

Selection, Treatment, and Switching Models

KPT, Ch. 14

Matlab scripts: `mod9_select_data`, `mod9_select`, `mod9_select_door2door`,
`mod9_treat_data`, `mod9_treat`, `mod9_treat_naive`,
`mod9_treatJobTrain`, `mod9_treatJobTrainNaive`,
`mod9_switch_data`, `mod9_switch`

Matlab functions: `gs_select`, `gs_treat`, `gs_switch`, `iw_sig11`

We will encounter two new estimation challenges for these models: A restricted covariance matrix, and drawing latent data from multivariate densities with inequality restrictions.

Sample Selection Model

Basic features

The basic sample selection model (often also referred to as “Heckman selection model” or “Type II Tobit”) is structured as follows:

$$\begin{aligned}
 y_{1i}^* &= \mathbf{x}'_{1i} \boldsymbol{\beta}_1 + \varepsilon_{1i} & y_{2i}^* &= \mathbf{x}'_{2i} \boldsymbol{\beta}_2 + \varepsilon_{2i} \\
 \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} &\sim n\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}\right) & \boldsymbol{\Sigma} &= \begin{bmatrix} 1 & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix} \\
 y_{1i} &= 1 & \text{if } & y_{1i}^* > 0, & = 0 & \text{otherwise} \\
 y_{2i} &= y_{2i}^* & \text{if } & y_{1i}^* > 0, & \text{unobserved} & \text{otherwise}
 \end{aligned} \tag{1}$$

Since the outcome for the first equation (often called “participation equation” or “hurdle equation”) is only observed in binary form (as in the Probit), the variance for the first error term is not identified and arbitrarily set to one. This generates a constrained covariance matrix for $\boldsymbol{\varepsilon}_i = [\varepsilon_{1i} \ \varepsilon_{2i}]'$, which renders a generic *IW* density (as used e.g. for the SUR model) unsuitable for obtaining draws of $\boldsymbol{\Sigma}$. We will deal with this issue below.

Conceptually, the model structure in (1) arises because there are unobservables in the hurdle and outcome equation that are correlated. A classic example would be a woman’s decision to join the labor force, which may be a function of unobserved innate drive, ability, or “spunk”. This quality will also affect the outcome of interest, e.g. hourly wage. If this link between the two equations is ignored, results flowing solely from the outcome equation will suffer from sample selection bias.

For example, let’s assume our interest is in the marginal effects of regressors on (potential) female wage for the entire population of female workers (including women who may not currently be in the labor market). Thus, we are primarily interested in $\boldsymbol{\beta}_2$. Ignoring the selection problem, we would naively use

the second equation only, run a basic normal regression model, and estimate β_2 . Under sample selection, i.e. if $\sigma_{12} \neq 0$, these estimates will be biased. Intuitively, this bias arises because under sample selection women who work share an unobserved quality that increases the likelihood of passing the first hurdle. The sub-sample that's left for the second equation is thus systematically different for the general population, and will be affected differently by regressors than those who don't. Econometrically, the bias can be best illustrated when comparing conditional regression means or expectations.

For the naïve model we have

$$E(y_{2i}^*) = \mathbf{x}'_{2i} \beta_2 = \omega_{2i} \quad (2)$$

In reality, however, the second equation is based on a conditional expectation, i.e

$$\begin{aligned} E[y_{2i}^* | y_{1i}^* > 0] &= E[\omega_{2i} + \varepsilon_{2i} | y_{1i}^* > 0] = \omega_{2i} + E[\varepsilon_{2i} | y_{1i}^* > 0] = \\ &= \omega_{2i} + E[\varepsilon_{2i} | \varepsilon_{1i} > -\omega_{1i}] = \omega_{2i} + \sigma_{12} \frac{\phi(\omega_{1i})}{\Phi(\omega_{1i})} \quad \text{where } \omega_{1i} = \mathbf{x}'_{1i} \beta_1 \end{aligned} \quad (3)$$

The fraction in the last term (Inverse Mills ratio, IMR) is always positive. Thus, the relative magnitude of the conditional vs. the unconditional mean depends on the sign of the covariance. If $\sigma_{12} > 0$ the conditional mean exceeds the unconditional mean. Intuition for the labor market example: In that case an increased propensity to join the labor force is positively correlated with an increased wage. Since we condition the distribution of y_{2i}^* exclusively on positive outcomes for the hurdle, we would naturally expect the conditional mean to be higher. The opposite holds if $\sigma_{12} < 0$. In either case, the estimates for β_2 will adjust for the added term, and thus be different from the estimates generated by the naïve model. Specifically, if $\sigma_{12} > 0$ we would expect the naïve estimate for β_2 to be biased upwards to make up for the missing positive term, and downwards if $\sigma_{12} < 0$.

Thus, a key focus in the estimation of the sample selection model will lie on the covariance σ_{12} . A natural “restricted” candidate model would be one with $\sigma_{12} = 0$, in which case the two equations are independent. For our labor market example this would imply that women that work are no different (based on unobservables in equation 1) than women who don't, and we can interpret any inference flowing from equation 2 to hold for the general population.

Another important feature of the sample selection model is the “identification restriction”, which stipulates that the first equation must include at least one explanatory variable that is not also present in the second equation. This can often be tricky to satisfy as there is substantial overlap between \mathbf{x}_{1i} and \mathbf{x}_{2i} in many applications.

Estimation

As for the Probit and Tobit models we assume that explanatory data are observed for all individuals, regardless of the outcome for the first hurdle. Our estimation strategy will have features of the SUR model (since we have multiple linked equations per individual) and of the Probit and Tobit models (since we will work with latent dependent variables).

Denoting the number of regressors in the first and second equation as k_1 and k_2 , respectively, we can write the pair of latent observations for a given individual as:

$$\mathbf{y}_i^* = \mathbf{X}_i \boldsymbol{\beta} + \boldsymbol{\varepsilon}_i \quad \text{with} \quad \boldsymbol{\varepsilon}_i \sim n(\mathbf{0}, \boldsymbol{\Sigma}),$$

$$\mathbf{X}_i = \begin{bmatrix} \mathbf{x}'_{1i} & \mathbf{0} \\ \mathbf{0} & \mathbf{x}'_{2i} \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{bmatrix} \quad \boldsymbol{\varepsilon}_i = \begin{bmatrix} \varepsilon_{1i} \\ \varepsilon_{2i} \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix} \quad (4)$$

In analogy to the full model over all individuals can be written as

$$\mathbf{y}^* = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\mathbf{y}^* = \begin{bmatrix} \mathbf{y}_1^* \\ \mathbf{y}_1^* \\ \vdots \\ \mathbf{y}_n^* \end{bmatrix}_{(nx2) \times 1} \quad \mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_n \end{bmatrix}_{(nx2) \times k} \quad \boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{bmatrix} \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \boldsymbol{\varepsilon}_1 \\ \boldsymbol{\varepsilon}_2 \\ \vdots \\ \boldsymbol{\varepsilon}_n \end{bmatrix} \sim n(\mathbf{0}, \boldsymbol{\Omega}) \quad \text{with} \quad (5)$$

$$\boldsymbol{\Omega} = \begin{bmatrix} \boldsymbol{\Sigma} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma} & \dots & \mathbf{0} \\ \mathbf{0} & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \boldsymbol{\Sigma} \end{bmatrix}_{(nx2) \times (nx2)}$$

To construct the likelihood function for this model, we need to review some basic results for joint densities with inequality restrictions for one or more equations. Generically, the truncated joint density for two random variables y and z can be written as

$$f(y < a, z) = \int_{-\infty}^a f(y, z) dy = \int_{-\infty}^a f(y|z) f(z) dz = f(z) \int_{-\infty}^a f(y|z) dy \quad (6)$$

Also, if y and z are bivariate normal, the conditional densities are normal as well with the following moments:

$$z|y \sim n(\mu_{z,y}, \sigma_{z,y}^2) \quad \mu_{z,y} = \mu_z + \frac{\sigma_{yz}}{\sigma_y^2}(y - \mu_y) \quad \sigma_{z,y}^2 = \sigma_z^2 - \frac{\sigma_{yz}^2}{\sigma_y^2}$$

$$y|z \sim n(\mu_{y,z}, \sigma_{y,z}^2) \quad \mu_{y,z} = \mu_y + \frac{\sigma_{yz}}{\sigma_z^2}(z - \mu_z) \quad \sigma_{y,z}^2 = \sigma_y^2 - \frac{\sigma_{yz}^2}{\sigma_z^2} \quad (7)$$

For our sample likelihood there are *two possible observable outcomes*, $pr(y_{1i}=1, y_{2i})$ and $pr(y_{1i}=0)$. The second one can be immediately expressed as

$$pr(y_{1i}=0) = pr(y_{1i}^* \leq 0) = pr(\varepsilon_{1i} \leq -\mathbf{x}'_{1i}\boldsymbol{\beta}_1) = \Phi(-\mathbf{x}'_{1i}\boldsymbol{\beta}_1) \quad (8)$$

