

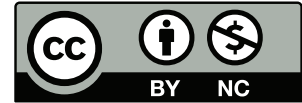
ECO475H1 S

Applied Econometrics II

Tianyu Du

February 24, 2019

This work is licensed under a Creative Commons “Attribution-NonCommercial 4.0 International” license.



Github Page https://github.com/TianyuDu/Spikey_UofT_Notes

Note Page TianyuDu.com/notes

Contents

1	Topic 1 Binary Outcome Model	3
1.1	Lecture 1. Jan. 10 2019	3
1.1.1	Model	3
1.1.2	Data Structure	3
1.1.3	Linear Probability Model	3
1.1.4	Model with Binomial Distribution	4
2	Lecture 3. Jan. 24 2019	5
2.1	Two Side Censoring MLE	5
2.2	Two Side Truncated MLE	6
3	Lecture 4. Jan. 31 2019	7
3.1	Tobit and Sample Selection	7
3.2	Heckman Estimation (Two-Step Procedure)	9
4	Binary Outcome with Continuous Endogenous Regressors:	
	Control Function Approach	10
4.1	Model	10
4.2	Maximum Likelihood Estimator	11
4.3	Control Function	12
5	Treatment Effect and Potential Outcome Model	12
5.1	Lecture 6. Feb. 18 2019	12
5.1.1	Model	12
5.1.2	Data Structure	13
5.1.3	Linear Models: Homogeneous Treatment Effect Problem	13
5.1.4	Linear Models: Endogeneity Problem	13
5.1.5	Comparing Two Models	14

5.1.6	POM Assumptions	14
-------	---------------------------	----

1 Topic 1 Binary Outcome Model

1.1 Lecture 1. Jan. 10 2019

1.1.1 Model

Interpretation Binary outcome models can be interpreted as rational individuals are making binary decisions by comparing the utilities resulting from each decision.

$$\begin{cases} u_1(\mathbf{x}, \varepsilon_1) & \text{utility from decision 1} \\ u_0(\mathbf{x}, \varepsilon_0) & \text{utility from decision 0} \end{cases} \quad (1.1)$$

and rationality of individuals suggests

$$y = \begin{cases} 1 & \text{if } u_1(\mathbf{x}_i, \varepsilon_1) \geq u_0(\mathbf{x}_i, \varepsilon_0) \\ 0 & \text{otherwise} \end{cases} \quad (1.2)$$

Assumption 1.1. By convention, the edge case that $u_0 = u_1$ is modelled and classified to be in the first case, but $\mathbb{P}[u_1 - u_0 = 0 | \mathbf{x}] = 0$ for continuous utilities, so the edge does not really matter.

Assumption 1.2. The difference in utilities can be captured by a linear function, and we define the **latent variable** y^* as

$$y^*(\mathbf{x}) \equiv u_1(\mathbf{x}, \varepsilon_1) - u_0(\mathbf{x}, \varepsilon_0) \quad (1.3)$$

$$= \mathbf{x}'\beta + \varepsilon \quad (1.4)$$

Remark 1.1 (About shapes of variables).

$$\mathbf{x}, \beta \in \mathbb{M}_{k \times 1}; y, d \in \mathbb{R} \quad (1.5)$$

Reduced Model Then the primary model (1.2) is reduced to

$$y = \mathbb{1}\{y^*(\mathbf{x}) \geq 0\} \quad (1.6)$$

1.1.2 Data Structure

For each individual, the observation contains two *primary variables*, (\mathbf{x}_i, y_i) . So the dataset would look like

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^N \quad (1.7)$$

1.1.3 Linear Probability Model

Problems

1. Out-of-range prediction;
2. Homogeneous marginal effect;
3. (Potentially) heteroskedasticity.

1.1.4 Model with Binomial Distribution

Procedures

1. Construct distribution and density functions.
2. Construct likelihood from density function, assuming iid of observations.
3. (MLE) estimates model parameters.
4. Predictions, both individual and average effect.

Conditional Density

$$f_Y(y|\mathbf{x}) = \begin{cases} P(\mathbf{x}) & \text{if } y = 1 \\ 1 - P(\mathbf{x}) & \text{if } y = 0 \end{cases} \quad (1.8)$$

$$= P(\mathbf{x})^d [1 - P(\mathbf{x})]^{1-d} \quad (1.9)$$

Assumption 1.3 (Distributions of ε). Generally, there are two assumptions regarding to the distribution of ε ,

- (i) **Standard Normal** $\varepsilon \sim \mathcal{N}(0, 1)^1$.
- (ii) **Gumbel Distribution** $F_\varepsilon(x) = \frac{e^x}{1+e^x}$.

