

# Censored and Truncated Outcome Panel Data Models

230347 Advanced Microeconometrics  
Tilburg University

Christoph Walsh

## Censored Data

- ▶  $y_{it}$  is censored when it is partly continuous but has positive probability mass at one or more points.
  - ▶ For example,  $y_{it}$  is continuous when  $y_{it} > 0$  but has a large mass at  $y_{it} = 0$ .
- ▶ We can sometimes think of the underlying model as:

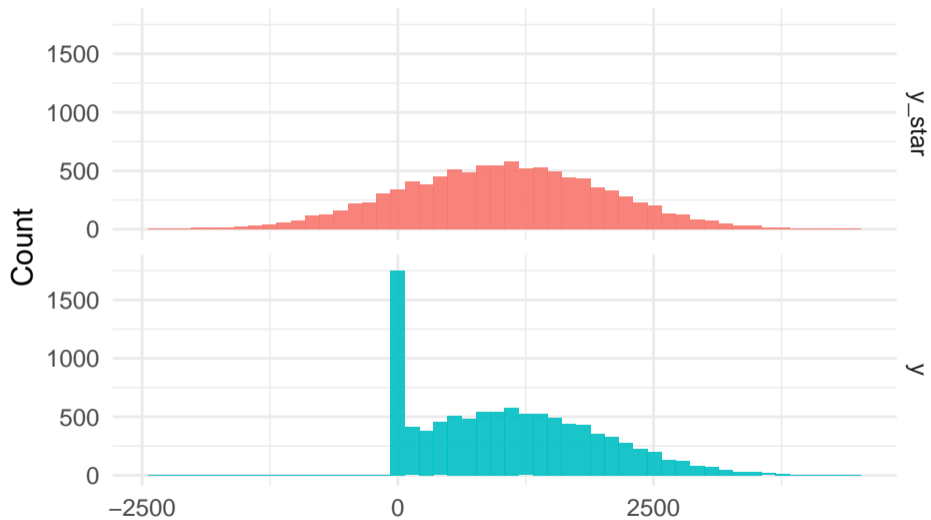
$$y_{it}^* = \alpha_i + \mathbf{x}'_{it} \boldsymbol{\beta} + \varepsilon_{it}$$

but we observe:

$$y_{it} = \begin{cases} y_{it}^* & \text{if } y_{it}^* > \underline{y} \\ \underline{y} & \text{if } y_{it}^* \leq \underline{y} \end{cases} \quad \text{or} \quad y_{it} = \begin{cases} \bar{y} & \text{if } y_{it}^* \geq \bar{y} \\ y_{it}^* & \text{if } y_{it}^* < \bar{y} \end{cases}$$

- ▶ For example, top-coded income.
- ▶ Other times we can think of  $\bar{y}$  or  $\underline{y}$  as a corner solution of an optimization problem.
  - ▶ For example, hours worked, firm expenditure on R&D.

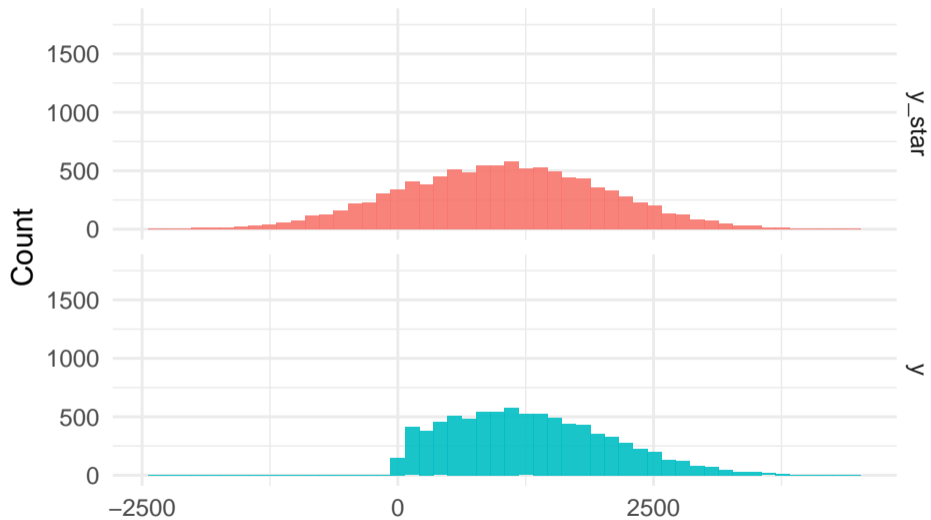
## Histogram of censored $y_{it}$ left of zero



# Truncated Data

- ▶ Our sample may be **truncated**, where our sample only has observations where  $y_{it} > \underline{y}$  or  $y_{it} < \bar{y}$
- ▶ For example, we may only observe people who work.

## Histogram of truncated $y_{it}$



## Models in this Topic

- ▶ Static Censored Random Effects
- ▶ Static Truncated Fixed Effects
- ▶ It is possible to estimate Static & Dynamic Censored Fixed Effects models, but we won't cover them here.