The first one takes a bit of extra work to bring it into a manageable form. We have

$$\begin{aligned} pr(y_{1i}=1, y_{2i}) &= f(\varepsilon_{1i} > -\mathbf{x}'_{1i}\boldsymbol{\beta}_1, \varepsilon_{2i}) = \\ &\int_{-\mathbf{x}'_{1i}\boldsymbol{\beta}_1}^{\infty} f(\varepsilon_{1i}, \varepsilon_{2i}) d\varepsilon_{1i} = \int_{-\mathbf{x}'_{1i}\boldsymbol{\beta}_1}^{\infty} f(\varepsilon_{1i} | \varepsilon_{2i}) f(\varepsilon_{2i}) d\varepsilon_{1i} = \\ &f(\varepsilon_{2i}) \int_{-\mathbf{x}'_{1i}\boldsymbol{\beta}_1}^{\infty} f(\varepsilon_{1i} | \varepsilon_{2i}) d\varepsilon_{1i} = \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{1}{2\sigma_2^2}(y_{2i} - \mathbf{x}'_{2i}\boldsymbol{\beta}_2)^2\right) \int_{-\mathbf{x}'_{1i}\boldsymbol{\beta}_1}^{\infty} f(\varepsilon_{1i} | \varepsilon_{2i}) d\varepsilon_{1i} \end{aligned} \quad (9)$$

where we use $\varepsilon_{2i} = y_{2i} - \mathbf{x}'_{2i}\boldsymbol{\beta}_2$ in the last expression. To solve $\int_{-\mathbf{x}'_{1i}\boldsymbol{\beta}_1}^{\infty} f(\varepsilon_{1i} | \varepsilon_{2i}) d\varepsilon_{1i}$, we note that

$$\begin{aligned} \int_{-\mathbf{x}'_{1i}\boldsymbol{\beta}_1}^{\infty} f(\varepsilon_{1i} | \varepsilon_{2i}) d\varepsilon_{1i} &= pr((\varepsilon_{1i} | \varepsilon_{2i}) > -\mathbf{x}'_{1i}\boldsymbol{\beta}_1) = pr\left(\frac{(\varepsilon_{1i} | \varepsilon_{2i}) - \mu_{1,2}}{\sigma_{1,2}} < \frac{\mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \mu_{1,2}}{\sigma_{1,2}}\right) = \\ &\Phi\left(\frac{\mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \mu_{1,2}}{\sigma_{1,2}}\right) \quad \text{where} \\ \mu_{1,2} &= 0 + \frac{\sigma_{12}}{\sigma_2^2}(\varepsilon_{2i} - 0) = \frac{\sigma_{12}}{\sigma_2^2}(y_{2i} - \mathbf{x}'_{2i}\boldsymbol{\beta}_2) \quad \sigma_{1,2} = \sqrt{1 - \frac{\sigma_{12}^2}{\sigma_2^2}} = \sqrt{1 - \rho^2} \end{aligned} \quad (10)$$

We can now write the sample likelihood concisely as

$$p(\mathbf{y} | \boldsymbol{\beta}, \boldsymbol{\Sigma}) = \prod_{i=1}^n (\Phi(-\mathbf{x}'_{1i}\boldsymbol{\beta}_1))^{1-y_{1i}} * \left(\frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{1}{2\sigma_2^2}(y_{2i} - \mathbf{x}'_{2i}\boldsymbol{\beta}_2)^2\right) \Phi\left(\frac{\mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \frac{\sigma_{12}}{\sigma_2^2}(y_{2i} - \mathbf{x}'_{2i}\boldsymbol{\beta}_2)}{\sqrt{1 - \rho^2}}\right) \right)^{y_{1i}} \quad (11)$$

The next step would be to assign priors to $\boldsymbol{\beta}$ and $\boldsymbol{\Sigma}$. However, this is where we run into the problem of a constrained covariance matrix. Several approaches have been suggested to handle this issue. For this simple problem with a single constrained variance, the *algorithm proposed by Nobile (2000)* provides fast and accurate results. Without going into the details for this algorithm, we can simply note that it allows you to draw from the Inverse Wishart as usual. The variance constraint is “automatically” taken care of. Thus, if you want to specify an IW prior with degrees of freedom ν_0 and scale matrix \mathbf{S}_0 , with the (1,1) element constrained to a constant (here “1”) you can use the Matlab function `iw_sig11(ν_0 , inv(\mathbf{S}_0), 1)`.

Thus, we can write our priors as

$p(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = p(\boldsymbol{\beta}) p(\boldsymbol{\Sigma})$ where

$$p(\boldsymbol{\beta}) = (2\pi)^{-k/2} |\mathbf{V}_0|^{-1/2} \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0)\right) \quad (12)$$

$$p(\boldsymbol{\Sigma}) = \left(2^{v_0} \pi^{1/2} \prod_{i=1}^2 \Gamma\left(\frac{v_0 + 1 - i}{2}\right)\right)^{-1} * |\mathbf{S}_0|^{v_0/2} |\boldsymbol{\Sigma}|^{-(v_0+3)/2} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{S}_0 \cdot \boldsymbol{\Sigma}^{-1})\right) I(\boldsymbol{\Sigma}_{11} = 1)$$

Setting $M = 2$ in the *IW* density. Combining the priors with the likelihood, and dropping all terms that are multiplicatively unrelated to our parameters of interest yields the posterior kernel

$$p(\boldsymbol{\beta}, \boldsymbol{\Sigma} | \mathbf{y}, \mathbf{X}) \propto |\boldsymbol{\Sigma}|^{-(v_0+3)/2} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{S}_0 \cdot \boldsymbol{\Sigma}^{-1})\right) \exp\left(-\frac{1}{2} \left((\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0)\right)\right) * \quad (13)$$

$$\prod_{i=1}^n (\Phi(-\mathbf{x}'_{1i} \boldsymbol{\beta}_1))^{1-y_{1i}} * \left(\Phi\left(\frac{\mathbf{x}'_{1i} \boldsymbol{\beta}_1 + \frac{\sigma_{12}}{\sigma_2^2} (y_{2i} - \mathbf{x}'_{2i} \boldsymbol{\beta}_2)}{\sqrt{1 - \rho^2}}\right)\right)^{y_{1i}}$$

As you can imagine by now, this would not yield “well-behaved” conditional posterior densities for either $\boldsymbol{\beta}$ or $\boldsymbol{\Sigma}$, since both are represented in the “unruly” likelihood function. We will turn again to the vector of latent outcomes, \mathbf{y}^* , for help. Conceptually, the joint augmented posterior kernel can be expressed as

$$p(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{y}^* | \mathbf{y}) \propto p(\boldsymbol{\beta}) p(\boldsymbol{\Sigma}) p(\mathbf{y}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}) p(\mathbf{y} | \mathbf{y}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}) \quad (14)$$

As for the Probit, Tobit, and OP, we need to take a closer look at the terms $p(\mathbf{y}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma})$ and $p(\mathbf{y} | \mathbf{y}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma})$ to assure that the augmentation will indeed make our life easier. The conditional prior $p(\mathbf{y}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma})$ describes essentially the likelihood of a basic SUR model, i.e.

$$p(\mathbf{y}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}) = \prod_{i=1}^n (2\pi)^{-1} |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})\right) = \quad (15)$$

$$(2\pi)^{-n} |\boldsymbol{\Sigma}|^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})\right)$$

The conditional likelihood function (the last term in (14)) can be written as

$$p(\mathbf{y} | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{y}^*) = p(\mathbf{y} | \mathbf{y}^*) = \prod_{i=1}^n \left(I(y_{1i} = 0) I(-\infty < y_{1i}^* \leq 0) + I(y_{1i} = 1, y_{2i} = y_{2i}^*) I(y_{1i}^* > 0, y_{2i}^*) \right) \quad (16)$$

As for the Probit and Tobit, our data augmentation step de-links the original likelihood from our main parameters of interest, which facilitates their drawing in the GS.

The augmented posterior kernel can now be written in its explicit form as

$$\begin{aligned}
& p(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{y}^* | \mathbf{y}, \mathbf{X}) \propto \\
& |\boldsymbol{\Sigma}|^{-(v_0+3+n)/2} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{S}_0 \cdot \boldsymbol{\Sigma}^{-1})\right) I(\Sigma_{11} = 1) * \exp\left(-\frac{1}{2} \left((\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0) \right)\right) * \\
& \exp\left(-\frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})\right) * \\
& \prod_{i=1}^n \left(I(y_{i1} = 0) I(-\infty < y_{i1}^* \leq 0) + I(y_{i1} = 1, y_{i2} = y_{i2}^*) I(y_{i1}^* > 0, y_{i2}^*) \right)
\end{aligned} \tag{17}$$

The conditional posterior kernel for $\boldsymbol{\beta}$ now takes the following form:

$$p(\boldsymbol{\beta} | \boldsymbol{\Sigma}, \mathbf{y}^*, \mathbf{X}) \propto \exp\left(-\frac{1}{2} \left((\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0) + \sum_{i=1}^n (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta}) \right)\right) \tag{18}$$

This is equivalent to the conditional posterior for the basic SUR model, and we can immediately derive the conditional posterior moments as:

$$\boldsymbol{\beta} | \boldsymbol{\Sigma}, \mathbf{y}^*, \mathbf{X} \sim n(\boldsymbol{\mu}_1, \mathbf{V}_1) \quad \text{with} \quad \mathbf{V}_1 = \left(\mathbf{V}_0^{-1} + \sum_{i=1}^n \mathbf{X}_i' \boldsymbol{\Sigma}^{-1} \mathbf{X}_i \right)^{-1} \quad \text{and} \quad \boldsymbol{\mu}_1 = \mathbf{V}_1 \left(\mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \sum_{i=1}^n \mathbf{X}_i' \boldsymbol{\Sigma}^{-1} \mathbf{y}_i^* \right) \tag{19}$$

