

Profesor: Héctor Bahamonde, PhD.

e: hector.bahamonde@uoh.cl

w: www.hectorbahamonde.com

Curso: MLE.

TA: Gonzalo Barria.

I. OUTCOMES CENSURADOS/TRUNCADOS

En esta clase pensaremos en las variables dependientes truncadas o censuradas y_* . *Censored* es diferente a *truncated*. Como veremos en esta clase,

1. **Censuradas:** Las variables dependientes censuradas y_i^* son aquellas donde a pesar de observar todo el dataset, sólo tenemos información parcial de y_i .
2. **Truncadas:** Las variables dependientes truncadas y_i^* son aquellas donde (por algún motivo) no vemos todo el dataset, y en consecuencia sólo tenemos información parcial de y_i .

También es posible ver las diferencias conceptuales entre variables truncadas y censuradas de manera gráfica:

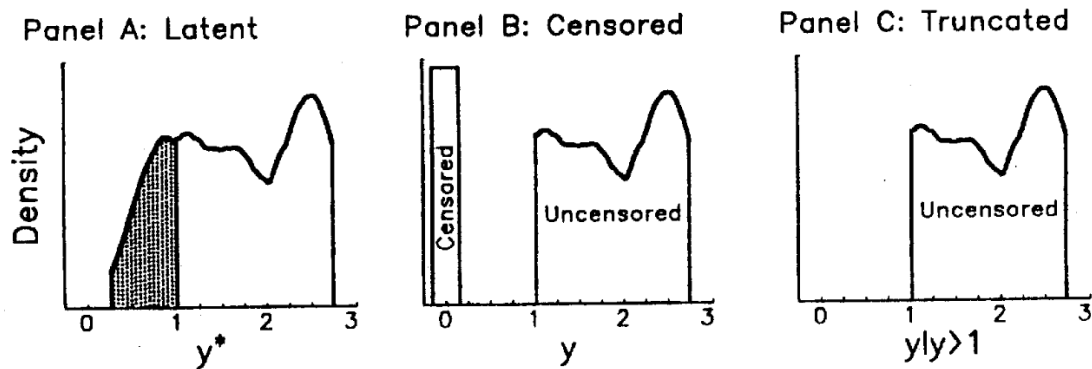


Figure 7.1. Latent, Censored, and Truncated Variables

Afortunadamente tenemos los modelos *tobit* para estimarlos. Como veremos, ambos procesos (*censored*/*truncated*) asumen supuestos distribucionales diferentes.

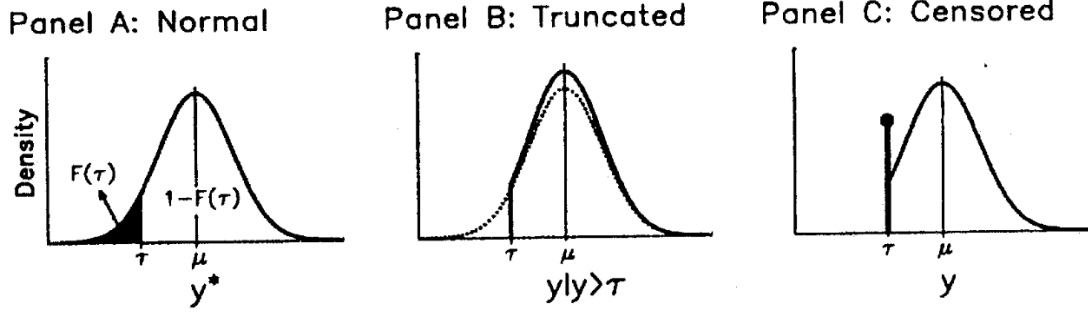


Figure 7.3. Normal Distribution With Truncation and Censoring

Como ves en la Figura 7.3, las distribuciones de y_i^* varían según el tipo de ignorancia con el que estamos lidiando. **Nota que uno de los supuestos de la especificación del modelo tobit es que y_i^* está normalmente distribuido.**

