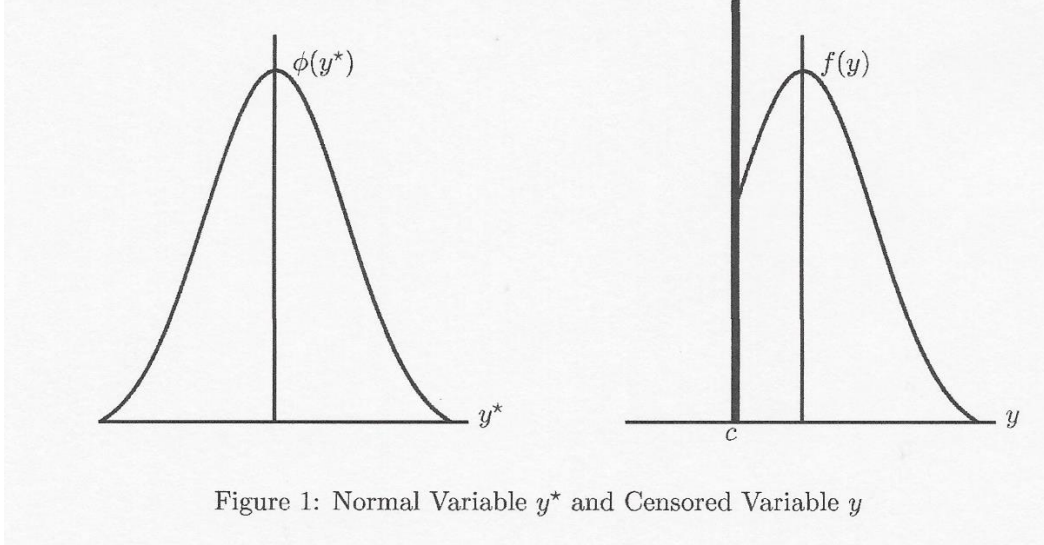


Figure 1: Normal Variable y^* and Censored Variable y

Let ϕ and Φ be the density and distribution functions of the standard normal distribution. The density function of y is

$$f(y) = \left[\Phi\left(\frac{c-\mu}{\sigma}\right) \right]^j \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{y-\mu}{\sigma}\right)^2} \right]^{1-j}, \quad (1)$$

where $j = 1$ if $y = c$ and $j = 0$, otherwise. This may be regarded as a mixture of a binary and a normal variable.

The mean and variance of y are (see, e.g., Greene, 2000, p. 907)

$$E(y) = \pi c + (1 - \pi)(\mu + \lambda\sigma), \quad (2)$$

$$\text{Var}(y) = (1 - \pi) \left[(1 - \delta) + (\alpha - \lambda)^2 \pi \right] \sigma^2, \quad (3)$$

where

$$\alpha = \frac{c - \mu}{\sigma}, \quad (4)$$

$$\pi = \Phi(\alpha), \quad (5)$$

$$\lambda = \frac{\phi(\alpha)}{1 - \phi(\alpha)}, \quad (6)$$

$$\delta = \lambda^2 - \lambda\alpha. \quad (7)$$

A consequence of (2) and (3) is that the sample mean and variance of y are not consistent estimates of μ and σ^2 . The bias of the mean $E(y) - \mu$ as a function of c is shown in Figure 2 for $\mu = 0$ and $\sigma = 1$.

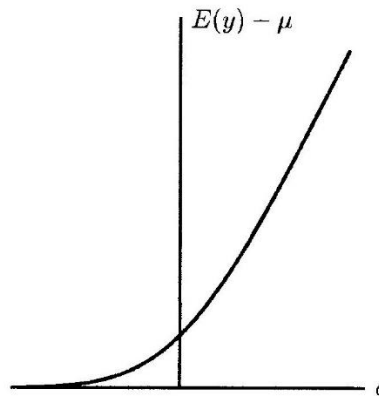


Figure 2: Bias $E(y) - \mu$ as a function of c

An observed variable y is censored above if

$$y = c \text{ if } y^* \geq c$$

$$= y^* \text{ otherwise,}$$

A variable can be censored both below and above. In all three cases, the mean μ and variance σ^2 can be estimated by maximum likelihood as described in the Appendix.

