

CHAPTER 4: CENSORING, TRUNCATION, AND SELECTION

Joan Llull

Microeconometrics
IDEA PhD Program

INTRODUCTION

Introduction

In this chapter we review models that deal with **censored**, **truncated** and **self-selected** data.

We consider a latent variable that is described by a **linear** model:

$$y^* = \mathbf{x}'\boldsymbol{\beta} + \varepsilon.$$

Although this assumption can be trivially relaxed, we **assume** (unless when otherwise noted) that:

$$\varepsilon|\mathbf{x} \sim \mathcal{N}(0, \sigma^2).$$

Truncation

We say that our sample is **(left) truncated** whenever:

$$y = \begin{cases} y^* & \text{if } y^* > 0 \\ - & \text{if } y^* \leq 0, \end{cases}$$

The **problem** with these data is that:

$$\mathbb{E}[y|\mathbf{x}] = \mathbb{E}[y^*|\mathbf{x}, y^* \text{ is observed}] = \mathbf{x}'\boldsymbol{\beta} + \mathbb{E}[\varepsilon|\mathbf{x}, \varepsilon > -\mathbf{x}'\boldsymbol{\beta}].$$

We analyze **left truncation**. Right truncation is analogous.

Censoring

We say that our sample is **(left) censored** whenever:

$$y = \begin{cases} y^* & \text{if } y^* > 0 \\ 0 & \text{if } y^* \leq 0. \end{cases}$$

The **problem** with these data is that:

$$\mathbb{E}[y|\mathbf{x}] = \mathbb{E}[\mathbb{E}(y|\mathbf{x}, d)|\mathbf{x}] = \Pr[\varepsilon > -\mathbf{x}'\boldsymbol{\beta}|\mathbf{x}](\mathbf{x}'\boldsymbol{\beta} + \mathbb{E}[\varepsilon|\mathbf{x}, \varepsilon > -\mathbf{x}'\boldsymbol{\beta}]),$$

where $d = \mathbb{1}\{y^* > 0\}$.

Selection

We say that our sample is **self-selected** whenever:

$$y = \begin{cases} y^* & \text{if } \mathbf{z}'\boldsymbol{\gamma} + \nu > 0 \\ - & \text{otherwise,} \end{cases}$$

and $d \equiv \mathbb{1}\{\mathbf{z}'\boldsymbol{\gamma} + \nu > 0\}$ is observed.

The **problem** with these data is that:

$$\mathbb{E}[y|\mathbf{x}] = \mathbb{E}[y^*|\mathbf{x}, d = 1] = \mathbf{x}'\boldsymbol{\beta} + \mathbb{E}[\varepsilon|\mathbf{z}'\boldsymbol{\gamma} + \nu > 0].$$

Endogenous selection whenever $\mathbb{E}[\varepsilon|\mathbf{z}'\boldsymbol{\gamma} + \nu > 0] \neq 0$.

CENSORING AND TRUNCATION. THE TOBIT MODEL

Maximum Likelihood Estimation (truncation)

The individual likelihood of an observation from a **(left) truncated** sample is:

$$g(y|\mathbf{x}, y > 0) = \frac{f(y|\mathbf{x})}{\Pr[y > 0|\mathbf{x}]} = \frac{f(y|\mathbf{x})}{1 - F(0|\mathbf{x})}.$$

Hence, the **log-likelihood** function is:

$$\mathcal{L}_N(\theta) = \sum_{i=1}^N \{\ln f(y_i|\mathbf{x}_i) - \ln(1 - F(0|\mathbf{x}_i))\}.$$

In the **Tobit** model (given by **normality** assumption):

$$\mathcal{L}_N(\theta) = \sum_{i=1}^N \left\{ -\frac{1}{2} \ln 2\pi - \frac{1}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} (y_i - \mathbf{x}_i' \boldsymbol{\beta})^2 - \ln \Phi \left(\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma} \right) \right\}.$$

Maximum Likelihood Estimation (censoring)