For the conditional posterior of $\boldsymbol{\Sigma}$ we have

$$\begin{aligned}
& p(\boldsymbol{\Sigma} | \boldsymbol{\beta}, \mathbf{y}^*, \mathbf{X}) \propto \\
& |\boldsymbol{\Sigma}|^{-(v_0+3+n)/2} \exp\left(-\frac{1}{2} \left(\text{tr}(\mathbf{S}_0 \cdot \boldsymbol{\Sigma}^{-1}) + \sum_{i=1}^n (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta}) \right)\right) I(\Sigma_{11} = 1)
\end{aligned} \tag{20}$$

This is again analogous to the SUR model (except for the variance restriction), leading to

$$\begin{aligned}
& \boldsymbol{\Sigma} | \boldsymbol{\beta}, \mathbf{y}^*, \mathbf{X} \sim IW(v_1, \mathbf{S}_1) I(\Sigma_{11} = 1) \quad \text{with} \\
& v_1 = v_0 + n \quad \mathbf{S}_1 = \mathbf{S}_0 + \sum_{i=1}^n (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})(\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})'
\end{aligned} \tag{21}$$

The conditional posterior kernel for latent data vector \mathbf{y}^* is given by

$$\begin{aligned}
& p(\mathbf{y}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{y}, \mathbf{X}) \propto \\
& \exp\left(-\frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})\right) * \\
& \prod_{i=1}^n \left(I(y_{i1} = 0) I(-\infty < y_{i1}^* \leq 0) + I(y_{i1} = 1, y_{i2} = y_{i2}^*) I(y_{i1}^* > 0, y_{i2}^*) \right)
\end{aligned} \tag{22}$$

At the individual level this implies

$$\begin{aligned}
& p(\mathbf{y}_i^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{y}, \mathbf{X}_i) \propto \\
& \exp\left(-\frac{1}{2}(\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1}(\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})\right)^* \\
& \left(I(y_{i1} = 0)I(-\infty < y_{i1}^* \leq 0) + I(y_{i1} = 1, y_{i2} = y_{i2}^*)I(y_{i1}^* > 0, y_{i2}^*)\right)
\end{aligned} \tag{23}$$

For the $y_{i1} = 0$ case we have

$$\begin{aligned}
& p(\mathbf{y}_i^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, y_{i1} = 0, \mathbf{X}_i) \propto \\
& \exp\left(-\frac{1}{2}(\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1}(\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})\right)^* I(-\infty < y_{i1}^* \leq 0)
\end{aligned} \tag{24}$$

This implies a draw of $\mathbf{y}_i^* = \begin{bmatrix} y_{i1}^* & y_{i2}^* \end{bmatrix}'$ from a joint normal density, with y_{i1}^* truncated to the negative domain, i.e. from

$$p(y_{i1}^* < 0, y_{i2}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) = \int_{-\infty}^0 f(y_{i1}^*, y_{i2}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) dy_{i1}^* \tag{25}$$

Where $f(\cdot)$ denotes the bivariate normal density. This can be computationally implemented as follows:

1. Using $p(y_{i1}^* < 0, y_{i2}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) = p(y_{i2}^* | y_{i1}^* < 0, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) p(y_{i1}^* < 0 | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i)$, we can first draw y_{i1}^* from its marginal normal density with mean $\mathbf{x}'_{i1}\boldsymbol{\beta}_1$ and variance 1, truncated from above at zero.
2. To obtain draws from $p(y_{i2}^* | y_{i1}^* < 0, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i)$ first note that $y_{i2}^* | y_{i1}^* \sim n(\mathbf{x}'_{i2}\boldsymbol{\beta}_2 + \sigma_{12}(y_{i1}^* - \mathbf{x}'_{i1}\boldsymbol{\beta}_1), \sigma_{22}(1 - \rho^2))$ where y_{i1}^* are the draws we just obtained in the previous step. Then draw y_{i2}^* from this *untruncated* distribution (no truncation is needed here, since, conditional on $y_{i1}^* < 0$, y_{i2}^* can be anything!)

For the $y_{i1} = 1$ case we have

$$\begin{aligned}
& p(\mathbf{y}_i^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, y_{i1} = 1, \mathbf{X}_i) \propto \\
& \exp\left(-\frac{1}{2}(\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1}(\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})\right)^* I(y_{i1}^* > 0, y_{i2}^*)
\end{aligned} \tag{26}$$

This implies a draw of $\mathbf{y}_i^* = \begin{bmatrix} y_{1i}^* & y_{2i}^* \end{bmatrix}'$ from a joint normal density, with y_{1i}^* truncated to the positive domain, i.e. from

$$p(y_{1i}^* > 0, y_{2i}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) = \int_0^{\infty} f(y_{1i}^*, y_{2i}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) dy_{1i}^* \quad (27)$$

However, since in this case y_{2i} is observed, we have

$$p(y_{1i}^* > 0, y_{2i}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) = p(y_{1i}^* > 0 | y_{2i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) p(y_{2i}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) = p(y_{1i}^* > 0 | y_{2i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i). \text{ Thus, we}$$

simply draw y_{1i}^* from its conditional normal density with mean $\mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \frac{\sigma_{12}}{\sigma_2^2}(y_{2i}^* - \mathbf{x}'_{2i}\boldsymbol{\beta}_2)$ and variance

$1 - \rho^2$, truncated from below at zero.

Predictive Constructs of Interest

The first predictive construct of interest for the basic selection model is usually the probability of mastering the first hurdle, i.e

$$pr(y_{1i}^* > 0 | \mathbf{x}_{1p}, \boldsymbol{\beta}_1) = \Phi(\mathbf{x}'_{1p}\boldsymbol{\beta}_1) = \Phi(\omega_{1p}) \quad (28)$$

Note that even though σ_{12} does not figure in this expression, the resulting PPD for a naïve probit model (ignoring the correlation across equations) will still be misleading since the posterior density of $\boldsymbol{\beta}_1$ will be biased. However, the degree of inferential error for this first predictive outcome of interest is likely smaller than the error for the second outcome, which is a direct function of σ_{12} . It is the expected value for the second equation outcome, conditional on mastering the first equation hurdle, i.e.

$$E[y_{2p} | y_{1p}^* > 0, \mathbf{x}_{1p}, \mathbf{x}_{2p}, \boldsymbol{\beta}, \boldsymbol{\Sigma}] = \omega_{2p} + \sigma_{12} \frac{\phi(\omega_{1p})}{\Phi(\omega_{1p})} \quad \text{where } \omega_{1p} = \mathbf{x}'_{1p}\boldsymbol{\beta}_1, \quad \omega_{2p} = \mathbf{x}'_{2p}\boldsymbol{\beta}_2 \quad (29)$$

As mentioned above, a simple regression model (ignoring both the the link with the first equation AND the left side truncation in the outcome variable) would estimate this as

$$E[y_{2p} | \mathbf{x}_{2p}, \boldsymbol{\beta}_2] = \omega_{2p} \quad (30)$$

A Tobit model would produce (see earlier lecture notes)

$$E[y_{2p} | y_{2p}^* > 0, \mathbf{x}_{2p}, \boldsymbol{\beta}_2] = \omega_{2p} + \sqrt{\sigma_{22}} \frac{\phi(\omega_{2p})}{\Phi(\omega_{2p})} \quad (31)$$

i.e. it would correctly censor y_2 to the left, but attribute the censoring rule to the same regressors as those that affect the conditional mean of the outcome variable.

The third construct that can be of interest in some applications is the expected value of the *latent* variable for the second equation, i.e.

$$E\left[y_{2p}^* \mid \mathbf{x}_{2p}, \boldsymbol{\beta}_2\right] = \omega_{2p} \quad (32)$$

Naturally, a naïve regression model or a Tobit would produce the same expression, but would suffer from biased posterior results for $\boldsymbol{\beta}_2$.

The fourth construct of occasional interest is the marginal effect of a given regressor on the conditional expectation. There are three separate expressions, depending if the regressor appears only in the first, only in the second, or both equations. The corresponding marginal effects (evaluated at the sample mean $\bar{\mathbf{x}}$) are as follows:

$$\begin{aligned} \frac{\partial E\left[y_2 \mid y_1^* > 0, \bar{\mathbf{x}}, \boldsymbol{\beta}, \boldsymbol{\Sigma}\right]}{\partial(x_{1,j})} &= -\beta_{1,j}\sigma_{12} \frac{\phi(\bar{\omega}_1)}{\Phi(\bar{\omega}_1)} \left(\frac{\phi(\bar{\omega}_1)}{\Phi(\bar{\omega}_1)} + \bar{\omega}_1 \right) \\ \frac{\partial E\left[y_2 \mid y_1^* > 0, \bar{\mathbf{x}}, \boldsymbol{\beta}, \boldsymbol{\Sigma}\right]}{\partial(x_{2,j})} &= \beta_{2,j} \\ \frac{\partial E\left[y_2 \mid y_1^* > 0, \bar{\mathbf{x}}, \boldsymbol{\beta}, \boldsymbol{\Sigma}\right]}{\partial(x_{1,j}, x_{2,j})} &= \beta_{2,j} - \beta_{1,j}\sigma_{12} \frac{\phi(\bar{\omega}_1)}{\Phi(\bar{\omega}_1)} \left(\frac{\phi(\bar{\omega}_1)}{\Phi(\bar{\omega}_1)} + \bar{\omega}_1 \right) \quad \text{with} \quad \bar{\omega}_1 = \bar{\mathbf{x}}_1' \boldsymbol{\beta}_1 \end{aligned} \quad (33)$$

Naturally, a Tobit model can only produce $\frac{\partial E\left[y_p \mid y_p^* > 0, \bar{\mathbf{x}}, \boldsymbol{\beta}, \sigma^2\right]}{\partial x_j} = \beta_j \Phi\left(\frac{\bar{\mathbf{x}}_j \boldsymbol{\beta}_j}{\sigma}\right)$.