2. Censored normal regression

Consider estimating the regression equation

$$y^* = \alpha + \boldsymbol{\gamma}' \mathbf{x} + z, \tag{8}$$

where α is an intercept term and $\boldsymbol{\gamma}$ is a vector of regression coefficients on the explanatory variables \mathbf{x} . The error term z is assumed to be normally distributed with mean 0 and variance ψ^2 . If y^* is observed as y throughout its entire range, the estimation of (8) is straightforward. However, if the observed variable y is censored below or above, then ordinary least squares (OLS) estimates of y on \mathbf{x} are biased. However, α and $\boldsymbol{\gamma}$ can be estimated by maximum likelihood as described in the Appendix and these maximum likelihood estimates are unbiased in large samples.

3. PRELIS implementation

I illustrate the case of 3 censored variables and 4 explanatory variables. Let Y_1 Y_2 Y_3 be the names of the censored variables and let X_1 X_2 X_3 X_4 be the name of the regressors.

Censored regression of y_1 on x_1, x_2, x_3, x_4 is obtained by the PRELIS command

```
CR Y1 on X1 X2 X3 X4
```

One can select any subset of y -variables and any subset of x -variables to be included in the equation. Thus, one can obtain the regression for all the censored variables simultaneously. For example, the command

```
CR Y1 Y2 Y3 on X1 X2 X3 X4
```

will estimate three regression equations, namely the regression equation of each y_i on all x_j . Note the word on (or ON) separating the censored variables from the regressors.

One can have several CR commands in the same input file. For example,

```
CR Y1 on X1  
CR Y1 on X1 X2  
CR Y1 on X1 X2 X3  
CR Y1 on X1 X2 X3 X4
```

will introduce one regressor at a time in the order x_1, x_2, x_3, x_4 .

General rules:

- All y and x -variables appearing on CR lines must be declared continuous before the first CR command, or else they must have at least 16 different values.
- A censored regression equation can only be estimated from raw data. If there are missing values in the data, the estimation will be based on all cases without missing values on the variables included in the regression equation. Thus, the number of cases used depends on the variables included in the equation. Alternatively, one can impute missing values by multiple imputation before estimating the regression equation. (See du Toit & du Toit, (2001, pp. 165-170). Note that this assumes multivariate normality and missingness at random).
- If several regression equations are estimated, the regression residuals in each equation are assumed to be uncorrelated.

4. Examples

I give four examples of censored regression. The starting point for each of these is a PRELIS system file (PSF file), see du Toit & du Toit (2001, pp. 384-385). You can create a PSF file from data in other formats such as Excel, SPSS and SAS. If the data is in text (ASCII) format it can be read into the PSF file directly from the Windows interface, or alternatively the PSF file can be created by running a simple PRELIS command file of the form

```
da ni = number-of-variables  
la  
labels-of-variables  
ra=filename for the data in text (ASCII) form  
ou ra=filename.PSF
```

If the data requires a variable format statement, include the format as the first line(s) in the ASCII data file. In the examples to follow, LSF files are used. Note that from LISREL 11 PSF files are no longer used. As a result, the zip file accompanying this example contains the necessary data in LSF format.

4.1 Examples 1 and 2

I begin with two small examples based on generated data. Both of them consists of two variables y and x and a sample of 10000 observations. In the first example, y is censored below and in the second example y is censored both below and above. The data are in files **cr1.lsf** and **cr2.lsf**, respectively. Both of these were generated from the true regression line $E(y^* | x) = 5 + x$. In the first example, y was censored below at 3 and in the second example y was censored above at 6, in addition.

To estimate the censored regression of y on z for the first example, use the following PRELIS command file (**cr1.pr1**):

```
Censored Regression: Example 1
sy=cr1.lsf
cr Y on X
ou
```

The output file reveals the following.

Total Sample Size(N) = 10000

Univariate Summary Statistics for Continuous Variables

Variable	Mean	St. Dev.	Skewness	Kurtosis	Minimum	Freq.	Maximum	Freq.
Y	5.654	2.689	0.745	-0.439	3.000	3041	16.372	1
X	-0.008	2.913	-0.010	-1.222	-4.999	1	4.999	1

Test of Univariate Normality for Continuous Variables

Variable	Skewness		Kurtosis		Skewness and Kurtosis	
	Z-Score	P-Value	Z-Score	P-Value	Chi-Square	P-Value
Y	27.295	0.000	-8.957	0.000	825.252	0.000
X	-0.402	0.687	-24.942	0.000	622.268	0.000

Variable Y is censored below

It has 3041 (30.41%) values = 3.000

Estimated Mean and Standard Deviation based on 10000 complete cases.