I. Censored Data

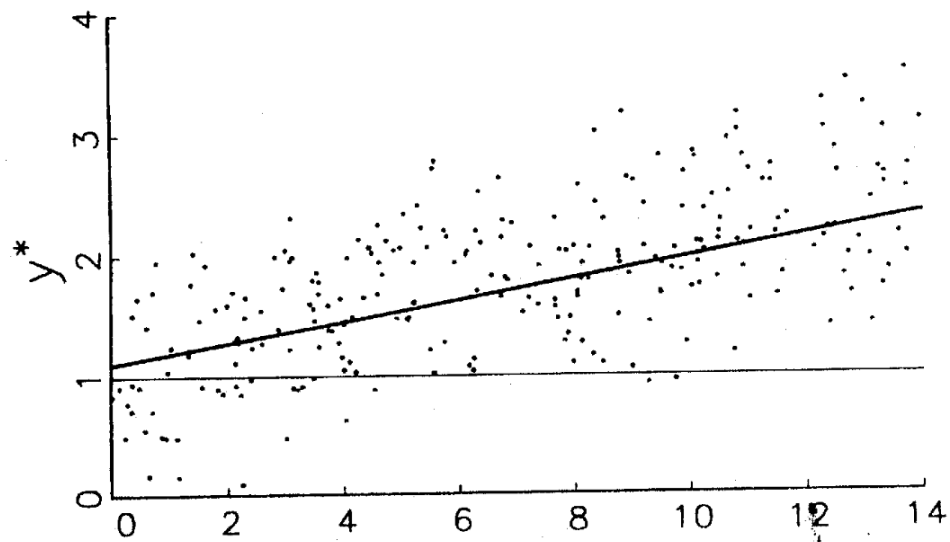
Las variables dependientes censuradas y_i^* son aquellas donde a pesar de observar todo el dataset, sólo tenemos información parcial de y_i . El típico ejemplo es cuando en las encuestas preguntan por los ingresos, y las escalas están truncadas: *Menos de \$100* (1) o *Más de \$1.000* (10). En el primer caso, respondentes que ganan \$25 serán codificados como “1” al igual que alguien que gana \$95. Esto se llama *left-censoring*. A la derecha (*right-censoring*) sería el caso de alguien que gane \$1.500: ella será codificada como “10” al igual que alguien que gane \$999. Nota que observas todo el dataset: es sólo que la estructura de la variable “censura” cierta información, convirtiendo y_i en y_i^* .

Continuaremos motivando el modelo tobit vía modelos latentes. Existe una variable no censurada o truncada y_i que no podemos observar. En vez, observamos la variable censurada y_i^* , donde,

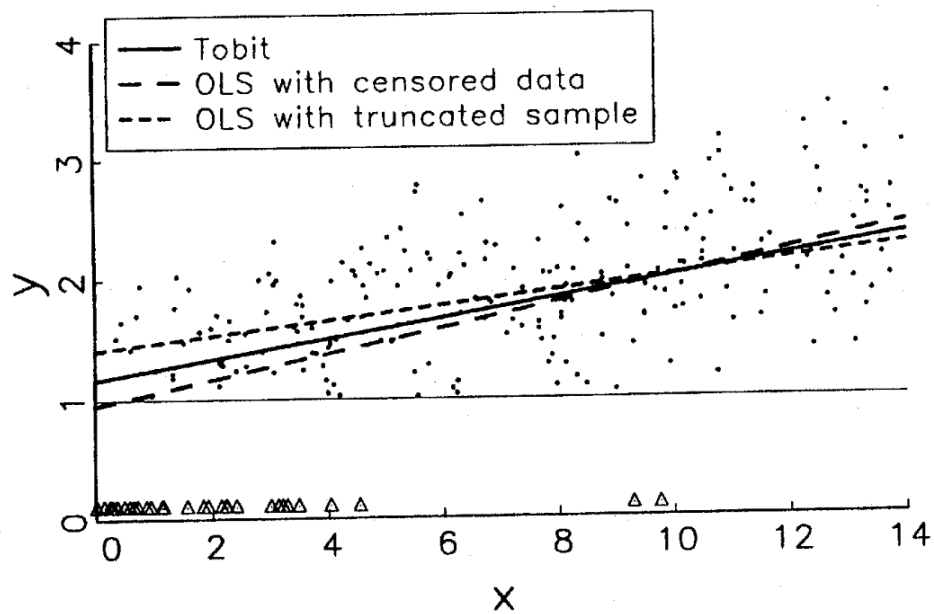
$$y_i = \begin{cases} y_i^* & \text{si } y_i^* > 1 \\ 0 & \text{si } y_i^* \leq 1 \end{cases}$$

En este ejemplo seguimos la Figura 7.2 para el caso específico de que veamos y_i^* sólo cuando $y_i^* \geq 1$, pero *no* veamos nada cuando $y_i^* \leq 1$. Obviamente, el umbral puede variar según la especificación.