For 0/1 regressors, use

$$\begin{aligned} E\left[y_2 \mid y_1^* > 0, \bar{\mathbf{x}}_{-j}, x_j = 1, \boldsymbol{\beta}, \boldsymbol{\Sigma}\right] - E\left[y_2 \mid y_1^* > 0, \bar{\mathbf{x}}_{-j}, x_j = 0, \boldsymbol{\beta}, \boldsymbol{\Sigma}\right] = \\ (\bar{\omega}_2 \mid x_{2,j} = 1) + \sigma_{12} \frac{\phi(\bar{\omega}_1 \mid x_{1,j} = 1)}{\Phi(\bar{\omega}_1 \mid x_{1,j} = 1)} - (\bar{\omega}_2 \mid x_{2,j} = 0) + \sigma_{12} \frac{\phi(\bar{\omega}_1 \mid x_{1,j} = 0)}{\Phi(\bar{\omega}_1 \mid x_{1,j} = 0)} \end{aligned} \quad (34)$$

In contrast, for the Tobit we used

$$E(y_p | y_p^* > 0, \bar{\mathbf{x}}_{-j}, x_j = 1) - E(y_p | y_p^* > 0, \bar{\mathbf{x}}_{-j}, x_j = 0) = \left(\bar{\mathbf{x}}'\boldsymbol{\beta} + \sigma \frac{\phi\left(\frac{\bar{\mathbf{x}}'\boldsymbol{\beta}}{\sigma}\right)}{\Phi\left(\frac{\bar{\mathbf{x}}'\boldsymbol{\beta}}{\sigma}\right)} \middle| x_j = 1 \right) - \left(\bar{\mathbf{x}}'\boldsymbol{\beta} + \sigma \frac{\phi\left(\frac{\bar{\mathbf{x}}'\boldsymbol{\beta}}{\sigma}\right)}{\Phi\left(\frac{\bar{\mathbf{x}}'\boldsymbol{\beta}}{\sigma}\right)} \middle| x_j = 0 \right) \quad (35)$$

Matlab implementation:

Scripts `mod9_select_data`, and `mod9_selec` implement this model using simulated data. An empirical application to door-t0-door fundraising is given in script `mod9_select_door2door`.

Treatment – Effect Model

The treatment-effect (TE) model has a similar structure as the Sample-Selection (SS) model. It focuses on the causal impact of a binary “treatment” on a continuous outcome. So to start, our main equation of interest can be written as

$$y_i = \mathbf{x}'_{2i}\boldsymbol{\beta}_2 + \gamma T_i + \varepsilon_{2i} \quad (36)$$

where $T_i = 1$ if individual i received the treatment, and 0 otherwise. If T_i was truly exogenous, we could stop right here and estimate the model as a generic normal regression. An exogenous treatment would be any outside intervention that is imposed on a given individual, rather than her choosing it. For example, a letter sent to randomly chosen households reminding them of designated watering days would be such a treatment. Or randomly choosing some wild horses to be trained before releasing them for adoption would be another.

However, in many cases the treatment is chosen by the individual herself. Examples are: Decision to enter a job training program, decision to start / quit smoking, decision to engage in criminal behavior, choice of schooling, etc. In these situations there could be the problem of *unobserved confounding*: If there exists an unobserved effect or characteristic that drives an individual towards a treatment decision, and that is correlated with unobservables in the outcome equation, our measure of the causal treatment effect (i.e. our estimate of γ) would be biased. A classic example would be that an individual with more “spunk” is more likely to enter a training program, but would have likely performed above average without the program as well (as measured by the dependent variable in the outcome equation). If mostly “talented” individuals take the training, and mostly “untalented” don’t, the causal impact of the treatment would be overestimated.

Econometrically, we can model this as

$$\begin{aligned}
T_i^* &= \mathbf{x}'_{1i} \boldsymbol{\beta}_1 + \varepsilon_{1i} \\
y_i &= \mathbf{x}'_{2i} \boldsymbol{\beta}_2 + \gamma T_i + \varepsilon_{2i} \\
\begin{bmatrix} \varepsilon_{1i} \\ \varepsilon_{2i} \end{bmatrix} &\sim n\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}\right) & \boldsymbol{\Sigma} &= \begin{bmatrix} 1 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix} \\
T_i &= 1 & \text{if } T_i^* > 0, & = 0 \text{ otherwise}
\end{aligned} \tag{37}$$

Note that we observe y_i regardless of that individual's treatment decision. Thus, in contrast to the SS model, data attrition via a "first hurdle" is not an issue in the TE model. We will also assume that all regressors are observed, regardless of the treatment decision. At the individual level we can write the model as

$$\begin{aligned}
\mathbf{y}_i^* &= \mathbf{X}_i \boldsymbol{\beta} + \boldsymbol{\varepsilon}_i \quad \text{with} \quad \boldsymbol{\varepsilon}_i \sim n(\mathbf{0}, \boldsymbol{\Sigma}), \\
\mathbf{y}_i^* &= \begin{bmatrix} T_i^* \\ y_i \end{bmatrix} \quad \mathbf{X}_i = \begin{bmatrix} \mathbf{x}'_{1i} & \mathbf{0} & 0 \\ \mathbf{0} & \mathbf{x}'_{2i} & T_i \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \\ \gamma \end{bmatrix} \quad \boldsymbol{\varepsilon}_i = \begin{bmatrix} \varepsilon_{1i} \\ \varepsilon_{2i} \end{bmatrix}
\end{aligned} \tag{38}$$

For the full sample of n individuals the model can be written as

$$\begin{aligned}
\mathbf{y}^* &= \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon} \\
\mathbf{y}^* &= \begin{bmatrix} \mathbf{y}_1^* \\ \mathbf{y}_2^* \\ \vdots \\ \mathbf{y}_n^* \end{bmatrix}_{(nx2) \times 1} \quad \mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_n \end{bmatrix}_{(nx2) \times k} \quad \boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \\ \gamma \end{bmatrix} \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \boldsymbol{\varepsilon}_1 \\ \boldsymbol{\varepsilon}_2 \\ \vdots \\ \boldsymbol{\varepsilon}_n \end{bmatrix} \sim n(\mathbf{0}, \boldsymbol{\Omega}) \quad \text{with} \\
\boldsymbol{\Omega} &= \begin{bmatrix} \boldsymbol{\Sigma} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma} & \dots & \mathbf{0} \\ \mathbf{0} & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \boldsymbol{\Sigma} \end{bmatrix}_{(nx2) \times (nx2)}
\end{aligned} \tag{39}$$

For our sample likelihood there are again *two possible observable outcomes*, $pr(T_i = 1, y_i)$ and $pr(T_i = 0, y_i)$. Building on our insights from the SS model, we can express these outcomes as

$$\begin{aligned}
pr(T_i = 1, y_i) &= f(\varepsilon_{1i} > -\mathbf{x}'_{1i} \boldsymbol{\beta}_1, \varepsilon_{2i}) = \\
f(\varepsilon_{2i}) \int_{-\mathbf{x}'_{1i} \boldsymbol{\beta}_1}^{\infty} f(\varepsilon_{1i} | \varepsilon_{2i}) d\varepsilon_{1i} &= \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{1}{2\sigma_2^2}(y_i - \mathbf{x}'_{2i} \boldsymbol{\beta}_2 - \gamma)^2\right) \Phi\left(\frac{\mathbf{x}'_{1i} \boldsymbol{\beta}_1 + \frac{\sigma_{12}}{\sigma_2^2}(y_i - \mathbf{x}'_{2i} \boldsymbol{\beta}_2 - \gamma)}{\sqrt{1 - \rho^2}}\right) \\
\rho &= \frac{\sigma_{12}}{\sqrt{\sigma_2^2}}
\end{aligned} \tag{40}$$

and

$$pr(T_i = 0, y_i) = f(\varepsilon_{1i} < -\mathbf{x}'_{1i}\boldsymbol{\beta}_1, \varepsilon_{2i}) =$$

$$f(\varepsilon_{2i}) \int_{-\infty}^{-\mathbf{x}'_{1i}\boldsymbol{\beta}_1} f(\varepsilon_{1i} | \varepsilon_{2i}) d\varepsilon_{1i} = \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{1}{2\sigma_2^2}(y_i - \mathbf{x}'_{2i}\boldsymbol{\beta}_2)^2\right) \Phi\left(\frac{-\mathbf{x}'_{1i}\boldsymbol{\beta}_1 - \frac{\sigma_{12}}{\sigma_2^2}(y_i - \mathbf{x}'_{2i}\boldsymbol{\beta}_2)}{\sqrt{1-\rho^2}}\right) \quad (41)$$