Mean = 4.963 (0.018)

Standard Deviation = 3.609 (0.003)

The table Univariate Summary Statistics for Continuous Variables shows that the smallest value of y in the sample is 3.000 and this occurs 3041 times. A consequence of this is that y is highly non-normal. This table also gives the mean and standard deviation of y as 5.654 and 2.689, respectively. These are *wrong* values. Taking censoring into account gives the maximum likelihood estimates of the mean and standard deviation of y as 4.963 and 3.609, respectively. If we take the standard error estimates of these estimates into account, it is clear that the ordinary mean and standard deviations are far outside of the 95% confidence interval. This demonstrates that the ordinary mean and standard deviation can be considerably biased in a variable is highly censored.

Test of Univariate Normality for Continuous Variables

Variable	Skewness		Kurtosis		Skewness and Kurtosis	
	Z-Score	P-Value	Z-Score	P-Value	Chi-Square	P-Value
GENDER	0.977	0.329	49.242	0.000	2425.729	0.000
AGE	7.815	0.000	1.159	0.246	62.418	0.000
YEARS	0.789	0.430	63.087	0.000	3980.631	0.000
CHILDREN	-8.278	0.000	-13.975	0.000	263.834	0.000
RELIGIOU	-0.898	0.369	-11.361	0.000	129.886	0.000
EDUCATIO	-2.492	0.013	-1.712	0.087	9.141	0.010
OCCUPATI	-6.741	0.000	-6.554	0.000	88.397	0.000
HAPPINES	-7.443	0.000	-1.058	0.290	56.515	0.000
AFFAIRS	14.813	0.000	8.018	0.000	283.709	0.000

Variable AFFAIRS is censored below
 It has 451 (75.04%) values = 0.000
 Estimated Mean and Standard Deviation based on 601 complete cases.
 Mean = -6.269 (0.056)
 Standard Deviation = 9.420 (0.007)

Estimated Censored Regression based on 601 complete cases.
 $AFFAIRS = 7.609 + 0.946 * GENDER - 0.193 * AGE + 0.533 * YEARS + 1.019 * CHILDREN$
 Standerr (3.936) (1.071) (0.0816) (0.148) (1.289)
 z-values 1.933 0.883 -2.362 3.610 0.790
 P-values 0.053 0.377 0.018 0.000 0.429

- 1.699*RELIGIOU + 0.0254*EDUCATIO + 0.213*OCCUPATI
 (0.409) (0.229) (0.324)
 -4.159 0.111 0.658
 0.000 0.912 0.510

- 2.273*HAPPINES + Error, R² = 0.220
 (0.419)
 -5.431
 0.000

Residual Covariance Matrix

	AFFAIRS
AFFAIRS	69.239

As judged by the *t*-values, the effects of GENDER, CHILDREN, EDUCATION, and OCCUPATION are not statistically significant. The number of extramarital affairs seem to increase with number of years of marriage, and decrease when age, religiousness, and happiness increase.

To illustrate the concept of a fit file, I consider entering one variable at a time in the regression equation. The order of variables corresponds to the size of *t*-values in the previous run. The PRELIS command file is **affairs2.prl**.

Sequential Censored Regression of Affairs
 sy=affairs2.lsf
 cr AFFAIRS on HAPPINESS

```

cr AFFAIRS on HAPPINESS RELIGIOUS
cr AFFAIRS on HAPPINESS RELIGIOUS YEARS
cr AFFAIRS on HAPPINESS RELIGIOUS YEARS AGE
cr AFFAIRS on HAPPINESS RELIGIOUS YEARS AGE GENDER
cr AFFAIRS on HAPPINESS RELIGIOUS YEARS AGE GENDER CHILDREN
cr AFFAIRS on HAPPINESS RELIGIOUS YEARS AGE GENDER CHILDREN OCCUPATION
cr AFFAIRS on HAPPINESS RELIGIOUS YEARS AGE GENDER CHILDREN OCCUPATION EDUCATION
ou xu

```

The xu on the ou line tells PRELIS to skip the results of the univariate data screening in the output. For each regression equation estimated, PRELIS produces a fit file with the same name as the PRELIS command file but with suffix **.FIT**. Abbreviated output showing the estimated regression equations are given below.