Panel A: Regression without Censoring



Panel B: Regression with Censoring and Truncation

**Figure 7.2.** Linear Regression Model With and Without Censoring and Truncation

Como lo muestra Long (1997, p. 189), los coeficientes β estarán mal estimados si se usan métodos lineales. Para ello existe el modelo `tobit`.

Supuestos distribucionales La distribución normal truncada:

$$f(y|y > \tau, \mu, \sigma) = \frac{\frac{1}{\sigma} \Phi\left(\frac{y_i^* - \mu}{\sigma}\right)}{1 - \Phi\left(\frac{\tau - \mu}{\sigma}\right)} \quad (1)$$

donde Φ es el *probability density function* de la distribución normal estándar, y τ es un “threshold” (o “umbral”) bajo el cual se censuran valores de y_i .

Es importante también notar que $\Phi\left(\frac{\tau - \mu}{\sigma}\right)$ en Equation 1 es la probabilidad de que y_i esté censurado (Long 1997, p. 195). Esto implica que la probabilidad de *no* estar censurado es $1 - \Phi\left(\frac{\tau - \mu}{\sigma}\right)$.

Como explica Long (1997, p. 199), existe un link entre el modelo `tobit` y el *probit*. Nota que no sólo ambos modelos estructurales son iguales, sino que también (y tal como lo mostramos en Equation 3), la probabilidad de que una observación sea censurada se calcula usando un modelo parecido al *probit*.

Estimación El modelo estructural `tobit` está dado por algo que ya debiera ser familiar,

$$y_i^* = \mathbf{x}_i \beta + e_i \quad (2)$$

donde e_i es homoeskedástico y $e_i \sim N(0, \sigma^2)$ (Long 1997, p. 206). Esto implica que un modelo `tobit` con *right-censoring* está dado por,

$$y_i = \begin{cases} y_i^* = \mathbf{x}_i \beta + e_i & \text{si } y_i^* < \tau \\ \tau & \text{si } y_i^* \geq \tau \end{cases}$$

Lo interesante de la especificación `tobit` es que modela (mediante supuestos distribucionales) la probabilidad de que un sujeto sea censurado, de manera tal que,

$$Pr(y_i^* \leq \tau | \mathbf{x}_i) = Pr(\epsilon_i \leq \tau - \mathbf{x}_i \beta | \mathbf{x}_i) = \Phi\left(\frac{\tau - \mathbf{x}_i \beta}{\sigma}\right) \quad (3)$$

Esta información es muy útil. Debido a que no sabemos los valores de $y_i^* \leq \tau$, nosotros usamos

la probabilidad de ser censurado como si fuera el likelihood (Long 1997, pp. 204–205). En otras palabras, debido a que no conocemos los valores de $y_i^* \leq \tau$, no podemos usar el alto de la curva normal para estimar el likelihood (ver Figura 7.8). Afortunadamente, debido a que conocemos la igualdad $y_i^* \leq \tau$, y que observamos \mathbf{x} , podemos calcular Equation 3.