## Censored Data: Panel Random Effects Tobit Model

- ▶ We consider the left-censored data case where  $\underline{y} = 0$ .
- ▶ We observe:

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

- ▶ Let  $d_{it} = \mathbb{1}\{y_{it} > 0\}$ .
- ▶ If  $\varepsilon_{it} \sim \mathcal{N}(0, \sigma_\varepsilon^2)$ , then using  $\varepsilon_{it} = y_{it}^* - \alpha_i - \mathbf{x}'_{it}\boldsymbol{\beta}$ , the joint conditional density of  $\mathbf{y}_i = (y_{i1}, \dots, y_{iT})$  is

$$f(\mathbf{y}_i | \mathbf{X}_i, \alpha_i, \boldsymbol{\beta}, \sigma_\varepsilon^2) = \prod_{t=1}^T \left[ \frac{1}{\sigma_\varepsilon} \phi\left(\frac{y_{it} - \alpha_i - \mathbf{x}'_{it}\boldsymbol{\beta}}{\sigma_\varepsilon}\right) \right]^{d_{it}} \left[ 1 - \Phi\left(\frac{\alpha_i + \mathbf{x}'_{it}\boldsymbol{\beta}}{\sigma_\varepsilon}\right) \right]^{1-d_{it}}$$

where  $\mathbf{X}_i = (\mathbf{x}_{i1}, \dots, \mathbf{x}_{iT})$  and  $\phi$  and  $\Phi$  are the pdf and cdf of the standard normal distribution respectively.

## Censored Data: Panel Random Effects Tobit Model

- ▶ If we model  $\alpha_i \sim \mathcal{N}(0, \sigma_\alpha^2)$ , then we can integrate out the  $\alpha_i$ :

$$f(\mathbf{y}_i | \mathbf{X}_i, \boldsymbol{\beta}, \sigma_\varepsilon^2, \sigma_\alpha^2) = \int_{-\infty}^{\infty} f(\mathbf{y}_i | \mathbf{X}_i, \alpha_i, \boldsymbol{\beta}, \sigma_\varepsilon^2) \frac{1}{\sqrt{2\pi\sigma_\alpha^2}} \exp\left(\frac{-\alpha_i^2}{2\sigma_\alpha^2}\right) d\alpha_i$$

- ▶ There is no closed-form solution for the likelihood and therefore needs to be computed using simulation methods.
- ▶ We can perform the same change of variables as with the probit random effects and approximate the integral with Gauss-Hermite quadrature.



## Truncated Fixed Effects: Only observe $y_{it}^*$ when $y_{it}^* > 0$

- ▶ When data are truncated, we cannot eliminate the fixed effects by differencing or mean differencing.
- ▶ For observed  $y_{it}$ :

$$\begin{aligned}y_{it} &= \mathbb{E} [y_{it}^* | \mathbf{x}_{it}, \alpha_i, y_{it}^* > 0] + \nu_{it} \\ &= \alpha_i + \mathbf{x}'_{it} \boldsymbol{\beta} + \mathbb{E} [\varepsilon_{it} | \varepsilon_{it} > -\alpha_i - \mathbf{x}'_{it} \boldsymbol{\beta}] + \nu_{it}\end{aligned}$$

- ▶ Consider the  $T = 2$  case. Taking differences:

$$\begin{aligned}y_{i2} - y_{i1} &= (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta} + \mathbb{E} [\varepsilon_{i2} | \varepsilon_{i2} > -\alpha_i - \mathbf{x}'_{i2} \boldsymbol{\beta}] - \\ &\quad \mathbb{E} [\varepsilon_{i1} | \varepsilon_{i1} > -\alpha_i - \mathbf{x}'_{i1} \boldsymbol{\beta}] + \nu_{i2} - \nu_{i1}\end{aligned}$$

- ▶ In general, this still depends on  $\alpha_i$  (unless  $\mathbf{x}_{i1} = \mathbf{x}_{i2}$ )

## Honoré (1992)

- ▶ Suppose we restricted our analysis to observations satisfying

$$y_{i1} \geq -(\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta} \quad \text{and} \quad y_{i2} \geq (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta}$$

- ▶ Suppose that  $(\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta} > 0$  ( $\exists$  similar argument for the opposite case). Then:

$$\begin{aligned} & \mathbb{E} [y_{i2} | \mathbf{x}_{i2}, \alpha_i, y_{i2} \geq (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta}] \\ &= \alpha_i + \mathbf{x}'_{i2} \boldsymbol{\beta} + \mathbb{E} [\varepsilon_{i2} | \varepsilon_{i2} \geq -\alpha_i - \mathbf{x}'_{i2} \boldsymbol{\beta} + (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta}] \\ &= \alpha_i + \mathbf{x}'_{i2} \boldsymbol{\beta} + \mathbb{E} [\varepsilon_{i2} | \varepsilon_{i2} \geq -\alpha_i - \mathbf{x}'_{i1} \boldsymbol{\beta}] \end{aligned}$$