Total Sample Size(N) = 601

Variable AFFAIRS is censored below
It has 451 (75.04%) values = 0.000
Estimated Mean and Standard Deviation based on 601 complete cases.
Mean = -6.269 (0.056)
Standard Deviation = 9.420 (0.007)

AFFAIRS = 4.635 - 2.734*HAPPINES + Error, R² = 0.134
Standerr (1.568) (0.423)
z-values 2.955 -6.462
P-values 0.003 0.000

AFFAIRS = 14.453 - 0.939*HAPPINES - 0.446*RELIGIOU + Error, R² = 0.0957
Standerr (0.672) (0.135) (0.126)
z-values 21.508 -6.979 -3.529
P-values 0.000 0.000 0.000

AFFAIRS = 5.434 - 2.281*HAPPINES - 1.731*RELIGIOU + 0.333*YEARS+ Error, R² = 0.213
Standerr (2.072) (0.411) (0.408) (0.0885)
z-values 2.623 -5.543 -4.244 3.761
P-values 0.009 0.000 0.000 0.000

AFFAIRS = 14.406 - 0.802*HAPPINES - 0.558*RELIGIOU + 0.176*YEARS - 0.0489*AGE
Standerr (0.897) (0.137) (0.128) (0.0422) (0.0249)
z-values 16.064 -5.858 -4.352 4.182 -1.969
P-values 0.000 0.000 0.000 0.000 0.049

+ Error, R² = 0.124

AFFAIRS = 8.904 - 2.280*HAPPINES - 1.716*RELIGIOU + 0.576*YEARS - 0.190*AGE
Standerr (2.673) (0.411) (0.406) (0.138) (0.0809)
z-values 3.331 -5.551 -4.224 4.188 -2.352
P-values 0.001 0.000 0.000 0.000 0.019

+ 1.377*GENDER + Error, R² = 0.223
(0.937)
1.470
0.141

AFFAIRS =	14.537	- 0.808* <td>- 0.556*RELIGIOU</td> <td>+ 0.194*YEARS</td> <td>- 0.0547*AGE</td>	- 0.556*RELIGIOU	+ 0.194*YEARS	- 0.0547*AGE
Standerr	(0.919)	(0.137)	(0.128)	(0.0474)	(0.0259)
z-values	15.813	-5.886	-4.333	4.096	-2.113
P-values	0.000	0.000	0.000	0.000	0.035

+ 0.219*GENDER	- 0.246*CHILDREN	+ Error, R ² = 0.123
(0.303)	(0.395)	
0.722	-0.622	
0.470	0.534	

AFFAIRS =	7.712	- 2.266* <td>- 1.701*RELIGIOU</td> <td>+ 0.533*YEARS</td> <td>- 0.192*AGE</td>	- 1.701*RELIGIOU	+ 0.533*YEARS	- 0.192*AGE
Standerr	(2.903)	(0.412)	(0.408)	(0.148)	(0.0814)
z-values	2.656	-5.497	-4.173	3.612	-2.362
P-values	0.008	0.000	0.000	0.000	0.018

+ 0.966*GENDER	+ 1.027*CHILDREN	+ 0.229*OCCUPATI
(1.055)	(1.286)	(0.290)
0.915	0.798	0.789
0.360	0.425	0.430

+ Error, R² = 0.221

AFFAIRS =	14.183	- 0.812* <td>- 0.546*RELIGIOU</td> <td>+ 0.192*YEARS</td> <td>- 0.0578*AGE</td>	- 0.546*RELIGIOU	+ 0.192*YEARS	- 0.0578*AGE
Standerr	(1.313)	(0.139)	(0.129)	(0.0475)	(0.0260)
z-values	10.803	-5.856	-4.239	4.040	-2.220
P-values	0.000	0.000	0.000	0.000	0.026

+ 0.0273*GENDER	- 0.149*CHILDREN	+ 0.116*OCCUPATI
(0.346)	(0.403)	(0.102)
0.0789	-0.369	1.133
0.937	0.712	0.257

- 0.000397*EDUCATIO	+ Error, R ² = 0.122
(0.0739)	
-0.00538	
0.996	

In this analysis, the AFFAIRS variable is treated as continuous. This means that the values 0, 1, 2, 3, 7, and 12 that this variable is assumed to be numbers on an interval scale. One could also treat AFFAIRS as censored both below and above, see file **affairs3.pri**. Another way is to treat AFFAIRS as an ordinal variable with six categories (it may be better to use the exact counts and treat this as a Poisson variable, but such data is not available to me) or as a binary variable where 0 is used in one category and all values larger than 0 are used in the other category. One can then use logistic or probit regression. To do so with PRELIS, use **affairs1.lsf**, include the lines

```

Probit Regression of Affairs
sy=affairs1.lsf
co GENDER - HAPPINESS
or AFFAIRS
lr AFFAIRS on HAPPINESS RELIGIOUS YEARS AGE
ou