Note that, in either case, ε is symmetrically distributed, which means ε and $-\varepsilon$ have the identical density and distribution.

$$P(\mathbf{x}) \equiv \mathbb{P}[y^* \geq 0|\mathbf{x}] = \mathbb{P}[-\varepsilon \leq \mathbf{x}'\beta|\mathbf{x}] \equiv F_\varepsilon(\mathbf{x}'\beta) \quad (1.10)$$

$$\implies f_Y(y|\mathbf{x}) = (F_\varepsilon(\mathbf{x}'\beta))^y (1 - F_\varepsilon(\mathbf{x}'\beta))^{1-y} \quad (1.11)$$

Likelihood assuming observations are iid, the joint density, conditioned on model parameter β , of the collection of observations $\{(y_i)\}_{i=1}^N$ is

$$f_{\{(Y_i)\}_{i=1}^N}(\{(y_i)\}_{i=1}^N|\beta, \{(\mathbf{x}_i)\}_{i=1}^N) = \prod_{i=1}^N f_Y(y_i|\mathbf{x}_i, \beta) \quad (1.12)$$

$$= \prod_{i=1}^N F_\varepsilon(\mathbf{x}_i'\beta)^{y_i} (1 - F_\varepsilon(\mathbf{x}_i'\beta))^{1-y_i} \quad (1.13)$$

and the likelihood of collection of observations $\{(\mathbf{x}_i, y_i, d_i)\}_{i=1}^N$ is given parameter β is

$$\mathcal{L}(\beta|\{(\mathbf{x}_i, y_i, d_i)\}_{i=1}^N) = \prod_{i=1}^N F_\varepsilon(\mathbf{x}_i'\beta)^{y_i} (1 - F_\varepsilon(\mathbf{x}_i'\beta))^{1-y_i} \quad (1.14)$$

$$\implies \ln \mathcal{L} = \sum_{i=1}^N y_i \ln F_\varepsilon(\mathbf{x}_i'\beta) + (1 - y_i) \ln(1 - F_\varepsilon(\mathbf{x}_i'\beta)) \quad (1.15)$$

¹In the normal distribution case, we cannot estimate the variance of $\hat{\beta}_{MLE}$ later if we do not assume the variance to be 1.

Maximum Likelihood Estimation First order conditions

$$\frac{\partial \ln \mathcal{L}}{\partial \beta} = \sum_{i=1}^N y_i \frac{F'_\varepsilon(\mathbf{x}'_i \beta)}{F_\varepsilon(\mathbf{x}'_i \beta)} \mathbf{x}'_i - (1 - y_i) \frac{F'_\varepsilon(\mathbf{x}'_i \beta)}{1 - F_\varepsilon(\mathbf{x}'_i \beta)} \mathbf{x}'_i \quad (1.16)$$

$$= \sum_{i=1}^N \frac{y_i F'_\varepsilon(\mathbf{x}'_i \beta) - F'_\varepsilon(\mathbf{x}'_i \beta) F_\varepsilon(\mathbf{x}'_i \beta)}{F_\varepsilon(\mathbf{x}'_i \beta) (1 - F_\varepsilon(\mathbf{x}'_i \beta))} \mathbf{x}'_i = 0 \quad (1.17)$$

Prediction - Marginal Effect With binary outcome model, according to mean of Binomial distribution,

$$\mathbb{E}[y|\mathbf{x}] = P(\mathbf{x}) \equiv F_\varepsilon(\mathbf{x}' \hat{\beta}_{MLE}) \quad (1.18)$$

Then the marginal effect is

$$\frac{\partial \mathbb{E}[y|\mathbf{x}]}{\partial x_k} = F'_\varepsilon(\mathbf{x}' \hat{\beta}_{MLE}) \hat{\beta}_{j,MLE} \quad (1.19)$$

In contrast to the linear probability model, in which the marginal effect is assumed to be homogeneous. In BOM prediction, the $F'_\varepsilon(\mathbf{x}' \hat{\beta}_{MLE})$ term in the marginal effect acts as a source of **heterogeneity** marginal effect.

Remark 1.2. Given $f_\varepsilon \geq 0$, $\hat{\beta}_{MLE}$ reports the **sign** of marginal effect, but it **provides no quantitative implication**.

Prediction - Average Marginal Effect There are two methods to estimate the average marginal effect, these two methods generate different estimations unless the density function of ε is linear.

$$(i) \overline{ME}_j = \frac{1}{N} \sum_{i=1}^N F'_\varepsilon(\mathbf{x}_i \hat{\beta}) \hat{\beta}_j$$

$$(ii) \overline{ME}_j = F'_\varepsilon(\bar{\mathbf{x}}_i \hat{\beta}) \hat{\beta}_j$$

2 Lecture 3. Jan. 24 2019

2.1 Two Side Censoring MLE

Consider the latent dependent variable

$$Y^* = \mathbf{x}' \beta + \epsilon \quad (2.1)$$

where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$.