The individual likelihood of an observation from a **(left) censored** sample is:

$$g(y|\mathbf{x}, y > 0) = f(y|\mathbf{x})^d F(0|\mathbf{x})^{1-d} = \begin{cases} f(y|\mathbf{x}) & \text{if } y^* > 0 \\ F(0|\mathbf{x}) & \text{if } y^* \leq 0. \end{cases}$$

Hence, the **log-likelihood** function is:

$$\mathcal{L}_N(\theta) = \sum_{i=1}^N \{d_i \ln f(y_i|\mathbf{x}_i) + (1 - d_i) \ln(F(0|\mathbf{x}_i))\}.$$

In the **Tobit** model (given by **normality** assumption):

$$\mathcal{L}_N(\theta) = \sum_{i=1}^N \left\{ d_i \left(-\frac{1}{2} \ln 2\pi - \frac{1}{2} \ln \sigma^2 - \frac{(y_i - \mathbf{x}'_i \boldsymbol{\beta})^2}{2\sigma^2} \right) + (1 - d_i) \ln \left(1 - \Phi \left(\frac{\mathbf{x}'_i \boldsymbol{\beta}}{\sigma} \right) \right) \right\}.$$

Potential inconsistency of MLE

As usual, consistency requires that the likelihood is **correctly specified**.

Potentially more severe here. Intuition can be seen in the **FOC** for β in the **truncated data**:

$$\frac{\partial \mathcal{L}_N}{\partial \beta} = \sum_{i=1}^N \frac{1}{\sigma^2} \left(y_i - \mathbf{x}_i' \beta - \sigma \lambda \left(\frac{\mathbf{x}_i' \beta}{\sigma} \right) \right) \mathbf{x}_i = \mathbf{0},$$

where $\lambda(z) = \phi(z)/\Phi(z) = \mathbb{E}[\epsilon | \epsilon > -z]$ if $\epsilon \sim \mathcal{N}(0, 1)$.

Heteroscedastic or **non-normal** errors lead to biased estimates.

Alternatives (I): Heckman two-step procedure

This is usually used for **self-selected** samples —however, censoring is a **special case**.

Relies on:

$$\mathbb{E}[y|\mathbf{x}, y > 0] = \mathbf{x}'\boldsymbol{\beta} + \sigma\lambda\left(\frac{\mathbf{x}'\boldsymbol{\beta}}{\sigma}\right).$$

Two steps:

- Estimate $\boldsymbol{\alpha} = \boldsymbol{\beta}/\sigma$ from a **Probit** over $d = \mathbb{1}\{y^* > 0\}$.
- Use $\hat{\boldsymbol{\alpha}}$ to compute $\lambda(\mathbf{x}'\hat{\boldsymbol{\alpha}})$, and include it as a control in a **linear regression** estimated with **uncensored observations** to identify $\boldsymbol{\beta}$ and σ .

Alternatives (II): Median Regression

If censoring is **below the median** of y we can still make inference on the median.

If the distribution is **symmetric** (e.g. normal distribution), the mean and the median coincide.

The Censored Least Absolute Deviations estimator is:

$$\hat{\beta}_{CLAD} = \arg \min_{\beta} N^{-1} \sum_{i=1}^N |y_i - \max(\mathbf{x}_i' \beta, 0)|.$$

Alternatives (III): Symmetrically Trimmed Mean

Key assumption: the distribution of $\varepsilon|\mathbf{x}$ is **symmetric** around 0.

We **do not consider** observations with $\mathbf{x}'\boldsymbol{\beta} < 0$.