```

and replace cr (censored regression) by lr (logistic regression) or pr (probit regression). To use AFFAIRS as a binary variable, include the line

```
re AFFAIRS old=1-12 new=1
```

see files **affairs4.prl** – **affairs7.prl**. The results obtained from these files are not directly comparable because the variable y^* is scaled differently for different methods. However, they can be made comparable by multiplying the regression coefficients and their standard error estimates by a suitable scale factor. The probit regressions (columns 4 and 6 in Table 1) are scaled such that the error variance is 1 (standard parameterization, see Jöreskog, 2002). Using this as a standard, we must scale the other solutions by $1/\hat{\psi}$. If AFFAIRS is treated as censored below (column 2 in the table) this scale factor is $1/\sqrt{69.031} = 0.12036$, see file **affairs2.out**. If AFFAIRS is treated as censored above and below (column 3 in the table), this scale factor is $1/\sqrt{123.39} = 0.09002$, see file **affairs3.out**. For the logistic regressions (column 5 and 7 in the table) the scale factor is $\sqrt{3}/\pi = 0.55133$, because the variance of the standard logistic distribution is $\pi^2/3$. The t -values are not affected by this scaling. After this scaling the results are shown in Table 1.

Table 1: Estimated regression coefficients with different methods

Variable	Censored		Ordinal		Binary	
	Below	& Above	Probit	Logit	Probit	Logit
HAPPINESS	-0.273 (0.049)	-0.281 (0.052)	-0.284 (0.049)	-0.278 (0.047)	-0.270 (0.052)	-0.254 (0.049)
RELIGIOUS	-0.207 (0.049)	-0.208 (0.051)	-0.209 (0.049)	-0.200 (0.028)	-0.187 (0.052)	-0.181 (0.049)
YEARS	0.065 (0.016)	0.067 (0.017)	0.067 (0.016)	0.067 (0.016)	0.058 (0.017)	0.056 (0.016)
AGE	-0.019 (0.009)	-0.020 (0.010)	-0.021 (0.010)	-0.23 (0.009)	-0.020 (0.010)	-0.019 (0.010)

It is seen that all methods give similar results. For most practical purposes these results are the same. The binary methods do not make use of all information in the AFFAIRS variable. Nevertheless, the results are very close to the other methods which make use of all available information. Which of these methods should be used to estimate the model? This is not fully continuous (as if it were an amount of money spent). Neither is it fully ordinal as if the responses were classified as never, sometimes, and often. It is somewhere in between. Censored regression is a method for continuous variables and probit and logistic regressions are methods for ordinal variables.

4.3 Example 4: Reading and spelling tests

The file **readspel.lsf** contains scores on 11 reading and spelling tests for 90 school children used in a study of the meta-phonological character of the Swedish language. It is of particular interest to predict one of these tests, V23, using the other 10 variables as predictors and to determine which of these variables are the best predictors. However, a data screening of **readspel.lsf** reveals that V23 may be censored both below and above. Hence, we must use censored regression to estimate the prediction equation. The PRELIS command file is **readspel1.prl**.

```

Eleven Reading and Spelling Tests
sy=READSPEL.LSF
co all
cr V23 on V01 - V22 V24 V25
ou

```

The output shows

```

Variable V23 is censored below and above.
It has 3 ( 3.33%) values = 3.000 and 21 (23.33%) values = 16.000

```

```

Estimated Mean and Standard Deviation based on 90 complete cases.
Mean = 13.142 (0.286)
Standard Deviation = 4.397 (0.021)

```

Estimated Censored Regression based on 90 complete cases.

V23 =	- 1.724	- 0.0731*V01	+ 0.122*V02	+ 0.384*V07	- 0.129*V08
Standerr	(2.414)	(0.0815)	(0.0900)	(0.262)	(0.213)
z-values	-0.714	-0.897	1.359	1.463	-0.605
P-values	0.475	0.370	0.174	0.144	0.545

	+ 0.0954*V09	+ 0.117*V10	+ 0.208*V21	+ 0.0679*V22	+ 0.0276*V24
	(0.0688)	(0.0838)	(0.177)	(0.149)	(0.163)
	1.388	1.403	1.178	0.456	0.170
	0.165	0.161	0.239	0.648	0.865

	+ 0.208*V25	+ Error, R ² = 0.491
	(0.158)	
	1.319	
	0.187	

Residual Correlation Matrix

	V23

V23	1.000

Residual Covariance Matrix

	V23

V23	9.848

None of the predictors are statistically significant. However, in terms of the *t*-values, the most important predictors seem to be V02, V07, V09, and V10. Using only these as predictors (see file **readspel2.prl**) gives the following prediction equation.

