

# 1 Truncation, Censoring, and Selectivity

## 1.1 Truncation

Consider the case of a regression model with a truncated sample. We assume

$$y_i = X_i\beta + u_i,$$

where  $u_i$  is normally distributed with variance  $\sigma$  and the “OLS-assumptions” are satisfied. Data with  $y_i > K$  are discarded for some number  $K$  (which is often normalized to 0 in textbooks).

The probability that an observation from a distribution with density  $f$  is in a small interval of length  $\Delta y$  around  $y_i$  is  $f(y_i)\Delta y$ . (Strictly speaking it would be  $\int_{y_i - \frac{\Delta y}{2}}^{y_i + \frac{\Delta y}{2}} f(y)dy$ ). While densities are not probabilities, it is much easier to use the shorthand of talking about the probability of  $y_i$ . So, because we only have a truncated sample the probability of observing  $y_i$  in the truncated sample is the unconditional probability divided by the probability that  $y < K$  as an application of  $P(A|B) = P(A \cap B)/P(B)$ . Here,  $A$  is  $[y_i - \Delta y/2, y_i + \Delta y/2]$  and the probability of  $A$  is  $f(y_i)\Delta y = f(y_i - X_i\beta)\Delta y$  when  $y_i < K$  (so  $B$  here is the set  $y_i < K$  and the density is zero outside the set  $B$ ). For  $f$  being the normal density (with  $\phi$  denoting the standard normal) we have the probability in the numerator being  $\frac{1}{\sigma}\phi(\frac{y_i - X_i\beta}{\sigma})$ . We have that the probability in the denominator is  $P(y_i < K) = P(X_i\beta + u_i < K) = \Phi(\frac{K - X_i\beta}{\sigma})$ . In total we have the truncated density (the limit of  $\Delta y$  going to zero) for observation  $i$ :

$$\phi\left(\frac{y_i - X_i\beta}{\sigma}\right) / \left(\sigma\Phi\left(\frac{K - X_i\beta}{\sigma}\right)\right).$$

As you can convince yourself, this is a density (positive and integrating to unity). The log likelihood function (ignoring the  $\pi$  term) is

$$\sum_{i=1}^N -0.5 \log \sigma^2 - 0.5 \frac{(y_i - X_i\beta)^2}{\sigma^2} - \log \Phi\left(\frac{K - X_i\beta}{\sigma}\right).$$

## 1.2 Censoring

Consider the case of a regression model with a censored . We assume

$$y_i^0 = X_i\beta + u_i,$$

where  $u_i$  is normally distributed with variance  $\sigma$  and the “OLS-assumptions” are satisfied. Data with  $y_i^0 > K$  are transformed to  $y_i = K$ .

Here I managed to make it confusing in the 9/11 lecture. The probability that an observation is in a small interval of length  $\Delta y = \Delta u$  around  $y_i$  is  $f(y_i)\Delta y$  is NOT conditional.  $f(u_i)\Delta u = f(y_i - X_i\beta)\Delta u$  is the probability of being in the  $\Delta y$  interval when  $y < K$ . The only other value  $y$  can take is  $y_i = K$  and the probability of this is  $P(y_i^0 > K) = 1 - \Phi(\frac{K - X_i\beta}{\sigma})$ .

The log-likelihood function is therefore

$$\sum_{i=1}^N I\{y_i < K\} * [-0.5 \log \sigma^2 - 0.5 \frac{(y_i - X_i\beta)^2}{\sigma^2}] + I\{y_i = K\} * \log (1 - \Phi(\frac{K - X_i\beta}{\sigma})) .$$

Davidson and MacKinnon point out that you can add and subtract  $\log \Phi(\frac{K - X_i\beta}{\sigma})$ , in which case the likelihood has the form of a sum of a truncated likelihood and a Probit likelihood (with  $\sigma$  identified from the first part). Conceptually this is writing the first part as  $P(A) = P(A|B)P(B)$  where  $A$  here is the probability of falling in a small interval around  $y_i$  and  $B$  is the event  $y_i < K$ . This is not important and may instead be confusing.

## 1.3 Selection

The general normal selection model is one where  $y$  is observed based on some outcome  $z$  which we model as a Probit. There are a huge number of applications of this. Say,  $y$  is the GPA of a student at, say, Rice, and  $z$  is literally the probability of getting selected (admitted) [ignore that students may decline]. U.S. college admission typically depend on a large number of variables such as which state you came from, whether your parents are alumni, and on and on. Assuming you have a sample of students, you might have data for many of these variables but not others. For example, you would likely

not observe the quality of the student's essay and this would go into the error term in the admissions equation. If the quality of the students essay also is correlated with the students performance, you would have more efficient inference taking into account that now the errors in the GPA equation is correlated with the error in the selection equation. More importantly, you will get bias if you do not control for this.

Assume that

$$y_i^0 = X_i\beta + u_i,$$

and

$$z_i^0 = W_i\gamma + v_i.$$

We assume the error terms are normal and independent across individuals (or whatever the  $i$  index stands for). We observe

$$z_i = 1 \text{ if } z_i^0 > 0; \text{ 0 otherwise}$$

and

$$y_i = y_i^0 \text{ if } z_i = 1.$$

If individual  $i$  is not selected, we do not observe  $y_i$ .