- ▶ Since  $(\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta} > 0$ , the restriction doesn't bind for  $y_{i1}$ :

$$\begin{aligned} \mathbb{E} [y_{i1} | \mathbf{x}_{i1}, \alpha_i, y_{i1} \geq -(\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta}] &= \mathbb{E} [y_{i1} | \mathbf{x}_{i1}, \alpha_i, y_{i1} \geq 0] \\ &= \alpha_i + \mathbf{x}'_{i1} \boldsymbol{\beta} + \mathbb{E} [\varepsilon_{i1} | \varepsilon_{i1} \geq -\alpha_i - \mathbf{x}'_{i1} \boldsymbol{\beta}] \end{aligned}$$

## Honoré (1992)

- ▶ If we assume the  $\varepsilon_{it} | \mathbf{x}_{it}, \alpha_i$  are iid, then:

$$\mathbb{E} [\varepsilon_{i1} | \varepsilon_{i1} \geq -\alpha_i - \mathbf{x}'_{i1} \boldsymbol{\beta}] = \mathbb{E} [\varepsilon_{i2} | \varepsilon_{i2} \geq -\alpha_i - \mathbf{x}'_{i1} \boldsymbol{\beta}]$$

- ▶ Therefore

$$\mathbb{E} [y_{i1} | \mathbf{x}_{i1}, \alpha_i, y_{i1} \geq -(\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta}] = \alpha_i + \mathbf{x}'_{i1} \boldsymbol{\beta} + \mathbb{E} [\varepsilon_{i1} | \varepsilon_{i1} \geq -\alpha_i - \mathbf{x}'_{i1} \boldsymbol{\beta}]$$

$$\mathbb{E} [y_{i2} | \mathbf{x}_{i2}, \alpha_i, y_{i2} \geq (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta}] = \alpha_i + \mathbf{x}'_{i2} \boldsymbol{\beta} + \mathbb{E} [\varepsilon_{i1} | \varepsilon_{i1} \geq -\alpha_i - \mathbf{x}'_{i1} \boldsymbol{\beta}]$$

- ▶ Together:

$$\begin{aligned} \mathbb{E} [y_{i2} - y_{i1} | \mathbf{x}_{i1}, \mathbf{x}_{i2}, \alpha_i, y_{i1} \geq -(\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta}, y_{i2} \geq (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta}] \\ = (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \boldsymbol{\beta} \end{aligned}$$

which no longer depends on the fixed effect  $\alpha_i$ .

- ▶ This only requires the iid assumption. We don't assume anything about the distribution of  $\varepsilon_{it}$ .

## Honoré (1992): Estimation when $T = 2$

- ▶ If we knew the true  $\beta$ , we could estimate it with OLS in the model:

$$y_{i2} - y_{i1} = (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \beta + \nu_{i2} - \nu_{i1}$$

using the sample where:

- ▶  $y_{i1} \geq -(\mathbf{x}_{i2} - \mathbf{x}_{i1})' \beta$
- ▶  $y_{i2} \geq (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \beta$
- ▶ However, we do not know  $\beta$ .

## Honoré (1992): Estimation

- ▶ Honoré (1992) proposes the following objective:

$$\begin{aligned}\hat{\beta} = \arg \min_{\beta} & \sum_{i=1}^N \left\{ [y_{i2} - y_{i1} - (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \beta]^2 \right. \\ & \times \mathbb{1} \{ y_{i1} \geq -(\mathbf{x}_{i2} - \mathbf{x}_{i1})' \beta, y_{i2} \geq (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \beta \} \\ & + y_{i1}^2 \mathbb{1} \{ y_{i1} \geq -(\mathbf{x}_{i2} - \mathbf{x}_{i1})' \beta, y_{i2} < (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \beta \} \\ & \left. + y_{i2}^2 \mathbb{1} \{ y_{i1} < -(\mathbf{x}_{i2} - \mathbf{x}_{i1})' \beta, y_{i2} \geq (\mathbf{x}_{i2} - \mathbf{x}_{i1})' \beta \} \right\}\end{aligned}$$

## Honoré (1992): Estimation

- ▶ Why the 2nd and 3rd term?
- ▶ Consider the single regressor case.
- ▶ Suppose we estimated  $\beta$  by minimizing:

$$\sum_{i=1}^N [y_{i2} - y_{i1} - (x_{i2} - x_{i1})\beta]^2 \mathbb{1}\{y_{i1} \geq -(x_{i2} - x_{i1})\beta, y_{i2} \geq (x_{i2} - x_{i1})\beta\}$$

- ▶ By setting  $\beta$  sufficiently large or small, no  $y_{i1}$  and  $y_{i2}$  will satisfy  $y_{i1} \geq -(x_{i2} - x_{i1})\beta$  and  $y_{i2} \geq (x_{i2} - x_{i1})\beta$  simultaneously for any  $i$ .
- ▶ The objective function would then be zero, its lowest possible value.
- ▶ The inclusion of the 2nd and 3rd term excludes these trivial solutions.

# Reading and References

- ▶ Cameron and Trivedi 23.5 for Random effects Tobit.
- ▶ Hsiao 8.4 and Honoré (1992) for Truncated Fixed Effects.

## References:

HONORÉ, B. E. (1992): "Trimmed LAD and least squares estimation of truncated and censored regression models with fixed effects," *Econometrica*, 533–565.