Estimated Censored Regression based on 90 complete cases.

V23 =	1.652	+ 0.210*V02	+ 0.283*V07	+ 0.0775*V09	+ 0.169*V10
Standerr	(1.843)	(0.0663)	(0.151)	(0.0656)	(0.0801)
z-values	0.896	3.170	1.877	1.181	2.114
P-values	0.370	0.002	0.060	0.238	0.035

+ Error, R² = 0.462

Residual Covariance Matrix

	V23

V23	10.392

Here it is seen that the effects V02 and V10 are statistically significant.

Appendix: computational notes

The estimation of a censored regression equation is described in Chapter 6 of Maddala (1983) for the case of a variable that is censored below at 0. The development outlined here covers the cases when the observed variable y is censored below, censored above, censored both below and above, and not censored at all. It also covers the case when there are no regressors in which case one can estimate the mean and standard deviation of y .

Changing notation slightly from Section 2, consider the estimation of the regression equation

$$y^* = \alpha^* + \boldsymbol{\gamma}^* \mathbf{x} + z, \quad (9)$$

where α^* is the intercept term, $\boldsymbol{\gamma}^*$ is the vector of regression coefficients, and \mathbf{x} the regressors. The error term z is assumed to be normally distributed with mean 0 and variance ψ^{*2} . If there are no regressors, the second term in (9) is not included.

The observed variable

$$\begin{aligned} y &= c_1 \text{ if } y^* \leq c_1 \\ &= y^* \text{ if } c_1 \leq y^* \leq c_2 \\ &= c_2 \text{ if } y^* \geq c_2, \end{aligned}$$

where c_1 and c_2 are constants. If y is censored below set $c_2 = +\infty$. If y is censored above set $c_1 = -\infty$. If y is not censored set both $c_1 = -\infty$ and $c_2 = +\infty$.

Let (y_i, \mathbf{x}_i) be the observed values of y and \mathbf{x} of case i in a random sample of N independent observations. The likelihood of (y_i, \mathbf{x}_i) is

$$L_i = \left[\Phi \left(\frac{c_1 - \alpha^* - \boldsymbol{\gamma}^* \mathbf{x}_i}{\psi^*} \right) \right]^{j_{1i}} \left[\frac{1}{\sqrt{2\pi\psi^*}} e^{-\frac{1}{2} \left(\frac{y_i - \alpha^* - \boldsymbol{\gamma}^* \mathbf{x}_i}{\psi^*} \right)^2} \right]^{1-j_{1i}-j_{2i}} \left[1 - \Phi \left(\frac{c_2 - \alpha^* - \boldsymbol{\gamma}^* \mathbf{x}_i}{\psi^*} \right) \right]^{j_{2i}},$$

where $j_{1i} = 1$ if $y = c_1$ and $j_{1i} = 0$ otherwise and $j_{2i} = 1$ if $y = c_2$ and $j_{2i} = 0$ otherwise. Note that j_{1i} and j_{2i} cannot be 1 simultaneously.

The log likelihood is

$$\ln L = \sum_{i=1}^N \ln L_i.$$

This is to be maximized with respect to the parameter vector $\boldsymbol{\theta}^* = (\alpha^*, \boldsymbol{\gamma}^*, \psi^*)$.

First and second derivatives of $\ln L$ with respect to $\boldsymbol{\theta}^*$ are very complicated. They will be considerably simplified and the maximization of $\ln L$ will be considerably more efficient if another parameterization due to Tobin (1858) is used.

This parameterization uses the parameter vector $\boldsymbol{\theta}' = (\alpha, \boldsymbol{\gamma}', \psi)$ instead of $\boldsymbol{\theta}^*$, where $\alpha = \alpha^* / \psi^*$, $\boldsymbol{\gamma}' = \boldsymbol{\gamma}^* / \psi^*$, and $\psi = 1 / \psi^*$.