We can now write the sample likelihood concisely as

$$p(\mathbf{y} | \boldsymbol{\beta}, \boldsymbol{\Sigma}) = \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{1}{2\sigma_2^2}(y_{2i} - \mathbf{x}'_{2i}\boldsymbol{\beta}_2 - \gamma)^2\right) \Phi\left(\frac{\mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \frac{\sigma_{12}}{\sigma_2^2}(y_{2i} - \mathbf{x}'_{2i}\boldsymbol{\beta}_2 - \gamma)}{\sqrt{1-\rho^2}}\right) \right)^{T_i} *$$

$$\left(\frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{1}{2\sigma_2^2}(y_{2i} - \mathbf{x}'_{2i}\boldsymbol{\beta}_2)^2\right) \Phi\left(\frac{-\mathbf{x}'_{1i}\boldsymbol{\beta}_1 - \frac{\sigma_{12}}{\sigma_2^2}(y_{2i} - \mathbf{x}'_{2i}\boldsymbol{\beta}_2)}{\sqrt{1-\rho^2}}\right) \right)^{1-T_i} \quad (42)$$

Using the same priors as for the SS model the augmented posterior kernel takes the following form:

$$p(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{y}^* | \mathbf{y}, \mathbf{X}) \propto$$

$$|\boldsymbol{\Sigma}|^{-(\nu_0+3+n)/2} \exp\left(-\frac{1}{2}tr(\mathbf{S}_0 \cdot \boldsymbol{\Sigma}^{-1})\right) I(\Sigma_{11} = 1) * \exp\left(-\frac{1}{2}\left((\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0)\right)\right) *$$

$$\exp\left(-\frac{1}{2}\sum_{i=1}^n (\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})\right) * \quad (43)$$

$$\prod_{i=1}^n \left(I(T_i = 0, y_i) I(T_i^* < 0, y_i) + I(T_i = 1, y_i) I(T_i^* > 0, y_i) \right)$$

The conditional posteriors for $\boldsymbol{\beta}$ and $\boldsymbol{\Sigma}$ take the exact same form as for the SS model. The conditional posterior for the latent data is given as

$$p(\mathbf{y}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{y}, \mathbf{X}) \propto$$

$$\exp\left(-\frac{1}{2}\sum_{i=1}^n (\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})\right) * \quad (44)$$

$$\prod_{i=1}^n \left(I(T_i = 0, y_i) I(T_i^* < 0, y_i) + I(T_i = 1, y_i) I(T_i^* > 0, y_i) \right)$$

At the individual level this implies

$$\begin{aligned}
& p(\mathbf{y}_i^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, y_i, \mathbf{X}_i) \propto \\
& \exp\left(-\frac{1}{2}(\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})\right) * \\
& I(T_i = 0, y_i) I(T_i^* < 0, y_i) + I(T_i = 1, y_i) I(T_i^* > 0, y_i)
\end{aligned} \tag{45}$$

For the $T_i = 0$ case we have

$$\begin{aligned}
& p(\mathbf{y}_i^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, T_i = 0, \mathbf{X}_i) \propto \\
& \exp\left(-\frac{1}{2}(\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})\right) * I(T_i^* < 0, y_i)
\end{aligned} \tag{46}$$

This implies a draw of $\mathbf{y}_i^* = [T_i^* \quad y_i]'$ from a joint normal density, with T_i^* truncated to the negative domain, i.e. from

$$p(T_i^* < 0, y_i | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) = \int_{-\infty}^0 f(T_i^*, y_i | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) dT_i^* \tag{47}$$

Where $f(\cdot)$ denotes the bivariate normal density. This can be computationally implemented as follows:

Using $p(T_i^* < 0, y_i | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) = p(T_i^* < 0 | y_i, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) p(y_i | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) = p(T_i^* < 0 | y_i, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i)$, since y_i is observed. Thus, we simply draw T_i^* from its conditional normal density with mean

$$\mathbf{x}'_{1i} \boldsymbol{\beta}_1 + \frac{\sigma_{12}}{\sigma_2^2} (y_i - \mathbf{x}'_{2i} \boldsymbol{\beta}_2) \text{ and variance } 1 - \rho^2, \text{ truncated from above at zero.}$$

Similarly, for the $T_i = 1$ case we have

$$\begin{aligned}
& p(\mathbf{y}_i^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, T_i = 1, \mathbf{X}_i) \propto \\
& \exp\left(-\frac{1}{2}(\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})\right) * I(T_i^* > 0, y_i)
\end{aligned} \tag{48}$$

This implies a draw of $\mathbf{y}_i^* = [T_i^* \quad y_i]'$ from a joint normal density, with T_i^* truncated to the positive domain, i.e. from

$$p(T_i^* > 0, y_i | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) = \int_0^{\infty} f(T_i^*, y_i | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) dT_i^* \tag{49}$$

This can be computationally implemented by simply drawing T_i^* from its conditional normal density with mean $\mathbf{x}'_{1i} \boldsymbol{\beta}_1 + \frac{\sigma_{12}}{\sigma_2^2} (y_i - \mathbf{x}'_{2i} \boldsymbol{\beta}_2 - \gamma)$ and variance $1 - \rho^2$, truncated from below at zero.

Matlab implementation:

See Matlab scripts / functions `mod9_treat_data`, `mod9_treat`, and `gs_treat` for the implementation of this model using simulated data. An empirical example using a popular job training data set is given in scripts `mod9_treatJobTrain` and `mod9_treatJobTrainNaive`.

Switching Regression Model

This model is closely related to the previous two. It is also known as the “Roy Model” (after Roy, 1951) and “type 5 Tobit” (as listed in Amemiya, 1985). It is a three equation model, where the first equation determines an endogenous treatment or selection outcome, and the remaining two equations describe the main outcome of interest for the treated and untreated sample, respectively. Thus, we now allow the outcome equation to differ (in terms of explanatory variables and / or coefficients) depending on the treatment decision.

Econometrically, this model is interesting since it shows how in Bayesian analysis we can sometimes learn about a parameter that is not directly identified.

The model can be written as

$$\begin{aligned}
 y_{1i}^* &= \mathbf{x}'_{1i} \boldsymbol{\beta}_1 + \varepsilon_{1i} \\
 y_{2i}^* &= \mathbf{x}'_{2i} \boldsymbol{\beta}_2 + \varepsilon_{2i} \\
 y_{3i}^* &= \mathbf{x}'_{3i} \boldsymbol{\beta}_3 + \varepsilon_{3i}
 \end{aligned}$$

$$\begin{bmatrix} \varepsilon_{1i} \\ \varepsilon_{2i} \\ \varepsilon_{3i} \end{bmatrix} \sim n \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma} \right) \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & \sigma_{12} & \sigma_{13} \\ \sigma_{12} & \sigma_{22} & \sigma_{23} \\ \sigma_{13} & \sigma_{23} & \sigma_{33} \end{bmatrix} \tag{50}$$

$$\begin{aligned}
 y_{1i} &= 1 & \text{if } & y_{1i}^* > 0, & & = 0 \text{ otherwise} \\
 y_{2i} &= y_{2i}^* & \text{if } & y_{1i} = 0 & & = \text{unobserved otherwise} \\
 y_{3i} &= y_{3i}^* & \text{if } & y_{1i} = 1 & & = \text{unobserved otherwise}
 \end{aligned}$$

The new element is that different outcome equations apply to different individuals. As for the two preceding models, the first equation must contain at least one regressor that is not included in the other two equations for identification purposes. As before, we assume that explanatory variables are observed for all individuals and equations, regardless of the first hurdle result. Although not necessary, most existing applications set $\mathbf{x}_{2i} = \mathbf{x}_{3i}$.

Learning about σ_{23}

Since *either* equ. 2 *or* equ. 3 apply to a given individual (i.e. a given individual cannot be treated AND untreated at the same time), the covariance between equations 2 and 3 (σ_{23}) does not enter the sample

likelihood, i.e. it is not identified. Yet, we will be able to derive its posterior density, based on information through priors (weak, if priors are vague) and – due to the error correlation – information via remaining model parameters.