Therefore, given fixed \mathbf{x} ,

$$Y^* \sim \mathcal{N}(\mathbf{x}' \beta, \sigma^2) \quad (2.2)$$

Define parameter set

$$\theta \equiv (\beta, \sigma) \quad (2.3)$$

²In general, we can assume the error variance to be σ^2 when the dependent variable is *quantitative*, but with *qualitative* dependent variables, we assume $\varepsilon \sim \mathcal{N}(0, 1)$ since we don't have sufficient information to estimate the error variance.

The observable variable is

$$Y = \begin{cases} U & \text{if } Y^* \geq U \\ Y^* & \text{if } Y^* \in (L, U) \\ L & \text{if } Y^* \leq L \end{cases} \quad (2.4)$$

Let $f_Y(y|\mathbf{x}, \boldsymbol{\beta}) : [L, U] \rightarrow [0, 1]$ be the probability measure of Y .

Let $y \in [L, U]$,

$$f_Y(y|\mathbf{x}, \boldsymbol{\beta}) = \begin{cases} \mathbb{P}(Y^* \geq U|\mathbf{x}, \boldsymbol{\beta}) & \text{if } y \geq U \\ f_{Y^*}(y|\mathbf{x}, \boldsymbol{\beta}) & \text{if } y \in (L, U) \\ \mathbb{P}(Y^* \leq L|\mathbf{x}, \boldsymbol{\beta}) & \text{if } y \leq L \end{cases} \quad (2.5)$$

$$= \begin{cases} 1 - F_{Y^*}(U|\mathbf{x}, \boldsymbol{\beta}) & \text{if } y \geq U \\ f_{Y^*}(y|\mathbf{x}, \boldsymbol{\beta}) & \text{if } y \in (L, U) \\ F_{Y^*}(L|\mathbf{x}, \boldsymbol{\beta}) & \text{if } y \leq L \end{cases} \quad (2.6)$$

Define indicator $(d_1(y), d_2(y), d_3(y))$ as

$$d_1(y) \equiv \mathcal{I}(y \geq U) \quad (2.7)$$

$$d_2(y) \equiv \mathcal{I}(y \in (L, U)) \quad (2.8)$$

$$d_3(y) \equiv \mathcal{I}(y \leq L) \quad (2.9)$$

Then the probability measure of Y can be expressed as

$$f_Y(y|\mathbf{x}, \boldsymbol{\beta}) = (1 - F_{Y^*}(U|\mathbf{x}, \boldsymbol{\beta}))^{d_1} \times f_{Y^*}(y|\mathbf{x}, \boldsymbol{\beta})^{d_2} \times F_{Y^*}(L|\mathbf{x}, \boldsymbol{\beta})^{d_3} \quad (2.10)$$

Suppose samples are i.i.d., the joint density is

$$f_{Y_1, \dots, Y_N}(y_1, \dots, y_N|\mathbf{X}, \boldsymbol{\beta}) = \prod_{i=1}^N f_Y(y_i|\mathbf{x}_i, \boldsymbol{\beta}) \quad (2.11)$$

The log-likelihood is

$$\mathcal{L}_N(\boldsymbol{\theta}|\mathbf{X}) = \sum_{i=1}^N \left\{ d_{1,i} \times \ln(1 - F_{Y^*}(U|\mathbf{x}_i, \boldsymbol{\beta})) + d_{2,i} \times \ln(f_{Y^*}(y_i|\mathbf{x}_i, \boldsymbol{\beta})) + d_{3,i} \times \ln(F_{Y^*}(L|\mathbf{x}_i, \boldsymbol{\beta})) \right\} \quad (2.12)$$

Finally, solving

$$\hat{\boldsymbol{\theta}}_{MLE} = (\hat{\boldsymbol{\beta}}_{MLE}, \hat{\sigma}_{MLE}) = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \mathcal{L}_N(\boldsymbol{\theta}) \quad (2.13)$$

2.2 Two Side Truncated MLE

Suppose the observations are truncated with lower and upper bounds L and U .

Let the latent dependent variable be

$$Y^* = \mathbf{x}'\boldsymbol{\beta} + \epsilon \quad (2.14)$$

and

$$\epsilon \sim \mathcal{N}(0, \sigma^2) \quad (2.15)$$

which implies, for given \mathbf{x} ,

$$Y^* \sim \mathcal{N}(\mathbf{x}'\boldsymbol{\beta}, \sigma^2) \quad (2.16)$$

Define parameter set

$$\boldsymbol{\theta} \equiv \{\boldsymbol{\beta}, \sigma\} \quad (2.17)$$

Observable random variable Y is

$$Y = \begin{cases} Y^* & \text{if } Y^* \in (L, U) \\ -- & \text{if } Y^* \notin (L, U) \end{cases} \quad (2.18)$$