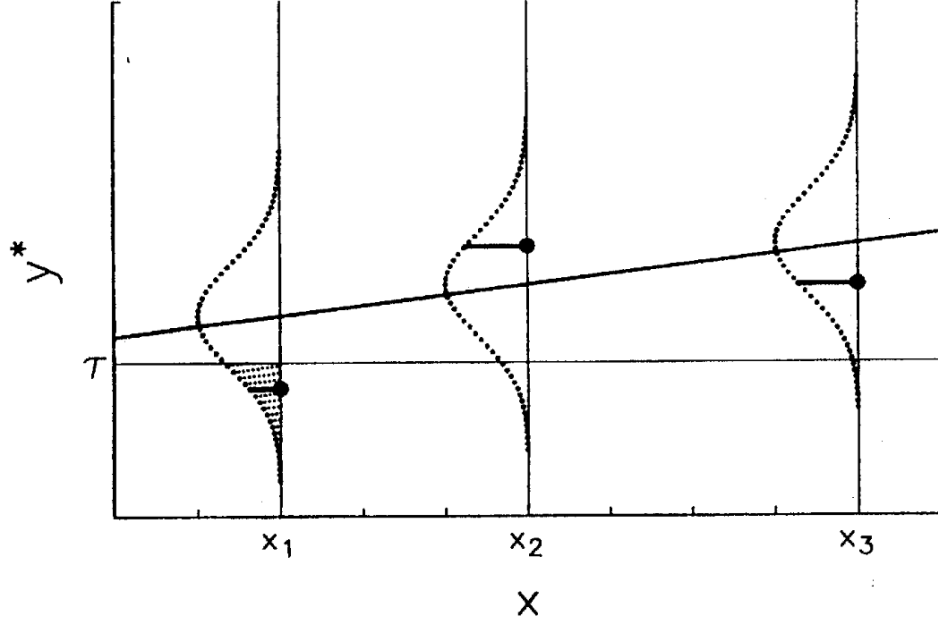


Figure 7.8. Maximum Likelihood Estimation for the Tobit Model

En consecuencia, la función likelihood para *uncensored outcomes* (“U”) está dada por,

$$\ln L_U(\beta, \sigma^2) = \sum_U \ln \frac{1}{\sigma} \Phi - \left(\frac{y_i - \mathbf{x}_i \beta}{\sigma} \right) \quad (4)$$

mientras que la función likelihood para *censored outcomes* (“C”) está dada por,

$$\ln L_C(\beta, \sigma^2) = \sum_C \ln \Phi - \left(\frac{\tau - \mathbf{x}_i \beta}{\sigma} \right) \quad (5)$$

II. Truncated Data

El modelo estructural de los procesos truncados es igual. Lo mismo ocurre con su likelihood. La única diferencia es que el likelihood debe ser ajustado por el área de la distribución normal que ha

sido truncada. **Nota, en consecuencia, que el supuesto de la normalidad es fundamental en el modelo tobit.**

II. PROGRAMACIÓN

Greene (2011, p. 811) explica que la base de datos **Affairs** pregunta el número de relaciones extra maritales en un año.

Carguemos los datos:

```
library(AER)
data("Affairs")
```

Hagamos un resumen de la variable dependiente,

```
table(Affairs$affairs)

##
##   0   1   2   3   7  12
## 451  34  17  19  42  38
```

Como ves, los códigos “7” y “12” están censurados. El primero (“7”) significa 4-10 *affairs*, y “12” una categoría superior.

El paquete de R que usaremos se llama **AER**, y la función se llama **tobit**.

```
tobit.m <- tobit(affairs ~ age + yearsmarried + religiousness + occupation + rating,
               right = 4,
               data = Affairs)
```

Fíjate que hemos declarado que tenemos *right-censoring* desde el código “4”. Hagamos una tabla.

```
p_load(texreg)
texreg(tobit.m) # usa "screenreg" no "texreg".
```

	Model 1
(Intercept)	7.90** (2.80)
age	-0.18* (0.08)
yearsmarried	0.53*** (0.14)
religiousness	-1.62*** (0.42)
occupation	0.32 (0.25)
rating	-2.21*** (0.45)
Log(scale)	2.07*** (0.11)
AIC	1014.09
BIC	1044.88
Log Likelihood	-500.04
Deviance	581.31
Total	601
Left-censored	451
Uncensored	70
Right-censored	80
Wald Test	42.56