The details about this indirect learning are given in Lie et al (2004). I will synthesize their argument:

First note that for Σ to be a well-behaved, positive definite variance matrix, we need

$$\begin{aligned}
 |\Sigma| &= \sigma_{22}\sigma_{33} \left(\left(1 - \frac{\sigma_{12}^2}{\sigma_{22}}\right) \left(1 - \frac{\sigma_{13}^2}{\sigma_{33}}\right) - \left(\frac{\sigma_{23}}{\sqrt{\sigma_{22}\sigma_{33}}} - \frac{\sigma_{12}\sigma_{13}}{\sqrt{\sigma_{22}\sigma_{33}}} \right)^2 \right) > 0 && \rightarrow \\
 \left(1 - \frac{\sigma_{12}^2}{\sigma_{22}}\right) \left(1 - \frac{\sigma_{13}^2}{\sigma_{33}}\right) &> \left(\frac{\sigma_{23}}{\sqrt{\sigma_{22}\sigma_{33}}} - \frac{\sigma_{12}\sigma_{13}}{\sqrt{\sigma_{22}\sigma_{33}}} \right)^2 && \rightarrow \quad (51) \\
 \sigma_{12}\sigma_{13} - \sqrt{\sigma_{22}\sigma_{33}} \left(\left(1 - \frac{\sigma_{12}^2}{\sigma_{22}}\right) \left(1 - \frac{\sigma_{13}^2}{\sigma_{33}}\right) \right)^{1/2} &< \sigma_{23} < \sigma_{12}\sigma_{13} + \sqrt{\sigma_{22}\sigma_{33}} \left(\left(1 - \frac{\sigma_{12}^2}{\sigma_{22}}\right) \left(1 - \frac{\sigma_{13}^2}{\sigma_{33}}\right) \right)^{1/2}
 \end{aligned}$$

Thus, conditional on the other elements of Σ the unidentified covariance is bound between

$$\begin{aligned}
 \underline{\sigma}_{23} &= \sigma_{12}\sigma_{13} - \sqrt{\sigma_{22}\sigma_{33}} \left(\left(1 - \frac{\sigma_{12}^2}{\sigma_{22}}\right) \left(1 - \frac{\sigma_{13}^2}{\sigma_{33}}\right) \right)^{1/2} && \text{and} \\
 \bar{\sigma}_{23} &= \sigma_{12}\sigma_{13} + \sqrt{\sigma_{22}\sigma_{33}} \left(\left(1 - \frac{\sigma_{12}^2}{\sigma_{22}}\right) \left(1 - \frac{\sigma_{13}^2}{\sigma_{33}}\right) \right)^{1/2} && (52)
 \end{aligned}$$

As described in Li et al. the other model parameters and the data add no further contribution to learning about σ_{23} once we condition its posterior on the other elements of Σ , i.e.

$$p(\sigma_{23} | \Sigma_{-23}, \boldsymbol{\theta}, \mathbf{y}, \mathbf{X}) = p(\sigma_{23} | \Sigma_{-23}) \quad (53)$$

where Σ_{-23} captures all other terms in Σ and $\boldsymbol{\theta}$ collects all other model parameters. Now define $R(\sigma_{23})$ the joint parameter space of Σ_{-23} , given σ_{23} , that assures that Σ remains positive definite. We can then write the marginal posterior of σ_{23} as:

$$\begin{aligned}
p(\sigma_{23} | \mathbf{y}, \mathbf{X}) &= \int_{R(\Sigma_{-23})} \int_{\boldsymbol{\theta}} p(\sigma_{23}, \Sigma_{-23}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) d\boldsymbol{\theta} d\Sigma_{-23} = \\
&\int_{R(\Sigma_{-23})} \int_{\boldsymbol{\theta}} p(\sigma_{23} | \Sigma_{-23}, \boldsymbol{\theta}, \mathbf{y}, \mathbf{X}) p(\Sigma_{-23}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) d\boldsymbol{\theta} d\Sigma_{-23} = \\
&\int_{R(\Sigma_{-23})} \int_{\boldsymbol{\theta}} p(\sigma_{23} | \Sigma_{-23}) p(\Sigma_{-23}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) d\boldsymbol{\theta} d\Sigma_{-23} = \\
&\int_{R(\Sigma_{-23})} p(\sigma_{23} | \Sigma_{-23}) p(\Sigma_{-23} | \mathbf{y}, \mathbf{X}) \left(\int_{\boldsymbol{\theta}} p(\boldsymbol{\theta} | \Sigma_{-23}, \mathbf{y}, \mathbf{X}) d\boldsymbol{\theta} \right) d\Sigma_{-23} = \\
&\int_{R(\Sigma_{-23})} p(\sigma_{23} | \Sigma_{-23}) p(\Sigma_{-23} | \mathbf{y}, \mathbf{X}) d\Sigma_{-23}
\end{aligned} \tag{54}$$

where we use result (53) in the third line. The final result shows that the marginal posterior of σ_{23} can be expressed as a function solely of the remaining elements in Σ . Furthermore, if the joint posterior of Σ_{-23} is highly informative or ‘tight ‘ around some point $\hat{\Sigma}_{-23}$, then

$$\begin{aligned}
p(\sigma_{23} | \mathbf{y}, \mathbf{X}) &\approx p(\sigma_{23} | \hat{\Sigma}_{-23}) \quad \text{or} \quad pr(\hat{\sigma}_{23} \leq \sigma_{23} \leq \bar{\sigma}_{23} | \mathbf{y}, \mathbf{X}) \approx 1 \quad \text{where} \\
\hat{\sigma}_{23} &= \sigma_{23} | \Sigma_{-23} = \hat{\Sigma}_{-23}, \quad \bar{\sigma}_{23} = \bar{\sigma}_{23} | \Sigma_{-23} = \hat{\Sigma}_{-23}
\end{aligned} \tag{55}$$

are posterior estimates of the unidentified covariances given in (52).

Thus, as stated in Li et al, information conveyed from the data on the identified elements of Σ ‘spills over’ and revises our belief about the conditional support of σ_{23} . In other words, the marginal prior and posterior for σ_{23} will generally differ, suggesting that learning has taken place. This can be easily verified by plotting both prior and posterior after running the Gibbs Sampler, as we will show below. Poirier (1998) provides other examples of learning about non-identified parameters.

Estimation

Before we work out the details for our Gibbs Sampler we need to review the formulas for bi-conditional normal densities, i.e. normal densities for a random variable, conditioned on *two* other normal variates. These densities will be needed to draw latent data in the Gibbs Sampler: Assume random vector \mathbf{x} is multinormal with mean $\boldsymbol{\mu}$ and variance matrix Σ . Partition \mathbf{x} arbitrarily into sub-vectors \mathbf{x}_1 and \mathbf{x}_2 , and

Σ conformably into $\begin{bmatrix} \Sigma_{11} & \Sigma'_{12} \\ \Sigma_{12} & \Sigma_{22} \end{bmatrix}$. As shown e.g. in Greene, Theorem B.7, p. 1014, the conditional densities for each sub-vector are then given as:

$$\begin{aligned}
\mathbf{x}_1 | \mathbf{x}_2 &\sim n(\boldsymbol{\mu}_{1,2}, \boldsymbol{\Sigma}_{1,2}) & \mathbf{x}_2 | \mathbf{x}_1 &\sim n(\boldsymbol{\mu}_{2,1}, \boldsymbol{\Sigma}_{2,1}) \\
\boldsymbol{\mu}_{1,2} &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}'_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2) & \boldsymbol{\mu}_{2,1} &= \boldsymbol{\mu}_2 + \boldsymbol{\Sigma}'_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) \\
\boldsymbol{\Sigma}_{1,2} &= \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}'_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12} & \boldsymbol{\Sigma}_{2,1} &= \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}'_{12}
\end{aligned} \tag{56}$$

For our GS we will require the conditional densities for $y_2^* | \begin{bmatrix} y_1^* \\ y_3^* \end{bmatrix} = y_2^* | \mathbf{y}_{13}^*$, and $y_3^* | \begin{bmatrix} y_1^* \\ y_2^* \end{bmatrix} = y_3^* | \mathbf{y}_{12}^*$. For

completeness, I also added the density for $y_1^* | \begin{bmatrix} y_2^* \\ y_3^* \end{bmatrix} = y_1^* | \mathbf{y}_{23}^*$.

Let

$$\begin{aligned}
E(y_1^*) &= \mathbf{x}'_1 \boldsymbol{\beta}_1 = \omega_1, & E(y_2^*) &= \mathbf{x}'_2 \boldsymbol{\beta}_2 = \omega_2, & E(y_3^*) &= \mathbf{x}'_3 \boldsymbol{\beta}_3 = \omega_3, \\
E(\mathbf{y}_{12}^*) &= \begin{bmatrix} \omega_1 \\ \omega_2 \end{bmatrix} = \boldsymbol{\omega}_{12}, & E(\mathbf{y}_{13}^*) &= \begin{bmatrix} \omega_1 \\ \omega_3 \end{bmatrix} = \boldsymbol{\omega}_{13}, & E(\mathbf{y}_{23}^*) &= \begin{bmatrix} \omega_2 \\ \omega_3 \end{bmatrix} = \boldsymbol{\omega}_{23}, \\
\boldsymbol{\Sigma}_{1r} &= \begin{bmatrix} \sigma_{12} \\ \sigma_{13} \end{bmatrix}, & \boldsymbol{\Sigma}_{2r} &= \begin{bmatrix} \sigma_{12} \\ \sigma_{23} \end{bmatrix}, & \boldsymbol{\Sigma}_{3r} &= \begin{bmatrix} \sigma_{13} \\ \sigma_{23} \end{bmatrix} \\
\boldsymbol{\Sigma}_{12} &= \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix} & \boldsymbol{\Sigma}_{13} &= \begin{bmatrix} \sigma_{11} & \sigma_{13} \\ \sigma_{13} & \sigma_{33} \end{bmatrix} & \boldsymbol{\Sigma}_{23} &= \begin{bmatrix} \sigma_{22} & \sigma_{23} \\ \sigma_{23} & \sigma_{33} \end{bmatrix}
\end{aligned} \tag{57}$$