Multiplication of (9) by $\psi = 1/\psi^*$ gives

$$\psi y^* = \alpha + \gamma' \mathbf{x} + \nu, \quad (10)$$

where $\nu = \psi z = z/\psi^*$ which is $N(0,1)$. Then

$$\begin{aligned} y = c_1 &\leftrightarrow y^* \leq c_1 \leftrightarrow \psi y^* \leq \psi c_1 \leftrightarrow \nu \leq \psi c_1 - \alpha - \gamma' \mathbf{x}, \\ y = c_2 &\leftrightarrow y^* \geq c_2 \leftrightarrow \psi y^* \geq \psi c_2 \leftrightarrow \nu \geq \psi c_2 - \alpha - \gamma' \mathbf{x}. \end{aligned}$$

Hence the likelihood L_i becomes

$$L_i = \left[\Phi(\psi c_1 - \alpha - \gamma' \mathbf{x}_i) \right]^{j_{1i}} \left[\frac{1}{\sqrt{2\pi\psi^*}} e^{-\frac{1}{2}(\psi y_i - \alpha - \gamma' \mathbf{x}_i)^2} \right]^{1-j_{1i}-j_{2i}} \left[1 - \Phi(\psi c_2 - \alpha - \gamma' \mathbf{x}_i) \right]^{j_{2i}}.$$

Let

$$\delta_i = \psi y_i - \alpha - \gamma' \mathbf{x}_i. \quad (11)$$

Then $\ln L_i$ becomes

$$\ln L_i = -\ln \sqrt{2\pi} + (1 - j_{1i} - j_{2i}) \left(\ln \psi - \frac{1}{2} \delta_i^2 \right) + j_{1i} \ln \Phi(\delta_i) + j_{2i} \ln [1 - \Phi(\delta_i)]. \quad (12)$$

First and second derivatives of $\ln L_i$ are straightforward by noting that $\partial \delta_i / \partial \alpha = -1$, $\partial \delta_i / \partial \gamma = -\mathbf{x}_i$ and $\partial \delta_i / \partial \psi = y_i$.

Furthermore, $\Phi'(t) = \phi(t)$, $\phi'(t) = -t\phi(t)$ and if $A(t) = \phi(t) / \Phi(t)$, then $A'(t) = -A(t)[t + A(t)] = B(t)$, say.

Omitting index i , the required derivatives are

$$\begin{aligned} \partial \ln L / \partial \alpha &= (1 - j_1 - j_2) \delta - j_1 A(\delta) + j_2 A(-\delta) \\ \partial \ln L / \partial \gamma &= (1 - j_1 - j_2) \delta \mathbf{x} - j_1 A(\delta) \mathbf{x} + j_2 A(-\delta) \mathbf{x} \\ \partial \ln L / \partial \psi &= (1 - j_1 - j_2)(1/\psi - \delta y) + j_1 A(\delta) y - j_2 A(-\delta) y \\ \partial^2 \ln L / \partial \alpha \partial \alpha &= -(1 - j_1 - j_2) + j_1 B(\delta) + j_2 B(-\delta) \\ \partial^2 \ln L / \partial \gamma \partial \alpha &= -(1 - j_1 - j_2) \mathbf{x} + j_1 B(\delta) \mathbf{x} + j_2 B(-\delta) \mathbf{x} \\ \partial^2 \ln L / \partial \gamma \partial \gamma &= -(1 - j_1 - j_2) \mathbf{x} \mathbf{x}' + j_1 B(\delta) \mathbf{x} \mathbf{x}' + j_2 B(-\delta) \mathbf{x} \mathbf{x}' \\ \partial^2 \ln L / \partial \psi \partial \alpha &= -(1 - j_1 - j_2) y + j_1 B(\delta) y + j_2 B(-\delta) y \\ \partial^2 \ln L / \partial \psi \partial \gamma &= -(1 - j_1 - j_2) y \mathbf{x}' + j_1 B(\delta) y \mathbf{x}' + j_2 B(-\delta) y \mathbf{x}' \\ \partial^2 \ln L / \partial \psi \partial \psi &= -(1 - j_1 - j_2)(1/\psi^2 + y^2) + j_1 B(\delta) y^2 + j_2 B(-\delta) y^2. \end{aligned}$$

Maximizing $\ln L$ is equivalent to minimizing the fit function $F(\boldsymbol{\theta}) = -\ln L$. Let $g(\boldsymbol{\theta}) = \partial F / \partial \boldsymbol{\theta}$ be the gradient vector and $H(\boldsymbol{\theta}) = \partial^2 F / \partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'$ be the Hessian matrix. Amemiya (1973) proved that \mathbf{H} is positive definite everywhere,

