Econometrics I

Professor William Greene Stern School of Business Department of Economics



Econometrics I

Part 19 – Sample Selection Two Step Estimation

******* This Book is Seriously Flawed, November 19, 2012

By Dennis Hanseman (Cincinnati, OH United States) - See all my reviews

Amazon Verified Purchase (What's this?)

This review is from: Applied Econometrics for Health Economists: A Practical Guide (Paperback)

After paying over \$30 for a 115-page book, I was shocked to find that it was seriously flawed. All of the analysis in the book is based on data from the British "Health and Lifestyle Survey". As Jones points out in Chapter 2, this Survey employed a complex sample design that incorporated stratification, clustering, and -- presumably -- unequal probabilities of selection. Rather than taking the design characteristics into account, Jones analyzes this data as if it came from a simple random sample. As a result, his estimates are likely to be biased, with overstated significance levels.

Dueling Selection Biases – From two emails, same day.

- "I am trying to find methods which can deal with data that is non-randomised and suffers from selection bias."
- "I explain the probability of answering questions using, among other independent variables, a variable which measures knowledge breadth. Knowledge breadth can be constructed only for those individuals that fill in a skill description in the company intranet. <u>This is</u> <u>where the selection bias comes from.</u>

Samples and Populations

- Consistent estimation
 - The sample is randomly drawn from the population
 - Sample statistics converge to their population counterparts
- A presumption: The 'population' is the population of interest.

Implication: If the sample is randomly drawn from a specific subpopulation, statistics converge to the characteristics of that subpopulation

Nonrandom Sampling

- Simple nonrandom samples: Average incomes of airport travelers → mean income in the population as a whole?
- Survivorship: Time series of returns on business performance. Mutual fund performance. (Past performance is no guarantee of future success.)
- Attrition: Drug trials. Effect of erythropoetin on quality of life survey.
- Self-selection:
 - Labor supply models
 - Shere Hite's (1976) "The Hite Report" 'survey' of sexual habits of Americans. "While her books are ground-breaking and important, they are based on flawed statistical methods and one must view their results with skepticism."

The Crucial Element

- Selection on the unobservables
 - Selection into the sample is based on both observables and unobservables.
 - All the observables are accounted for.
 - Unobservables in the selection rule also appear in the model of interest (or are correlated with unobservables in the model of interest).
- "Selection Bias" = the bias due to not accounting for the unobservables that link the equations.

Heckman's Canonical Model

A behavioral model:

Offered wage $= o^* = \beta'x + v$ (x = age, experience, educ...) Reservation wage $= r^* = \delta'z + u$ (z = age, kids, family stuff) Labor force participation:

$$\begin{split} \mathsf{LFP} &= 1 \text{ if } o^* \geq \mathsf{r}^*, \ 0 \text{ otherwise} \\ \mathsf{Prob}(\mathsf{LFP}=1) = \Phi\Big[(\beta' \mathsf{x} \cdot \delta' \mathsf{z})/\sqrt{\sigma_\mathsf{v}^2 + \sigma_\mathsf{u}^2}\Big] \\ \mathsf{Desired Hours} &= \mathsf{H}^* = \gamma' \mathbf{w} + \varepsilon \\ \mathsf{Actual Hours} &= \mathsf{H}^* \text{ if } \mathsf{LFP} = 1 \\ & \mathsf{unobserved if } \mathsf{LFP} = 0 \\ \varepsilon \text{ and } \mathsf{u} \text{ are correlated. } \varepsilon \text{ and } \mathsf{v} \text{ might be correlated.} \\ \mathsf{What is } \mathsf{E}[\mathsf{H}^* \mid \mathsf{w}, \mathsf{LFP} = 1]? \text{ Not } \gamma' \mathsf{w}. \end{split}$$