Constructing the distribution for Y , note that F_Y is only defined on $y \in (L, U)$,

$$F_Y(y|\mathbf{x}, \boldsymbol{\theta}) = \mathbb{P}(Y < y|\mathbf{x}, \boldsymbol{\theta}) \quad (2.19)$$

$$= \frac{\mathbb{P}(Y^* < y \wedge Y^* \in (L, U)|\mathbf{x}, \boldsymbol{\theta})}{\mathbb{P}(Y^* \in (L, U)|\mathbf{x}, \boldsymbol{\theta})} \quad (2.20)$$

$$= \frac{\mathbb{P}(Y^* \in (L, y)|\mathbf{x}, \boldsymbol{\theta})}{\mathbb{P}(Y^* \in (L, U)|\mathbf{x}, \boldsymbol{\theta})} \quad (2.21)$$

$$= \frac{F_{Y^*}(y|\mathbf{x}, \boldsymbol{\theta}) - F_{Y^*}(L|\mathbf{x}, \boldsymbol{\theta})}{F_{Y^*}(U|\mathbf{x}, \boldsymbol{\theta}) - F_{Y^*}(L|\mathbf{x}, \boldsymbol{\theta})} \quad (2.22)$$

Then construct the density of Y

$$f_Y(y|\mathbf{x}, \boldsymbol{\theta}) = \frac{\partial F_Y(y|\mathbf{x}, \boldsymbol{\theta})}{\partial y} \quad (2.23)$$

$$= \frac{f_{Y^*}(y|\mathbf{x}, \boldsymbol{\theta})}{F_{Y^*}(U|\mathbf{x}, \boldsymbol{\theta}) - F_{Y^*}(L|\mathbf{x}, \boldsymbol{\theta})} \quad (2.24)$$

The sample log-likelihood is

$$\mathcal{L}_N(\boldsymbol{\theta}) = \sum_{i=1}^N \ln(f_{Y^*}(y_i|\mathbf{x}_i, \boldsymbol{\theta})) - \ln(F_{Y^*}(U|\mathbf{x}_i, \boldsymbol{\theta}) - F_{Y^*}(L|\mathbf{x}_i, \boldsymbol{\theta})) \quad (2.25)$$

and the estimator is given by

$$\hat{\boldsymbol{\theta}}_{MLE} = \{\hat{\boldsymbol{\beta}}_{MLE}, \hat{\sigma}_{MLE}\} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \mathcal{L}_N(\boldsymbol{\theta}) \quad (2.26)$$

3 Lecture 4. Jan. 31 2019

3.1 Tobit and Sample Selection

Model the *observable* variables in Tobit model with sample selection are determined by both **outcome equation** and **selection equation**.

$$y_i = \begin{cases} \mathbf{x}_i'\boldsymbol{\beta} + \epsilon_i & \text{if } \mathbf{w}_i'\boldsymbol{\gamma} + v_i > 0 \\ \mathbf{x} & \text{otherwise} \end{cases} \quad (3.1)$$

where **unmeasurable errors** are assumed to follow joint normal distribution,

$$\begin{pmatrix} \epsilon_i \\ v_i \end{pmatrix} \sim \mathcal{N}(\mathbf{0}, \begin{bmatrix} \sigma_\epsilon^2 & \rho\sigma_\epsilon \\ \rho\sigma_\epsilon & 1 \end{bmatrix}) \quad (3.2)$$

Lemma 3.1. Pending Review, the variance of decomposed e If (ϵ, v) follows joint normal distribution, then there exists $e \perp v$ and $e \sim \mathcal{N}(0, 1)$ such that

$$\frac{\epsilon}{\sigma_\epsilon} = \rho v + e \quad (3.3)$$

Theorem 3.1. Let

$$(X, Y) \sim \mathcal{N}(\mathbf{0}, \begin{bmatrix} 1 & \rho\sigma_Y \\ \rho\sigma_Y & \sigma_Y^2 \end{bmatrix}) \quad (3.4)$$

Then Y can be decomposed into

$$Y = \rho\sigma_Y X + \eta \quad (3.5)$$

where $\eta \perp X$ and $\eta \sim \mathcal{N}(0, \sigma_Y^2 - \rho^2\sigma_Y^2)$.