The fit function $F(\boldsymbol{\theta})$ is minimized using a Newton-Raphson procedure which converges very fast. The starting values $\boldsymbol{\theta}_0$ are the parameters estimated by OLS. Successive estimates are obtained by the formula

$$\boldsymbol{\theta}_{s+1} = \boldsymbol{\theta}_s = \mathbf{H}_s^{-1} \mathbf{g}_s,$$

where $\mathbf{g}_s = \mathbf{g}(\boldsymbol{\theta}_s)$ and $\mathbf{H}_s = \mathbf{H}(\boldsymbol{\theta}_s)$.

Let $\hat{\boldsymbol{\theta}} = \left(\hat{\alpha}, \hat{\boldsymbol{\gamma}}, \hat{\psi} \right)$ be the maximum likelihood estimates of $\boldsymbol{\theta}$. The asymptotic covariance matrix of $\hat{\boldsymbol{\theta}}$ is $\mathbf{E} = \mathbf{H}^{-1}(\boldsymbol{\theta})$ evaluated at the true parameter $\boldsymbol{\theta}$. Since the transformation from $\boldsymbol{\theta}$ to $\boldsymbol{\theta}^*$ is one-to-one, the maximum likelihood estimates of $\boldsymbol{\theta}^*$ is where $\hat{\alpha}^* = \hat{\alpha} / \hat{\psi}$, $\hat{\boldsymbol{\gamma}}^* = \hat{\boldsymbol{\gamma}} / \hat{\psi}$, and $\hat{\psi}^* = 1 / \hat{\psi}$.

To obtain the asymptotic covariance matrix of $\hat{\boldsymbol{\theta}}^*$, we evaluate the matrix $\partial \boldsymbol{\theta}^* / \partial \boldsymbol{\theta}'$. This is

$$\partial \boldsymbol{\theta}^* / \partial \boldsymbol{\theta}' = (1 / \psi^2) \begin{vmatrix} \psi & \mathbf{0}' & -\alpha \\ \mathbf{0} & \psi \mathbf{1} & \boldsymbol{\gamma} \\ 0 & \mathbf{0}' & -1 \end{vmatrix} = \mathbf{A}(\boldsymbol{\theta}), \text{ say,}$$

where $\mathbf{0}$ and $\mathbf{1}$ are column vectors of zeros and ones, respectively. The asymptotic covariance matrix of $\hat{\boldsymbol{\theta}}^* = \mathbf{A} \mathbf{E} \mathbf{A}'$, where \mathbf{A} and \mathbf{E} are evaluated at the true parameter values. An estimate of the asymptotic of $\hat{\boldsymbol{\theta}}^* = \mathbf{A} \mathbf{E} \mathbf{A}'$ is obtained by using the estimated parameter values in \mathbf{A} and \mathbf{E} . Asymptotic standard error estimates of the parameter estimates are obtained as the square roots of the diagonal elements of this matrix.

5. References

- Amemiya, T. (1973). Regression analysis when the dependent variable is truncated normal. *Econometrica*, **41**, 997-1016.
- Du Toit, M. & Du Toit, S. (2001). *Interactive LISREL: User's Guide*. Chicago: Scientific Software International.
- Fair, R. (1978). A theory of extramarital affairs. *Journal of Political Economy*, **85**, 45-61.
- Goldberger, A.S. (1964). *Econometric theory*. New York: Wiley.
- Greene, W.H. (2000). *Econometric Analysis*. Fourth Edition. London: Prentice Hall International.
- Joreskog, K.G. (2002). Structural equation modeling with ordinal variables using LISREL. Available at https://ssicentral.com/wp-content/uploads/2021/04/lis_ordinal.pdf
- Maddala, G.S. (1983). *Limited Dependent and Qualitative Variables in Econometrics*. Cambridge University Press.
- Melenberg, B. and van Soest, A. (1996). Parametric and semi-parametric modeling of vacation expenditures. *Journal of Applied Econometrics*, **11**, 59-76.
- Quester, A. & Greene, W. (1982). Divorce risk and wives' labor supply behavioral. *Social Science Quarterly*, 16-27.
- Tobin, J. (1958). Estimation of relationships for limited dependent variables. *Econometrica*, **26**, 24-36.

Witte, A. (1980). Estimating an economic model of crime with individual data. *Quarterly Journal of Economic*, **94**, 57.84.