Standard Sample Selection Model

$$\begin{aligned} \mathbf{d}_{i}^{*} &= \boldsymbol{\alpha}' \mathbf{z}_{i} + \mathbf{u}_{i} \\ \mathbf{d}_{i} &= \mathbf{1}(\mathbf{d}_{i}^{*} > \mathbf{0}) \\ \mathbf{y}_{i}^{*} &= \boldsymbol{\beta}' \mathbf{x}_{i} + \boldsymbol{\epsilon}_{i} \\ \mathbf{y}_{i} &= \mathbf{y}_{i}^{*} \text{ when } \mathbf{d}_{i} = \mathbf{1}, \text{ unobserved otherwise} \\ (\mathbf{u}_{i}, \mathbf{v}_{i}) &\sim \text{Bivariate Normal}[(\mathbf{0}, \mathbf{0}), (\mathbf{1}, \boldsymbol{\rho} \sigma, \sigma^{2})] \\ \text{E}[\mathbf{y}_{i} \mid \mathbf{y}_{i} \text{ is observed}] &= \text{E}[\mathbf{y}_{i} | \mathbf{d}_{i} = \mathbf{1}] \\ &= \boldsymbol{\beta}' \mathbf{x}_{i} + \text{E}[\boldsymbol{\epsilon}_{i} \mid \mathbf{d}_{i} = \mathbf{1}] \\ &= \boldsymbol{\beta}' \mathbf{x}_{i} + \text{E}[\boldsymbol{\epsilon}_{i} \mid \mathbf{u}_{i} > -\boldsymbol{\alpha}' \mathbf{z}_{i}] \\ &= \boldsymbol{\beta}' \mathbf{x}_{i} + (\boldsymbol{\rho} \sigma) \frac{\boldsymbol{\phi}(\boldsymbol{\alpha}' \mathbf{z}_{i})}{\boldsymbol{\Phi}(\boldsymbol{\alpha}' \mathbf{z}_{i})} \\ &= \boldsymbol{\beta}' \mathbf{x}_{i} + \boldsymbol{\theta} \lambda_{i} \end{aligned}$$

Incidental Truncation u1,u2~N[(0,0),(1,.71,1)



19-10/39

Selection as a Specification Error

- $\Box E[y_i | \mathbf{x}_i, y_i \text{ observed}] = \boldsymbol{\beta}' \mathbf{x}_i + \theta \lambda_i$
- **D** Regression of y_i on \mathbf{x}_i omits λ_i .
 - λ_i will generally be correlated with \mathbf{x}_i if \mathbf{z}_i is.
 - z_i and x_i often have variables in common.
 - There is no specification error if $\theta = 0 \iff \rho = 0$
- **□** "Selection Bias" is plim $(\mathbf{b} \mathbf{\beta})$
- What is "selection bias..."

Control Function

Labor Force Participation

 $d^* = \boldsymbol{\alpha}' \mathbf{z} + \mathbf{u}$

What is u? Unmeasured factors that motivate LFP, u = (m,a)Desired Hours

 $\mathbf{H}^* = \mathbf{\beta}' \mathbf{x} + \mathbf{\varepsilon}$

What is ϵ ? Unmeasured factors that motivate H*, $\epsilon = (m,c)$

 $\epsilon = \rho u + w \quad \epsilon \text{ and } u \text{ share factors, m.}$

 $\mathbf{H}^* = \boldsymbol{\beta}' \mathbf{x} + \boldsymbol{\rho} \mathbf{u} + \mathbf{w}$

Regression of H* on x omits u. λ is the prediction of u.

Note, the problem goes away if $\rho = 0$.

Estimation of the Selection Model

- Two step least squares
 - Inefficient
 - Simple exists in current software
 - Simple to understand and widely used
- Full information maximum likelihood
 - Efficient
 - Simple exists in current software
 - Not so simple to understand widely misunderstood

Estimation

Heckman's two step procedure

- (1) Estimate the probit model and compute λ_i for each observation using the estimated parameters.
- (2) a. Linearly regress y_i on x_i and λ_i using the observed data
 - b. Correct the estimated asymptotic covariance matrix for the use of the estimated λ_i. (An application of Murphy and Topel (1984)
 Heckman was 1979) See text, pp. 953-955.

Variance of a Heckman's Two Step Estimator

The parameters in γ do have to be estimated using the probit equation. Rewrite (19-24) as

$$(y_i|z_i = 1, \mathbf{x}_i, \mathbf{w}_i) = \mathbf{x}'_i \boldsymbol{\beta} + \beta_\lambda \hat{\lambda}_i + v_i - \beta_\lambda (\hat{\lambda}_i - \lambda_i).$$