Proof. Let

$$\begin{pmatrix} \eta \\ X \end{pmatrix} = \begin{pmatrix} -\rho\sigma_Y & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} X \\ Y \end{pmatrix} \quad (3.6)$$

Note that, given $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$, $\mathbf{Ax} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\Sigma\mathbf{A}')$. Therefore

$$\begin{pmatrix} \eta \\ X \end{pmatrix} \sim \mathcal{N}(\mathbf{0}, \begin{pmatrix} -\rho\sigma_Y & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & \rho\sigma_Y \\ \rho\sigma_Y & \sigma_Y^2 \end{pmatrix} \begin{pmatrix} -\rho\sigma_Y & 1 \\ 1 & 0 \end{pmatrix}) \quad (3.7)$$

$$= \mathcal{N}(\mathbf{0}, \begin{pmatrix} 0 & \sigma_Y^2 - \rho^2\sigma_Y^2 \\ 1 & \rho\sigma_Y \end{pmatrix} \begin{pmatrix} -\rho\sigma_Y & 1 \\ 1 & 0 \end{pmatrix}) \quad (3.8)$$

$$= \mathcal{N}(\mathbf{0}, \begin{pmatrix} \sigma_Y^2 - \rho^2\sigma_Y^2 & 0 \\ 0 & 1 \end{pmatrix}) \quad (3.9)$$

Since the off-diagonal elements of $\mathbf{A}\Sigma\mathbf{A}'$ are all zero, therefore $\eta \perp X$, and $\mathbb{V}[\eta] = (1 - \rho^2)\sigma_Y^2$. ■

Expectation Define $\tilde{\mathbf{x}}_i \equiv [\mathbf{x}_i, \mathbf{w}_i]$, then the expected *observed* dependent variable is ³

$$\mathbb{E}[y|\mathbf{w}'_i\gamma + v_i > 0, \tilde{\mathbf{x}}] \quad (3.10)$$

$$= \mathbb{E}[\mathbf{x}'\beta + \epsilon|\mathbf{w}'_i\gamma + v_i > 0, \tilde{\mathbf{x}}] \quad (3.11)$$

$$= \mathbf{x}'\beta + \mathbb{E}[\epsilon|\mathbf{w}'_i\gamma + v_i > 0, \tilde{\mathbf{x}}] \quad (3.12)$$

$$= \mathbf{x}'\beta + \mathbb{E}[\rho v\sigma_\epsilon + e\sigma_\epsilon|\mathbf{w}'_i\gamma + v_i > 0, \tilde{\mathbf{x}}] \quad (3.13)$$

$$= \mathbf{x}'\beta + \rho\sigma_\epsilon\mathbb{E}[v|\mathbf{w}'_i\gamma + v_i > 0, \tilde{\mathbf{x}}] + \sigma_\epsilon\mathbb{E}[e|\mathbf{w}'_i\gamma + v_i > 0, \tilde{\mathbf{x}}] \quad (3.14)$$

$$= \mathbf{x}'\beta + \rho\sigma_\epsilon\mathbb{E}[v|\mathbf{w}'_i\gamma + v_i > 0, \tilde{\mathbf{x}}] \quad (3.15)$$

Remark 3.1. If $\rho = 0$ in equation (2.9), there is no sample selection problem and we can use OLS to estimate the outcome equation.

Lemma 3.2. If $X \sim \mathcal{N}(\mu, \sigma^2)$ then

$$\mathbb{E}[X|X > \alpha] = \mu + \sigma \frac{\phi(\frac{\alpha - \mu}{\sigma})}{1 - \Phi(\frac{\alpha - \mu}{\sigma})} \quad (3.16)$$

³For each variable, the i subscript is omitted in the derivation

(continue)

$$\dots = \mathbf{x}'\beta + \rho\sigma_\epsilon \mathbb{E}[v|v > -\mathbf{w}'\gamma, \tilde{\mathbf{x}}] \quad (3.17)$$

$$= \mathbf{x}'\beta + \rho\sigma_\epsilon \frac{\phi(-\mathbf{w}'\gamma)}{1 - \Phi(-\mathbf{w}'\gamma)} \quad (3.18)$$

$$= \mathbf{x}'\beta + \rho\sigma_\epsilon \frac{\phi(\mathbf{w}'\gamma)}{\Phi(\mathbf{w}'\gamma)} \quad (3.19)$$

$$= \mathbf{x}'\beta + \rho\sigma_\epsilon \lambda(\mathbf{w}'\gamma) \quad (3.20)$$

where $\lambda(x)$ is the **inverse Mill's ratio** of standard normal at x .

Marginal Effect Consider the case

$$\exists x_k \in \mathbf{x} \cap \mathbf{w} \quad (3.21)$$

for instance, x_k can be *wage taxation*. The marginal effect of x_k is

$$\frac{\partial \mathbb{E}[y|\mathbf{w}'\gamma + v > 0, \tilde{\mathbf{x}}]}{\partial x_k} = \frac{\partial \mathbf{x}'\beta + \rho\sigma_\epsilon \lambda(\mathbf{w}'\gamma)}{\partial x_k} \quad (3.22)$$

$$= \beta_k + \rho\sigma_\epsilon \lambda'(\mathbf{w}'\gamma)\gamma_k \quad (3.23)$$

$$(3.24)$$

where β_k measures the **direct effect** and $\lambda'(\mathbf{w}'\gamma)\gamma_k$ measures the **indirect effect** of x_k .