Denote the variance of  $u_i$  by  $\sigma^2$ . The variance of  $v_i$  is (as usual for a Probit model) not identified and it is normalized to 1. The covariance of two random variables can always be written as the correlation times the standard deviations of the variables. Here, it is convenient to label the correlation  $\rho$  and the covariance is then  $\rho\sigma$ .

We want to study the distribution of  $y_i$  conditional on  $z_i = 1$ . We therefore study the conditional distribution of  $u_i$  conditional on  $z_i = 1$ . This is somewhat difficult because  $z_i = 1$  is a set of  $v_i$ 's. The trick therefore is to use the identity

$$P(A|B) = P(B|A)P(A)/P(B).$$

In our application we want to write  $P(u_i|z_i = 1)$  as  $P(z_i = 1|u_i)P(u_i)/P(z_i = 1)$  because we can easily find the three terms involved.  $P(u_i)$  is just the normal density

and  $P(z_i = 1)$  is a Probit probability. That last term involves  $z_i$  which is a function of  $z_i^0$  and we know how to find conditional normals. The mean of  $z_i^0$  conditional on  $u_i$  is  $W_i\gamma + \frac{\rho\sigma}{\sigma^2}(u_i - 0) = W_i\gamma + \frac{\rho}{\sigma}u_i$ , by the usual formula for normal conditionals, and the conditional variance is  $1 - \frac{(\rho\sigma)^2}{\sigma^2} = 1 - \rho^2$ . So

$$P(u_i = 1|z_i) = \Phi\left(\frac{W_i\gamma + \frac{\rho}{\sigma}(y_i - X_i\beta)}{\sqrt{1 - \rho^2}}\right).$$

Now we can write the density for  $y_i$  conditional on  $z_i = 1$  as

$$P(z_i = 1|u_i) = \Phi\left(\frac{W_i\gamma + \frac{\rho}{\sigma}(y_i - X_i\beta)}{\sqrt{1 - \rho^2}}\right) * \frac{1}{\sigma}\phi\left(\frac{y_i - X_i\beta}{\sigma}\right)/\Phi(W_i\gamma).$$

The full likelihood becomes this probability  $P(u_i|z_i = 1)$  times the probability  $P(z_i = 1)$  “plus” the probability  $z_i = 0$ ; i.e.:

$$I(z_i = 1) * \left[\Phi\left(\frac{W_i\gamma + \frac{\rho}{\sigma}(y_i - X_i\beta)}{\sqrt{1 - \rho^2}}\right) * \frac{1}{\sigma}\phi\left(\frac{y_i - X_i\beta}{\sigma}\right)\right] + I(z_i = 0) * [1 - \Phi(W_i\gamma)].$$

The log-likelihood is after re-ordering a bit:

$$\sum_{i=1}^N I(z_i = 1) * \left[-0.5 \log \sigma^2 - 0.5 \frac{(y_i - X_i\beta)^2}{\sigma^2} + \log \Phi\left(\frac{W_i\gamma + \frac{\rho}{\sigma}(y_i - X_i\beta)}{\sqrt{1 - \rho^2}}\right)\right] + I(z_i = 0) * \log \Phi(-W_i\gamma).$$

Notice what happens if  $\rho = 0$ : you have a Probit model and an independent normal which you can estimate by least squares. True, it is strange that  $y$  is only observed when  $z = 1$ , but you do not have to adjust the least squares estimation in the case where the error term in the selection equation is not affecting the error term in the regression.

### 1.3.1 Heckman correction term for selection

Heckman was the first to consider correction for selection (in his thesis, I think) and this was the basis for his later Nobel prize.

The Heckman correction involves as two-step estimator. Assume you first estimate the Probit equation and then the regression (not recommended—it is always most efficient

to estimate the full system, but sometimes we do anyway, at least in a first exploration and earlier it may have been hard numerically to estimate the full system).

Consider the regression

$$y_i = X_i\beta + u_i,$$

where you ignore that  $y_i$  has been selected based on  $z$ . The problem now is that  $E u_i = 0$  if you observed all outcomes (including the ones that were not selected) but if  $E(u_i|v_i) = \rho\sigma v_i$ , which is easy to see using the standard formula for conditional normals, we are allowed to write  $u_i$  as  $\rho\sigma v_i + e_i$  where  $e_i = u_i - \rho\sigma v_i$  is independent of  $v_i$ . We have

$$y_i = X_i\beta + \rho\sigma v_i + e_i,$$

where  $e_i$  is independent of  $X_i$  but because  $v_i$  is instrumental in deciding whether  $y_i$  was observed, it is unlikely to have mean zero. In fact,  $E(v_i|z_i = 1) = E(v_i|W_i\gamma + v_i > 0) = \frac{\phi(W_i\gamma)}{\Phi(W_i\gamma)}$ , where ratio is called the inverse Mill's Ratio.<sup>1</sup> If you have estimated the first state you have an estimate  $\hat{\gamma}$  and you run the regression

$$y_i = X_i\beta + \kappa \frac{\phi(W_i\hat{\gamma})}{\Phi(W_i\hat{\gamma})} + e_i,$$

which is a consistent estimator of  $\beta$  (and approximately unbiased if  $\gamma$  is well estimated). Usually, economists do not attempt to extract the parameters  $\rho$  and  $\sigma$ .

---

<sup>1</sup>To derive the inverse Mill's ratio notice that  $\int x \exp(-x^2/2)dx = \int \exp(-x^2/2)(xdx)$  and do a change of variables to  $y = \frac{x^2}{2}$  with  $dy = xdx$ .