

① a) we have ε_{it} iid $N(0,1)$ conditional on X

$$\begin{cases} y_{it}^* = x_{it}'\beta + \varepsilon_{it} \\ y_{it} = \mathbb{1}_{y_{it}^* > 0} \end{cases}$$

$$\begin{aligned} \text{then } P(y_{it} = 1 | X) &= P(y_{it}^* > 0 | X) \\ &= P(y_{it}^* > 0 | x_{it}) \\ &= P(\varepsilon_{it} > -x_{it}'\beta | x_{it}) \\ &= \Phi(x_{it}'\beta) \end{aligned}$$

Thus this model can be fitted by MLE:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmax}} \sum_{i=1}^n \sum_{t=1}^{T_i} \left[y_{it} \ln \Phi(x_{it}'\beta) + (1-y_{it}) \ln (1-\Phi(x_{it}'\beta)) \right]$$

There are three main assumptions on ε_{it} :

① ε_{it} is independent over i, t . This means that the omitted variables should not display any time correlation, even for the same firm i . In particular, there is no time invariant unobserved heterogeneity. This is unlikely to be true as paying dividends may depend on unobservables of the firm (eg management practice) that are difficult to observe.

② We assume that ε_{it} are homoskedastic ie $V(\varepsilon_{it} | X) = 1$ for all i, t . this may not be reasonable if there is a lot of heterogeneity across firms.

③ We assume a specific pdf $N(\cdot, \cdot)$ for ε_{it} .

Violations of ①, ② or ③ leads to inconsistent estimators if we use the probit model. (pooled).

① b) xtprob in Stata

$$\text{Assume } \begin{cases} y_{it}^* = x_{it}'\beta + \varepsilon_{it} \\ y_{it} = \mathbb{1}_{y_{it}^* > 0} \\ \varepsilon_{it} = \alpha_i + \nu_{it} \end{cases}$$

We have a Balanced panel N observations of individuals for T periods.

The general formula is:

$$P(y_i | X_i) = P(y_{i1} \dots y_{iT} | X_i) = \int \dots \int_{a_{i1}}^{b_{i1}} \dots \int_{a_{iT}}^{b_{iT}} f(\varepsilon_{i1}, \dots, \varepsilon_{iT}) d\varepsilon$$

where $\left\{ \begin{array}{l} \text{if } y_{ij} = 0, a_{ij} = -\infty, b_{ij} = -x_{ij}'\beta \\ \text{if } y_{ij} = 1, a_{ij} = -x_{ij}'\beta, b_{ij} = +\infty \end{array} \right.$

This formula simplifies here if we assume that the observations (ν_{it}) are iid over t conditional on X_i, α_i . In xtprob, we assume $\nu_{it} \sim \mathcal{N}(0, 1) | X_i, \alpha_i$

then $P(y_i | X_i) = E(P(y_i | X_i, \alpha_i)) = \int \prod_{t=1}^T P(y_{it} | X_i, \alpha_i) f(\alpha_i) d\alpha_i$

We further assume that $\alpha_i | X_i \sim \mathcal{N}(0, \sigma_\alpha^2)$ iid over i .

Thus $P(y_i | X_i) = \int_{-\infty}^{+\infty} \prod_{t=1}^T P(y_{it} | X_i, \alpha_i) \frac{1}{\sigma_\alpha} \varphi\left(\frac{\alpha}{\sigma_\alpha}\right) d\alpha$

with $P(y_{it} | X_i, \alpha_i) = P(y_{it} | X_{it}, \alpha_i) = \Phi(x_{it}'\beta + \alpha_i)^{y_{it}} \times (1 - \Phi(x_{it}'\beta + \alpha_i))^{1-y_{it}}$

The observations are independent over i (individuals) so we get: (conditional on X):

$$L = \prod_{i=1}^N \left(\int_{-\infty}^{+\infty} \prod_{t=1}^T \Phi(x_{it}'\beta + \alpha) ^{y_{it}} (1 - \Phi(x_{it}'\beta + \alpha)) ^{1-y_{it}} \frac{1}{\sigma_\alpha} \varphi\left(\frac{\alpha}{\sigma_\alpha}\right) d\alpha \right)$$

① c) Suppose now that we have:

$$\begin{cases} y_{it}^* = x_{it}'\beta + \varepsilon_{it} \\ y_{it} = \mathbb{1}_{\{y_{it}^* > 0\}} \\ \varepsilon_{it} = \alpha_{it} + \nu_{it} \\ \nu_{it} = \rho \nu_{it-1} + \zeta_{it} \end{cases} \quad \zeta_{it} \text{ i.i.d. over } i, t \sim N(0,1)$$

Then as before we can write:

$$P(y_i | X_i) = P(y_{i1} \dots y_{iT} | X_i) = \int \prod_{j=1}^T \int_{a_j}^{b_j} f(\varepsilon_{i1}, \dots, \varepsilon_{iT}) d\varepsilon.$$

where $\begin{cases} \text{if } y_{ij} = 1, & a_j = -x_{ij}'\beta, & b_j = +\infty \\ \text{if } y_{ij} = 0, & a_j = -\infty, & b_j = -x_{ij}'\beta \end{cases}$

Here this T dimensional integral does not simplify as the value of ε_{it} depends both of the fixed heterogeneity α_i and the past value of ν_{it}, ν_{it-1} .

We need to use a simulation based procedure.

* The parameters of interest are: $\theta = \begin{cases} \beta, & \alpha \text{ the set of } \alpha_i, \\ \rho & \text{the autocorrelation factor of } \nu_{it}. \end{cases}$

* We could use SSML and define:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \frac{1}{N} \sum_{i=1}^N \ln \tilde{L}_i(\theta, R) = \underset{\theta}{\operatorname{argmax}} \tilde{L}(\theta)$$

where R is the number of simulations.

