

Model

Tobit model is a regression model where some of the dependent variable y_i^* observations i are censored. For censored observations all we can say is that the variable is less(left censoring) or more(right censoring) than some other value c_i . Let's indicate that variable is left censored with $l_i = 1, r_i = 0$, right censored with $r_i = 1, l_i = 0$ and is not censored when $r_i = 0 = l_i$. Let's call the latent variable before censoring y^* .

$$y_i = \begin{cases} y_i^* & \text{if } r_i = 0 = l_i, \\ c_i & \text{if } l_i = 1, r_i = 0, \\ c_i & \text{if } r_i = 1, l_i = 0 \end{cases} \quad (1)$$

We assume that the underlying generative model is:

$$y_i^* = f(\mathbf{x}_i, \boldsymbol{\beta}) + \epsilon \quad (2)$$

Where \mathbf{x}_i are independent variables and $\boldsymbol{\beta}$ are model parameters and ϵ is normally distributed $\epsilon \sim \mathcal{N}(0, \sigma^2)$. Then y_i^* is also normally distributed $y_i^* \sim \mathcal{N}(f(\mathbf{x}_i, \boldsymbol{\beta}), \sigma^2)$. After observing real data we try to find parameters $\boldsymbol{\beta}$ that fits the data. We do it using maximum likelihood estimation. We have 3 types of observation that will translate to parts of likelihood. For uncensored data we take probability density function at y_i observed values. Let pdf be the probability density function of standard normal distribution

$$\mathcal{L}_{1,i} = pdf\left(\frac{y_i - f(\mathbf{x}_i, \boldsymbol{\beta})}{\sigma}\right) \frac{1}{\sigma} \quad (3)$$

For left censored data we take probability that latent variable is less than observed variable. Let cdf be cumulative distribution function of standard normal distribution

$$\mathcal{L}_{2,i} = \mathcal{P}(y_i^* \leq y_i) = cdf\left(\frac{y_i - f(\mathbf{x}_i, \boldsymbol{\beta})}{\sigma}\right) \quad (4)$$

For right censored data we take probability that latent variable is more than observed variable

$$\mathcal{L}_{3,i} = \mathcal{P}(y_i^* > y_i) = 1 - cdf\left(\frac{y_i - f(\mathbf{x}_i, \boldsymbol{\beta})}{\sigma}\right) \quad (5)$$

Finally putting everything together

$$\mathcal{L} = \prod_{i=1}^N = \mathcal{L}_{1,i}^{(1-l_i)(1-r_i)} * \mathcal{L}_{2,i}^{(l_i)(1-r_i)} * \mathcal{L}_{3,i}^{(1-l_i)(r_i)} \quad (6)$$

Our goal is to find

$$\max_{\boldsymbol{\beta}} \mathcal{L} \quad (7)$$

Because its easier to work with sum and adding and because logarithm is monotone increasing we can solve equivalent problem by taking log of L

$$\max_{\boldsymbol{\beta}} \sum_{i=1}^N [\log(\mathcal{L}_{1,i})(1-l_i)(1-r_i)] * [\log(\mathcal{L}_{2,i})(l_i)(1-r_i)] * [\log(\mathcal{L}_{3,i})(1-l_i)(r_i)] \quad (8)$$