3.2 Heckman Estimation (Two-Step Procedure)

Step 1 Run a *probit* estimation on the selection equation.

MLE gives

- (i) An estimation $\hat{\gamma}_{MLE}$ captures the *indirect effect* of regressors in \mathbf{w} on y through the selection equation.

And compute

$$\hat{\lambda}(\mathbf{w}'\hat{\gamma}_{MLE}) \equiv \frac{\phi(\mathbf{w}'\hat{\gamma}_{MLE})}{\Phi(\mathbf{w}'\hat{\gamma}_{MLE})} \quad (3.25)$$

Step 2 Run OLS

$$y = \mathbf{x}'\beta + \rho\sigma_\epsilon \hat{\lambda} + \eta \text{ where } \mathbb{E}[\eta|\mathbf{x}, \hat{\lambda}] = 0 \quad (3.26)$$

OLS gives

- (i) An estimation $\hat{\beta}_{OLS}$ measures the *direct effect* of regressors in \mathbf{x} on y through the outcome equation.
- (ii) An estimation of $\widehat{\rho\sigma_\epsilon}$, given $\sigma_\epsilon > 0$, we can estimate the *sign* of ρ .

Special Case (i) Consider the special case where

$$\mathbf{w} = \mathbf{x} \quad (3.27)$$

$$\lambda(x) \text{ is linear} \quad (3.28)$$

then (2.14) and regression (2.20) can be written as

$$y = \mathbf{x}'\beta + \rho\sigma_\epsilon\mathbf{x}'\lambda(\gamma) + \eta \quad (3.29)$$

$$= \mathbf{x}'[\beta + \rho\sigma_\epsilon\lambda(\gamma)] + \eta \quad (3.30)$$

where $\beta + \rho\sigma_\epsilon\lambda(\gamma)$ represents the *mixed and non-separable* effect.

Special Case (ii) If

$$\mathbf{w} = [\mathbf{x}, z] \quad (3.31)$$

$$\lambda(x) \text{ is linear} \quad (3.32)$$

$$(3.33)$$

Let the coefficients of \mathbf{w} be $[\gamma, \theta]$, then

$$\lambda(\mathbf{w}[\gamma, \theta]) = \lambda(\mathbf{x}\gamma) + \lambda(z\theta) \quad (3.34)$$

$$= \mathbf{x}\lambda(\gamma) + z\lambda(\theta) \quad (3.35)$$

Then the regression can be rewritten as

$$y = \mathbf{x}'[\beta + \rho\sigma_\epsilon\lambda(\gamma)] + \rho\sigma_\epsilon z\lambda(\theta) + \eta \quad (3.36)$$

Remark 3.2. Therefore, if λ is linear, we need at least one exclusion variable to identify the direct and indirect effects. If λ is non-linear, it's *probably* fine.

4 Binary Outcome with Continuous Endogenous Regressors: Control Function Approach

4.1 Model

In ordinary binary outcome models, like Probit models, we assumed all regressors are *exogenous* ($Cov(x, \varepsilon) = 0$). But in many cases, we have some of the explanatory variables are endogenous. In this section, we are going to consider the case where the endogenous regressors are **continuous**.
Outcome Equation

$$y = \mathbb{I}\{\mathbf{x}_y\theta + \mathbf{w}\gamma + \varepsilon > 0\} \quad (4.1)$$

where

- (i) \mathbf{x}_y : exogenous observable characteristics.
- (ii) \mathbf{w} : endogenous observable regressors, which are continuous.

Similarly to the IV approach, we use another "auxiliary equation" to estimate \mathbf{w} :

$$\mathbf{w} = \mathbf{x}_w \eta + \sigma_w v \quad (4.2)$$

where the error terms in (3.1) and (3.2) follows

$$\begin{pmatrix} \varepsilon \\ v \end{pmatrix} \sim \mathcal{N}(\mathbf{0}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}) \quad (4.3)$$

Define $\tilde{\mathbf{x}} \equiv [\mathbf{x}_y, \mathbf{x}_w]$ as the set of usable regressors.

4.2 Maximum Likelihood Estimator

To estimate the model using MLE, we need to construct the likelihood function. By Bayesian Theorem,

$$f(y|\mathbf{w}, \tilde{\mathbf{x}}) = \frac{f(y, \mathbf{w}|\tilde{\mathbf{x}})}{f(\mathbf{w}|\tilde{\mathbf{x}})} \quad (4.4)$$

$$\iff f(y, \mathbf{w}|\tilde{\mathbf{x}}) = f(y|\mathbf{w}, \tilde{\mathbf{x}})f(\mathbf{w}|\tilde{\mathbf{x}}) \quad (4.5)$$