We can then derive the required conditional densities as

$$\begin{aligned}
y_1^* | \mathbf{y}_{23}^* &\sim n(\omega_{1,23}, \sigma_{1,23}) & \text{with} \\
\omega_{1,23} &= \omega_1 + \boldsymbol{\Sigma}'_{1r} \boldsymbol{\Sigma}_{23}^{-1} (\mathbf{y}_{23}^* - \boldsymbol{\omega}_{23}) & \sigma_{1,23} &= 1 - \boldsymbol{\Sigma}'_{1r} \boldsymbol{\Sigma}_{23}^{-1} \boldsymbol{\Sigma}_{1r} \\
y_2^* | \mathbf{y}_{13}^* &\sim n(\omega_{2,13}, \sigma_{2,13}) & \text{with} \\
\omega_{2,13} &= \omega_2 + \boldsymbol{\Sigma}'_{2r} \boldsymbol{\Sigma}_{13}^{-1} (\mathbf{y}_{13}^* - \boldsymbol{\omega}_{13}) & \sigma_{2,13} &= \sigma_{22} - \boldsymbol{\Sigma}'_{2r} \boldsymbol{\Sigma}_{13}^{-1} \boldsymbol{\Sigma}_{2r} \\
y_3^* | \mathbf{y}_{12}^* &\sim n(\omega_{3,12}, \sigma_{3,12}) & \text{with} \\
\omega_{3,12} &= \omega_3 + \boldsymbol{\Sigma}'_{3r} \boldsymbol{\Sigma}_{12}^{-1} (\mathbf{y}_{12}^* - \boldsymbol{\omega}_{12}) & \sigma_{3,12} &= \sigma_{33} - \boldsymbol{\Sigma}'_{3r} \boldsymbol{\Sigma}_{12}^{-1} \boldsymbol{\Sigma}_{3r}
\end{aligned} \tag{58}$$

For a given individual, we can write the system as for the two preceding models, plus an additional equation:

$$\begin{aligned}
\mathbf{y}_i^* &= \mathbf{X}_i \boldsymbol{\beta} + \boldsymbol{\varepsilon}_i & \text{with} & \boldsymbol{\varepsilon}_i \sim n(\mathbf{0}, \boldsymbol{\Sigma}), \\
\mathbf{y}_i^* &= \begin{bmatrix} y_{1i}^* \\ y_{2i}^* \\ y_{3i}^* \end{bmatrix} & \mathbf{X}_i &= \begin{bmatrix} \mathbf{x}'_{1i} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{x}'_{2i} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{x}'_{3i} \end{bmatrix} & \boldsymbol{\beta} &= \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \\ \boldsymbol{\beta}_3 \end{bmatrix} & \boldsymbol{\varepsilon}_i &= \begin{bmatrix} \varepsilon_{1i} \\ \varepsilon_{2i} \\ \varepsilon_{3i} \end{bmatrix}
\end{aligned} \tag{59}$$

Also, building on the two previous model, it is straightforward to derive the likelihood function as

$$\begin{aligned}
p(\mathbf{y} | \boldsymbol{\beta}, \boldsymbol{\Sigma}) &= \prod_{i=1}^n (f(\varepsilon_{1i} < -\mathbf{x}'_{1i}\boldsymbol{\beta}_1, \varepsilon_{2i}))^{1-y_{1i}} (f(\varepsilon_{1i} > -\mathbf{x}'_{1i}\boldsymbol{\beta}_1, \varepsilon_{3i}))^{y_{1i}} = \\
&\prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi\sigma_{22}}} \exp\left(-\frac{1}{2\sigma_{22}}(y_{2i} - \mathbf{x}'_{2i}\boldsymbol{\beta}_2)^2\right) \Phi\left(\frac{-\mathbf{x}'_{1i}\boldsymbol{\beta}_1 - \frac{\sigma_{12}}{\sigma_{22}}(y_{2i} - \mathbf{x}'_{2i}\boldsymbol{\beta}_2)}{\sqrt{1-\rho_{12}^2}}\right) \right)^{1-y_{1i}} * \\
&\left(\frac{1}{\sqrt{2\pi\sigma_{33}}} \exp\left(-\frac{1}{2\sigma_{33}}(y_{3i} - \mathbf{x}'_{3i}\boldsymbol{\beta}_3)^2\right) \Phi\left(\frac{\mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \frac{\sigma_{12}}{\sigma_{33}}(y_{3i} - \mathbf{x}'_{3i}\boldsymbol{\beta}_3)}{\sqrt{1-\rho_{13}^2}}\right) \right)^{y_{1i}} \text{ with} \\
\rho_{12} &= \frac{\sigma_{12}}{\sqrt{\sigma_{22}}} \quad \rho_{13} = \frac{\sigma_{13}}{\sqrt{\sigma_{33}}}
\end{aligned} \tag{60}$$

We will augment the model with latent data for all three equations. The unobserved latent data for equations two and three is also often referred to as “missing data”.

Using the same priors as for the SS and TE models the augmented posterior kernel takes the following form:

$$\begin{aligned}
p(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{y}^* | \mathbf{y}, \mathbf{X}) &\propto \\
|\boldsymbol{\Sigma}|^{-(v_0+4+n)/2} \exp\left(-\frac{1}{2}\text{tr}(\mathbf{S}_0 \cdot \boldsymbol{\Sigma}^{-1})\right) &I(\Sigma_{11} = 1) * \exp\left(-\frac{1}{2}\left((\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0)\right)\right) * \\
\exp\left(-\frac{1}{2}\sum_{i=1}^n (\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})\right) &* \\
\prod_{i=1}^n \left(I(y_{1i} = 0, y_{2i} = y_{2i}^*) I(y_{1i}^* < 0, y_{2i}^*) + I(y_{1i} = 1, y_{3i} = y_{3i}^*) I(y_{1i}^* > 0, y_{3i}^*) \right) &
\end{aligned} \tag{61}$$

The conditional posteriors for $\boldsymbol{\beta}$ and $\boldsymbol{\Sigma}$ take the exact same form as for the SS model (except for the “4” instead of “3” in $|\boldsymbol{\Sigma}|^{-(v_0+4+n)/2}$). The conditional posterior for the latent data is given as

$$\begin{aligned}
p(\mathbf{y}^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{y}, \mathbf{X}) &\propto \\
\exp\left(-\frac{1}{2}\sum_{i=1}^n (\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i\boldsymbol{\beta})\right) &* \\
\prod_{i=1}^n \left(I(y_{1i} = 0, y_{2i} = y_{2i}^*) I(y_{1i}^* < 0, y_{2i}^*) + I(y_{1i} = 1, y_{3i} = y_{3i}^*) I(y_{1i}^* > 0, y_{3i}^*) \right) &
\end{aligned} \tag{62}$$

At the individual level this implies

$$\begin{aligned}
& p(\mathbf{y}_i^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{y}, \mathbf{X}_i) \propto \\
& \exp\left(-\frac{1}{2}(\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})\right) * \\
& I(y_{1i} = 0, y_{2i} = y_{2i}^*) I(y_{1i}^* < 0, y_{2i}^*) + I(y_{1i} = 1, y_{3i} = y_{3i}^*) I(y_{1i}^* > 0, y_{3i}^*)
\end{aligned} \tag{63}$$

For the $y_{1i} = 0$ case we have

$$\begin{aligned}
& p(\mathbf{y}_i^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, y_{1i} = 0, \mathbf{X}_i) \propto \\
& \exp\left(-\frac{1}{2}(\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})\right) * I(y_{1i}^* < 0, y_{2i}^*)
\end{aligned} \tag{64}$$

Thus, we need draws of $\mathbf{y}_i^* = [y_{1i}^* \quad y_{2i}^* \quad y_{3i}^*]'$ conditional on: (i) $y_{1i}^* < 0$ and (ii) y_{2i}^* is observed.

In other words, we need to draw $[y_{1i}^* \quad y_{3i}^*]'$ | y_{2i}^* with y_{1i}^* truncated to the negative domain, i.e. from

$$p(y_{1i}^* < 0, y_{3i}^* | y_{2i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) = \int_{-\infty}^0 f(y_{1i}^*, y_{3i}^* | y_{2i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) dy_{1i}^* \tag{65}$$

Where $f(\cdot)$ denotes the bivariate normal density. This can be computationally implemented as follows:

1. Using $p(y_{1i}^* < 0, y_{3i}^* | y_{2i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) = p(y_{3i}^* | y_{1i}^* < 0, y_{2i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i) p(y_{1i}^* < 0 | y_{2i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i)$, we can first draw y_{1i}^* from its conditioned-on y_{2i}^* normal density with mean $\mathbf{x}'_{1i} \boldsymbol{\beta}_1 + \frac{\sigma_{12}}{\sigma_{22}} (\mathbf{y}_{2i} - \mathbf{x}'_{2i} \boldsymbol{\beta}_2)$ and variance $1 - \rho_{12}^2$, truncated from above at zero.
2. We then draw $y_{3i}^* | \mathbf{y}_{12}^*$ from its bi-conditional normal density, with moments given in (58).

For the $y_{1i} = 1$ case we have

$$\begin{aligned}
& p(\mathbf{y}_i^* | \boldsymbol{\beta}, \boldsymbol{\Sigma}, y_{1i} = 1, \mathbf{X}_i) \propto \\
& \exp\left(-\frac{1}{2}(\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i^* - \mathbf{X}_i \boldsymbol{\beta})\right) * I(y_{1i}^* > 0, y_{3i}^*)
\end{aligned} \tag{66}$$

Thus, we need draws of $\mathbf{y}_i^* = [y_{1i}^* \quad y_{2i}^* \quad y_{3i}^*]'$ conditional on: (i) $y_{1i}^* > 0$ and (ii) y_{3i}^* is observed.