* We will need to iterate over our choices for θ , eg $\theta^{(m)}$
 at $\theta^{(m)}$ we need to evaluate $\tilde{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \ln \tilde{L}_i(\theta, R)$

* ① We draw a T dimensional uniform random vectors R times $\Rightarrow \tilde{u}_{i1}^{T \times 1}, \dots, \tilde{u}_{iR}^{T \times 1}$

* ② Given that we know $\theta^{(m)}$ we define:

$$\tilde{\varepsilon}_i^r(\theta^{(m)}) = F_{\varepsilon}^{-1}(\tilde{u}_i^R, \theta^{(m)})$$

thanks to our assumptions that the pdf of $\varepsilon_i, F_{\varepsilon}$.

x ③ Knowledge of $\tilde{\epsilon}_i^r(\theta^{(m)})$ and suit gives us a simulated latent var $\tilde{y}_{it}^*(\theta^{(m)})$

②

$$\text{and } \tilde{y}_{it}(\theta^{(m)}) = \mathbb{1}(\tilde{y}_{it}^* > 0)$$

x ④ From our R simulations we can get the simulated likelihood $\tilde{P}_i(\theta^{(m)}, R)$ of a given vector of choices $y_i = (y_{i1} \dots y_{iT})$ and thus compute $\tilde{L}(\theta^{(m)})$

x ⑤ Keeping fixed our simulations (ie keeping the same R $T \times 1$ vectors $\tilde{u}_{i,r}$). We can iterate ②, ③ and ④ to get the values of $\tilde{L}(\theta^{(m)})$ and search for θ satisfying this objective criterion-

① d) unbalanced panel data, now we have T_i observations by individual i and T_i is not constant.

in a) this does not matter at all as we assumed $\tilde{\epsilon}_{it}$ iid over i and t and on X . (if unbalancing exog).

in b) the likelihood of a given individual i becomes $P(y_i | X_i) = \int \prod_{t=1}^{T_i} P(y_{it} | X_i, \alpha_i) f(\alpha) d\alpha$

if the process that determines the number of observations is exogenous

inc) same as in b we will need to adjust for ϵ_i of different sizes if the observations are not all exogenous $T_i \times 1$

Here we have a panel of firms and balancing is unlikely to be exogenous (firms' closure), given that we are interested in the financial "health" of the firms.

ec475

① d) we will need to model the process that determines the balance sheet if we think that it depends on the financial health of the firm.

② Observations $y = X\beta + \varepsilon$

$$\left\{ \begin{array}{l} A_1 \quad r(X) = k \\ A_2 \quad y_i = X_i\beta + \varepsilon_i, E(\varepsilon_i) = 0 \\ A_3 \quad E(\varepsilon_i | X) = 0 \\ A_4 \quad E(\varepsilon\varepsilon' | X) = \sigma_\varepsilon^2 I_s \\ A_5 \quad \varepsilon | X \sim \text{pdf}(0, \sigma_\varepsilon^2 I_s) \end{array} \right.$$

a) $\varepsilon | X \sim \mathcal{N}(0, \Omega)$ with $\Omega = \begin{pmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_s^2 \end{pmatrix}$

The assumption A_4 is violated.

The OLS estimator is unbiased (and consistent under additional assumptions) and no longer BLUE because the disturbance term is not homoskedastic. If we know $\sigma_1^2, \dots, \sigma_s^2$ we could use GLS to obtain a BLUE estimator. If we do not know them we could use FGLS to obtain an asymptotically efficient estimator.

b) $y_t = \beta_1 + \beta_2 y_{t-1} + \dots + \underbrace{\rho \varepsilon_{t-1}} + \varepsilon_t$

Then clearly A_3 is violated $E(\varepsilon_t | X) \neq 0$ as ε_t depends on ε_{t-1} as y_{t-1} .

2) b) as A_3 is violated OLS is not unbiased.
 moreover A_3 s.t. $E(\varepsilon_s | x_s) = 0$ is
 also violated so OLS is not even consistent.

(6)

c) This violates A_2 , we have in general:

$$y_i = F(x_i, \varepsilon_i, \beta)$$

where F is a non-linear function.

Then if unobserved heterogeneity is additive

$$y_i = F(x_i, \beta) + \varepsilon_i$$

we could use NLS or NLE methods.

In general, OLS will not estimate the true parameters β .

d) We have: $y_s^* = x_s' \beta + \varepsilon_s$
 and we observe $y_s = \begin{cases} y_s^* & \text{if } y_s^* \geq 1000 \\ 0 & \text{otherwise} \end{cases}$

This is a truncated sample.

Assuming $\varepsilon | X \sim N(0, \sigma_\varepsilon^2 I_s)$

we get
$$E(y_s^* | y_s^* \geq 1000, x_s) = x_s' \beta + E(\varepsilon_s | \varepsilon_s \geq 1000 - x_s' \beta, x_s)$$

$$= \sigma_\varepsilon^2 \left(\frac{\phi\left(\frac{1000 - x_s' \beta}{\sigma_\varepsilon}\right)}{1 - \Phi\left(\frac{1000 - x_s' \beta}{\sigma_\varepsilon}\right)} \right)$$

Again clearly A_2 is violated.

OLS will be biased and inconsistent for β . As in our sample, we need to model

and not $E(y_s^* | x_s)$.

$$E(y_s | x_s) = E(y_s^* | x_s, y_s^* \geq 1000)$$

We could use NLS or NLE.