By equation (3.2)

$$w|\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{x}_w \eta, \sigma_w^2) \quad (4.6)$$

$$\implies f(w|\tilde{\mathbf{x}}) = \frac{1}{\sqrt{2\pi}\sigma_w} e^{-\frac{(w - \mathbf{x}_w \eta)^2}{2\sigma_w^2}} \quad (4.7)$$

and to compute $f(y|w, \mathbf{x}_w \eta)$, since y is binary, we are going to compute $\mathbb{P}[y = 1|w, \tilde{\mathbf{x}}]$ first.

$$\mathbb{P}[y = 1|w, \tilde{\mathbf{x}}] = \mathbb{P}[-\varepsilon < \mathbf{x}_y \theta + w\gamma|w, \tilde{\mathbf{x}}] \quad (4.8)$$

$$= \mathbb{P}[-\varepsilon < \mathbf{x}_y \theta + w\gamma|v, \tilde{\mathbf{x}}] \quad (4.9)$$

Lemma 4.1. Given joint normal variables (ε, v) conditioned on $\tilde{\mathbf{x}}$ following

$$\begin{pmatrix} \varepsilon \\ v \end{pmatrix} | \tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}) \quad (4.10)$$

then

$$\varepsilon|v, \tilde{\mathbf{x}} \sim \mathcal{N}(\rho v, 1 - \rho^2) \quad (4.11)$$

which implies

$$\frac{\varepsilon - \rho v}{\sqrt{1 - \rho^2}} \sim \mathcal{N}(0, 1) \quad (4.12)$$

$$\implies \frac{-\varepsilon + \rho v}{\sqrt{1 - \rho^2}} \sim \mathcal{N}(0, 1) \quad (4.13)$$

by the symmetry of standard normal distribution.

Therefore,

$$\mathbb{P}[-\varepsilon < \mathbf{x}_y \theta + w\gamma|v, \tilde{\mathbf{x}}] \quad (4.14)$$

$$= \mathbb{P}[-\varepsilon + \rho v < \mathbf{x}_y \theta + w\gamma + \rho v|v, \tilde{\mathbf{x}}] \quad (4.15)$$

$$= \mathbb{P}\left[\frac{-\varepsilon + \rho v}{\sqrt{1 - \rho^2}} < \frac{\mathbf{x}_y \theta + w\gamma + \rho v}{\sqrt{1 - \rho^2}} | v, \tilde{\mathbf{x}}\right] \quad (4.16)$$

$$= \Phi\left(\frac{\mathbf{x}_y \theta + w\gamma + \rho v}{\sqrt{1 - \rho^2}}\right) \quad (4.17)$$

4.3 Control Function

Step 1 Run OLS on $w = \mathbf{x}_w\eta + \sigma_w v$, Obtain estimations $\hat{\eta}_{OLS}$, $\hat{\sigma}_{wOLS}$.

Step 2 Obtain estimation of v using the error terms and standard deviation in OLS results.

$$\hat{v} = \frac{w - \mathbf{x}_w\hat{\eta}_{OLS}}{\hat{\sigma}_{OLS}} \quad (4.18)$$

Step 3 Plug in \hat{v} and run **probit** model in (3.17),

$$\Phi\left(\frac{\mathbf{x}_y\theta + w\gamma + \rho v}{\sqrt{1 - \rho^2}}\right) \quad (4.19)$$

$$= \Phi\left(\frac{\mathbf{x}_y\theta}{\sqrt{1 - \rho^2}} + \frac{w\gamma}{\sqrt{1 - \rho^2}} + \frac{\rho v}{\sqrt{1 - \rho^2}}\right) \quad (4.20)$$

Define

$$\theta^* \equiv \frac{\theta}{\sqrt{1 - \rho^2}} \quad (4.21)$$

$$\gamma^* \equiv \frac{\gamma}{\sqrt{1 - \rho^2}} \quad (4.22)$$

$$\alpha^* \equiv \frac{\alpha}{\sqrt{1 - \rho^2}} \quad (4.23)$$

So the probit model can be written as

$$y = \mathbb{1}\{-\tilde{u} < \mathbf{x}_y\theta^* + w\gamma^* + v\alpha^*\} \quad (4.24)$$

where $\tilde{u} \sim \mathcal{N}(0, 1)$.

Once we have an estimation on α^* , ρ can be calculated with

$$\rho = \pm \sqrt{\frac{\alpha^*}{1 + \alpha^{*2}}} \quad (4.25)$$

5 Treatment Effect and Potential Outcome Model

5.1 Lecture 6. Fen. 18 2019

5.1.1 Model

$$Y = \begin{cases} Y_1 & \text{if } D = 1 \\ Y_0 & \text{if } D = 0 \end{cases} \quad (5.1)$$

which is equivalent to

$$Y = Y_1 D + Y_0 (1 - D) \quad (5.2)$$

Definition 5.1. (Y_0, Y_1) are called **potential outcomes**, only one of them will be realized and observed, depends on the *decision*, D , of individuals.