In other words, we need to draw $[y_{1i}^* \quad y_{2i}^*]'$ | y_{3i}^* with y_{1i}^* truncated to the positive domain, i.e. from

$$p\left(y_{1i}^* > 0, y_{2i}^* \mid y_{3i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i\right) = \int_0^{\infty} f\left(y_{1i}^*, y_{2i}^* \mid y_{3i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i\right) dy_{1i}^* \quad (67)$$

Where $f(\cdot)$ denotes the bivariate normal density. This can be computationally implemented as follows:

1. Using $p\left(y_{1i}^* > 0, y_{2i}^* \mid y_{3i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i\right) = p\left(y_{2i}^* \mid y_{1i}^* > 0, y_{3i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i\right) p\left(y_{1i}^* > 0 \mid y_{3i}^*, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X}_i\right)$, we can first draw y_{1i}^* from its conditioned-on y_{3i}^* normal density with mean $\mathbf{x}'_{1i}\boldsymbol{\beta}_1 - \frac{\sigma_{13}}{\sigma_{33}}(y_{3i}^* - \mathbf{x}'_{3i}\boldsymbol{\beta}_3)$ and variance $1 - \rho_{13}^2$, truncated from below at zero.
2. We then draw $y_{2i}^* \mid y_{1i}^*$ from its bi-conditional normal density, with moments given in (58).

Matlab examples: `mod9_switch_data`, `mod9_switch`

Post-Estimation

Two of the most important outcomes of interest in the program evaluation literature are the Average Treatment Effect (ATE) and the Treatment Effect on the Treated (TT). The first refers to the expected treatment effect for a randomly drawn person from the underlying population of interest. The second refers to the expected treatment effect on those who actually receive the treatment.

Let \mathbf{x}_p denote some common predictive settings for regressors in equations 2 and 3, and let

$\Delta_p = y_{3p}^* - y_{2p}^* = \mathbf{x}'_p(\boldsymbol{\beta}_3 - \boldsymbol{\beta}_2) + (\varepsilon_3 - \varepsilon_2)$. Then in a classical setting the Average Treatment Effect would be computed as

$$ATE_p = E(\Delta_p) = E\left(\mathbf{x}'_p(\boldsymbol{\beta}_3 - \boldsymbol{\beta}_2) + (\varepsilon_3 - \varepsilon_2)\right) = \mathbf{x}'_p(\boldsymbol{\beta}_3 - \boldsymbol{\beta}_2) \quad (68)$$

To compute the (classical) TT, we need to introduce additional notation. Assume that regressors for the treatment equation are set to \mathbf{z}_p , s.t. $y_{1p}^* = \mathbf{z}'_p\boldsymbol{\beta}_1 + \varepsilon_1$. Also let $\Delta_\varepsilon = \varepsilon_3 - \varepsilon_2$ with

$$E(\Delta_\varepsilon) = 0, V(\Delta_\varepsilon) = \sigma_{2-3}^2 = \sigma_{33} + \sigma_{22} - 2\sigma_{23}. \text{ Also note that}$$

$$\text{cov}(\varepsilon_1, \Delta_\varepsilon) = \sigma_{1,3-2} = \text{cov}(\varepsilon_1, (\varepsilon_3 - \varepsilon_2)) = \text{cov}(\varepsilon_1, \varepsilon_3) - \text{cov}(\varepsilon_1, \varepsilon_2) = \sigma_{13} - \sigma_{12}.$$

Then

$$TT_p = E(\Delta_p \mid y_{1p}^* > 0) = \mathbf{x}'_p(\boldsymbol{\beta}_3 - \boldsymbol{\beta}_2) + E(\Delta_\varepsilon \mid \varepsilon_1 < \mathbf{z}'_p\boldsymbol{\beta}_1) =$$

$$\mathbf{x}'_p(\boldsymbol{\beta}_3 - \boldsymbol{\beta}_2) + \sigma_{1,2-3} \frac{\phi(\mathbf{z}'_p\boldsymbol{\beta}_1)}{\Phi(\mathbf{z}'_p\boldsymbol{\beta}_1)} \quad (69)$$

Note that, not surprisingly, neither the expression for ATE nor TT include σ_{23} . The Bayesian approach has the advantage that learning about σ_{23} is possible during posterior analysis. This allows for the construction of *PPDs* for ATE and TT, which are much more informative than the simple expectations given in (68) and (69).

The PPD for ATE is formally given by

$$p(\Delta_p | \mathbf{x}_p, \mathbf{X}, \mathbf{y}) = \int_{\Gamma} p(\Delta_p | \mathbf{x}_p, \Gamma) p(\Gamma | \mathbf{y}, \mathbf{X}) d\Gamma \quad \text{with} \quad (70)$$

$$\Delta_p | \mathbf{x}_p, \Gamma \sim n(\mathbf{x}'_p (\boldsymbol{\beta}_3 - \boldsymbol{\beta}_2), \sigma_{22} + \sigma_{33} - 2\sigma_{23})$$

where Γ comprises all model parameters. The posterior mean of $\Delta_p | \mathbf{x}_p, \mathbf{X}, \mathbf{y}$ will be the analog of the classical ATE measure. However, to go beyond the simple expectation of this distribution and to simulate the entire distribution, knowledge of σ_{23} is required, as is evident from the formula.

Similarly, for the PPD of TT we have

$$p(\Delta_p | \mathbf{x}_p, \mathbf{z}_p, y_{1p}^* > 0, \mathbf{X}, \mathbf{y}) = \int_{\Gamma} p(\Delta_p | \mathbf{x}_p, \mathbf{z}_p, y_{1p}^* > 0, \Gamma) p(\Gamma | \mathbf{y}, \mathbf{X}) d\Gamma \quad \text{with}$$

$$p(\Delta_p | \mathbf{x}_p, \mathbf{z}_p, y_{1p}^* > 0, \Gamma) = \frac{p(\Delta_p, y_{1p}^* > 0 | \mathbf{x}_p, \mathbf{z}_p, \Gamma)}{p(y_{1p}^* > 0 | \mathbf{x}_p, \mathbf{z}_p, \Gamma)} =$$

$$\Phi(\mathbf{z}'_p \boldsymbol{\beta}_1)^{-1} \int_{z'_p \boldsymbol{\beta}_1}^{\infty} p(\Delta_p, \varepsilon_1 | \mathbf{x}_p, \mathbf{z}_p, \Gamma) d\varepsilon_1 =$$

$$\Phi(\mathbf{z}'_p \boldsymbol{\beta}_1)^{-1} \int_{z'_p \boldsymbol{\beta}_1}^{\infty} f(\varepsilon_1 | \Delta_p, \mathbf{x}_p, \mathbf{z}_p, \Gamma) p(\Delta_p | \mathbf{x}_p, \mathbf{z}_p, \Gamma) d\varepsilon_1 =$$

$$\frac{p(\Delta_p | \mathbf{x}_p, \mathbf{z}_p, \Gamma)}{\Phi(\mathbf{z}'_p \boldsymbol{\beta}_1)} \int_{z'_p \boldsymbol{\beta}_1}^{\infty} f(\varepsilon_1 | \Delta_p, \mathbf{x}_p, \mathbf{z}_p, \Gamma) d\varepsilon_1 =$$

$$\frac{p(\Delta_p | \mathbf{x}_p, \mathbf{z}_p, \Gamma)}{\Phi(\mathbf{z}'_p \boldsymbol{\beta}_1)} \Phi\left(\frac{\mathbf{z}'_p \boldsymbol{\beta}_1 - \frac{\sigma_{1,3-2}}{\sigma_{3-2}^2} (\Delta_p - \mathbf{x}'_p (\boldsymbol{\beta}_3 - \boldsymbol{\beta}_2))}{\sqrt{1 - \frac{\sigma_{1,3-2}^2}{\sigma_{3-2}^2}}}\right) \quad \text{with}$$

$$p(\Delta_p | \mathbf{x}_p, \mathbf{z}_p, \Gamma) = n(\mathbf{x}'_p (\boldsymbol{\beta}_3 - \boldsymbol{\beta}_2), \sigma_{22} + \sigma_{33} - 2\sigma_{23}) \quad (71)$$

To draw $p(\Delta_p | \mathbf{x}_p, \mathbf{z}_p, y_{1p}^* > 0, \Gamma)$, we first draw y_{1p}^* from its truncated marginal density, then we draw Δ_p from its conditional normal with mean $\mathbf{x}'_p (\boldsymbol{\beta}_3 - \boldsymbol{\beta}_2) + \sigma_{1,3-2} (y_{1p}^* - \mathbf{z}'_p \boldsymbol{\beta}_1)$ and variance $\sigma_{3-2}^2 - \sigma_{1,3-2}^2$

References

Li, M., D. J. Poirier and J. L. Tobias. 2004. Do dropouts suffer from dropping out? Estimation and prediction of outcome gains in a generalized selection model. *Journal of Applied Econometrics* **19**: 203-225.

Poirier, D. J. 1998. Revising beliefs in non-identified models. *Econometric Theory* **14**: 483-509.