Probability that $\mathbf{x}'\boldsymbol{\beta} + \varepsilon < 0$ (and hence the **observation is censored**) is equal to that of $\mathbf{x}'\boldsymbol{\beta} + \varepsilon > 2\mathbf{x}'\boldsymbol{\beta}$. Hence:

$$\mathbb{E}[\mathbb{1}\{\mathbf{x}'\boldsymbol{\beta} > 0\}(\min(y, 2\mathbf{x}'\boldsymbol{\beta}) - \mathbf{x}'\boldsymbol{\beta})\mathbf{x}] = \mathbf{0}.$$

This estimator delivers FOC that are **sample analogs** to the previous moment conditions:

$$\hat{\boldsymbol{\beta}}_{STM} = \arg \min_{\boldsymbol{\beta}} \frac{1}{N} \sum_{i=1}^N \left\{ \left[y_i - \max\left(\frac{y_i}{2}, \mathbf{x}'_i \boldsymbol{\beta}\right) \right]^2 + \mathbb{1}\{y_i > 2\mathbf{x}'_i \boldsymbol{\beta}\} \left[\frac{y_i^2}{4} - \max(0, \mathbf{x}'_i \boldsymbol{\beta}) \right]^2 \right\}.$$

SELECTION

The Sample Selection Model

The **model** is defined by:

$$\begin{aligned} y^* &= \mathbf{x}'\boldsymbol{\beta} + \varepsilon \\ d &= \mathbb{1}\{\mathbf{z}'\boldsymbol{\gamma} + \nu\} \end{aligned} \quad \left(\begin{array}{c} \varepsilon \\ \nu \end{array} \right) \Big| \mathbf{z}, \mathbf{x} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma^2 & \rho\sigma \\ \rho\sigma & 1 \end{bmatrix} \right).$$

We only **observe** $y = y^* \times d$.

Without loss of generality, \mathbf{x} is **included** in \mathbf{z} .

The **likelihood** of our sample is:

$$L_N(\theta) = \prod_{i=1}^N (1 - \Phi(\mathbf{z}'\boldsymbol{\gamma}))^{1-d} \{f(y^*|\mathbf{z}) \Pr(d=1|y^*, \mathbf{z})\}^d,$$

where:

$$f(y^*|\mathbf{z}) = \frac{1}{\sigma} \phi \left(\frac{y^* - \mathbf{x}'\boldsymbol{\beta}}{\sigma} \right) \quad \text{and} \quad \Pr(d=1|y^*, \mathbf{z}) = \Phi \left(\frac{\mathbf{z}'\boldsymbol{\gamma} + \frac{\rho}{\sigma}(y^* - \mathbf{x}'\boldsymbol{\beta})}{\sqrt{1 - \rho^2}} \right).$$

(Note that when evaluating this, we observe y^* as we only have to evaluate it when $d=1$).

Heckman two-step procedure (Heckit)

Same idea as we have seen for **Tobit**.

Relies on:

$$\mathbb{E}[y|\mathbf{x}, d = 1] = \mathbf{x}'\boldsymbol{\beta} + \rho\sigma\lambda(\mathbf{z}'\boldsymbol{\gamma}).$$

Two steps:

- Estimate $\boldsymbol{\gamma}$ from a **Probit** over $d = \mathbb{1}\{\mathbf{z}'\boldsymbol{\gamma} + \nu > 0\}$.
- Include $\lambda(\mathbf{z}'\hat{\boldsymbol{\gamma}})$ as a control in a **linear regression** with the observations for which the **outcome is observed** to estimate $\boldsymbol{\beta}$ and $\rho\sigma$.

Heckman two-step procedure: Remarks

Important remarks to make:

- **Standard errors** have to control for the particular form of heteroskedasticity and for using $\hat{\gamma}$: **bootstrap**.
- **Credible identification** requires excluded variables in \mathbf{z} :

$$\mathbb{E}[y|\mathbf{z}] \approx \mathbf{x}'\boldsymbol{\beta} + a + b(\mathbf{z}'\hat{\gamma}).$$

- It is a **LIML** estimation using, in the previous likelihood:

$$f(y^*, d = 1|\mathbf{z}) = f(y^*|d = 1, \mathbf{z}) \Pr(d = 1|\mathbf{z}).$$

- **Test** for selection: $\rho \neq 0$.
- Alternative: **semi-parametric** estimation of $g(\mathbf{z}'\boldsymbol{\gamma})$ vs $\lambda(\mathbf{z}'\boldsymbol{\gamma})$.