Remark 5.1. Rationality of individuals suggests that while people are making decision, they will choose the one associated with higher potential outcome. But in this model, individuals might not be perfectly informed about the potential outcomes.

Definition 5.2. Individual treatment effect is defined as

$$Y_1 - Y_0 \quad (5.3)$$

Each individual is associated with one realization of above random variable. So the individual treatment effect depends on individual characteristics X , and can be written as $(Y_1 - Y_0)(X)$. **Average treatment effect (ATE)** is defined as

$$\mathbb{E}[Y_1 - Y_0] \quad (5.4)$$

And **Conditional ATE** is defined as

$$\mathbb{E}[Y_1(X) - Y_2(X)|X = x] \equiv \alpha(x) \quad (5.5)$$

5.1.2 Data Structure

$$(X, D, Y) \rightarrow \{(x_i, d_i, y_i)\}_{i=1}^N \quad (5.6)$$

5.1.3 Linear Models: Homogeneous Treatment Effect Problem

Model (5.2) can be rewritten as

$$Y = (Y_1 - Y_0)D + Y_0 \quad (5.7)$$

the D and Y are observable, an OLS linear model would give

$$Y = \alpha D + \varepsilon \quad (5.8)$$

$$\text{where } \alpha \equiv Y_1 - Y_0, \varepsilon \equiv Y_0 \quad (5.9)$$

OLS methods treat the individual treatment effect as a constant α . Thus linear models assume a *homogenous treatment effect*, which is unrealistic.

Remark 5.2. Homogeneous treatment effect is more realistic if the decision is made on macro level.

Remark 5.3. If we only wish to extract the *correlation* instead of the *causal effect*, using OLS is fine.

5.1.4 Linear Models: Endogeneity Problem

Remark 5.4. This section needs to be revised, equation 5.11, 5.12

Suppose the treatment effect is heterogeneous, we can write (5.2) as

$$Y = (Y_1 - Y_0)D + Y_0 = \alpha(\textcolor{red}{X})D + Y_0(X) \quad (5.10)$$

assuming both potential outcomes are linear in individual characteristics

$$Y_1 = \alpha_1 + \beta_1 x + \varepsilon_1 \quad (5.11)$$

$$Y_0 = \alpha_0 + \beta_0 x + \varepsilon_0 \quad (5.12)$$

then (5.10) becomes

$$Y = [(\alpha_1 - \alpha_0) + (\beta_1 - \beta_0)x + (\varepsilon_1 - \varepsilon_0)]D + \alpha_0 + \beta_0 x + \varepsilon_0 \quad (5.13)$$

$$= \alpha^*(x)D + \alpha_0 + \beta_0 x + \underbrace{[\varepsilon_1 D + \varepsilon_0(1 - D)]}_{\varepsilon^*} \quad (5.14)$$

$$\text{where } \alpha^*(x) \equiv (\alpha_1 - \alpha_0) + (\beta_1 - \beta_0)x \quad (5.15)$$

Unless $\varepsilon_0 = \varepsilon_1$, the ε^* is correlated with D , which causes endogeneity problem.

Remark 5.5. In Linear models, we assumed $\mathbb{E}[\varepsilon|\mathbf{X}] = 0$, this is **not** a standard assumption, we can assume this because of the constant term in the linear model, which offsets the mean of error term.

5.1.5 Comparing Two Models

$$\begin{cases} \text{POM:} & Y = \alpha(X)D + Y_0(X) \\ \text{LM:} & Y = \alpha D + \beta x + \varepsilon \end{cases} \quad (5.16)$$

5.1.6 POM Assumptions

Assumption 5.1 (Conditional Independent Assumption; Randomized Assignment).

$$\{(Y_1, Y_0) \perp D\} | X \quad (5.17)$$

Interpretation (i)(CIA Aspect) For each individual/batch of individuals with the same characteristics $X = x^*$, the two potential outcomes are independent to their decision D .

Interpretation (ii)(RA Aspect) As mentioned before, rationality of individuals suggests individual will choose the decision associated with higher potential payoff. But we assume it's **not** the case, within each batch of individuals, each individual is assigned a random D , and is forced to take it.

Assumption 5.2 (Overlapping Matching Condition).

$$0 < \mathbb{P}[D = 1 | X = \mathbf{x}] < 1 \quad \forall \mathbf{x} \in \mathcal{D}(X) \quad (5.18)$$

Interpretation For every possible individual characteristics \mathbf{x}^* , there exists realizations of both potential outcomes. So the batch of observations $\{(\mathbf{x}, d, y) | \mathbf{x} = \mathbf{x}^*\}$ is inter-comparable. That is, for each \mathbf{x}^* , we have at least one realization with $d = 0$ and $d = 1$, so we can *construct the counterfactual*.