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$

Table 1: *Statistical models*

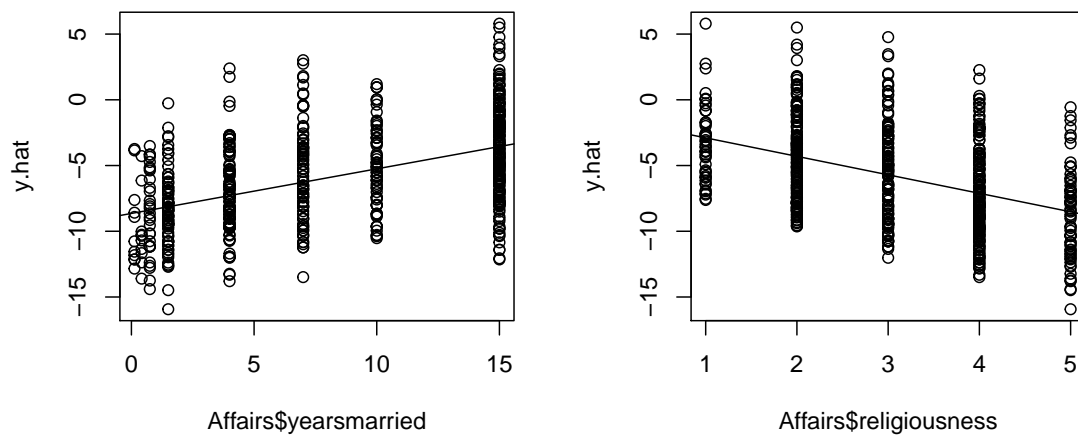
Interpretación Debido a que nuestra variable es continua, los modelos `tobit` se pueden interpretar igual que un modelo OLS. **Interpretemos la tabla.**

Estimemos las *predicted probabilities*,

```
y.hat = predict(tobit.m, interval = "prediction")
```

Ahora hagamos dos gráficos, que por propósitos demostrativos, sólo seleccionaremos dos variables.

```
par(mfrow=c(1,2))
plot(Affairs$yearsmarried, y.hat);abline(lm(y.hat ~ Affairs$yearsmarried))
plot(Affairs$religiousness, y.hat);abline(lm(y.hat ~ Affairs$religiousness))
```



Finalmente, hagamos un diagnóstico rápido. Veamos los residuos. Primero, hagamos el cálculo.

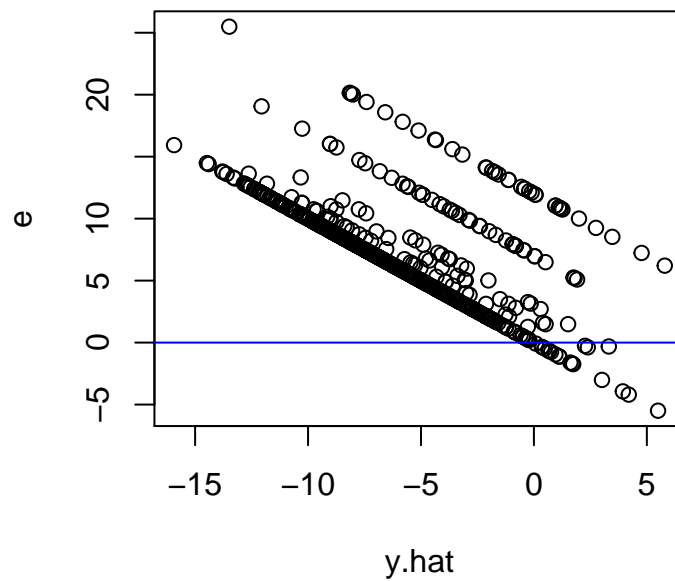
```
e = Affairs$affairs - y.hat
```

Veamos si tienen promedio cero.

```
mean(e)
## [1] 7.326525
```

Ahora veamos el gráfico. Siempre gráfica tus residuos.


```
plot(y.hat, e)
abline(h=0, col="blue")
```



```
knitr::purl('Tobit.Rnw')

## [1] "Tobit.R"

Stangle('Tobit.Rnw')

## Writing to file Tobit.R

## Error in match.arg(options$results, c("verbatim", "tex", "hide")): 'arg' should
be one of "verbatim", "tex", "hide"
```

REFERENCES

- Greene, William (2011). *Econometric Analysis*. 7th. Prentice Hall.
- Long, Scott (1997). *Regression Models for Categorical and Limited Dependent Variables*. Sage, p. 297.