In this form, we see that in the preceding expression we have ignored both an additional source of variation in the compound disturbance and correlation across observations; the same estimate of γ is used to compute $\hat{\lambda}_i$ for every observation. Heckman has shown that the earlier covariance matrix can be appropriately corrected by adding a term inside the brackets,

$$\mathbf{Q} = \hat{\rho}^2(\mathbf{X}'_* \hat{\boldsymbol{\Delta}} \mathbf{W}) \text{Est.Asy.Var}[\hat{\boldsymbol{\gamma}}](\mathbf{W}' \hat{\boldsymbol{\Delta}} \mathbf{X}_*) = \hat{\rho}^2 \hat{\mathbf{F}} \hat{\mathbf{V}} \hat{\mathbf{F}}',$$

where $\hat{\mathbf{V}} = \text{Est.Asy.Var}[\hat{\boldsymbol{\gamma}}]$, the estimator of the asymptotic covariance of the probit coefficients. Any of the estimators in (17-22) to (17-24) may be used to compute $\hat{\mathbf{V}}$. The complete expression is

Est.Asy.Var[**b**,
$$b_{\lambda}$$
] = $\hat{\sigma}_{\varepsilon}^{2}$ [**X**'_***X**_*]⁻¹[**X**'_*(**I** - \hat{\rho}^{2}\hat{\Delta})**X**_* + **Q**][**X**'_***X**_*]⁻¹.

This is the estimator that is embedded in contemporary software such as *Stata*. We note three useful further aspects of the two-step estimator:

1. This is an application of the two-step procedures we developed in Section 8.4.1 and 14.7 and that were formalized by Murphy and Topel (1985).⁴⁰

19-15/39

Application – Labor Supply

MROZ lak	oor supply data. Cross section, 753 observations
Use LFP	for binary choice, KIDS for count models.
LFP	= labor force participation, 0 if no, 1 if yes.
WHRS	= wife's hours worked. 0 if LFP=0
KL6	= number of kids less than 6
K618	= kids 6 to 18
WA	= wife's age
WE	= wife's education
WW	= wife's wage, 0 if LFP=0.
RPWG	= Wife's reported wage at the time of the interview
HHRS	= husband's hours
HA	= husband's age
HE	= husband's education
HW	= husband's wage
FAMINC	= family income
MTR	= marginal tax rate
WMED	<pre>= wife's mother's education</pre>
WFED	<pre>= wife's father's education</pre>
UN	= unemployment rate in county of residence
CIT	= dummy for urban residence
AX	= actual years of wife's previous labor market experience
AGE	= Age
AGESQ	= Age squared
EARNINGS	S= WW * WHRS
LOGE	= Log of EARNINGS
KIDS	= 1 if kids < 18 in the home.

Labor Supply Model

```
NAMELIST ; Z = One,KL6,K618,WA,WE,HA,HE $
NAMELIST ; X = One,KL6,K618,Age,Agesq,WE,Faminc $
PROBIT ; Lhs = LFP ; Rhs = Z ; Hold(IMR=Lambda) $
SELECT ; Lhs = WHRS ; Rhs = X $
REGRESS ; If [ LFP = 1] ; Lhs = WHRS ; Rhs = X $
REGRESS ; If [ LFP = 1] ; Lhs = WHRS ; Rhs = X,Lambda $
REGRESS ; If [ LFP = 1] ; Lhs = WHRS ; Rhs = X,Lambda $
Cluster = 1 $
```

Participation Equation

Binomial Dependent Log like Restricte Chi squar Significa McFadden Estimatic Inf.Cr.Al Results r	ependent variable LFP og likelihood function -461.37865 estricted log likelihood -514.87320 hi squared [6](P= .000) 106.98911 ignificance level .00000 cFadden Pseudo R-squared .1038985 stimation based on N = 753, K = 7 nf.Cr.AIC = 936.8 AIC/N = 1.244 esults retained for SELECTION model.							
LFP	Coefficient	Standard Error	z	Prob. z >Z *	95% Con Inte	fidence rval		
Constant KL6 K618 WA WE HA HE	Index function for 1.00265** 90400*** 05453 02602* .16039*** 01643 05191**	probabilit .49994 .11434 .04021 .01333 .02774 .01329 .02040	y 2.01 -7.91 -1.36 -1.95 5.78 -1.24 -2.54	.0449 .0000 .1751 .0508 .0000 .2165 .0110	.02277 -1.12811 13334 05214 .10603 04248 09190	1.98252 67989 .02428 .00009 .21475 .00962 01192		

Hours Equation

Sample Selection Model Probit selection equation based on LFP Selection rule is: Observations with LFP = 1 Results of selection:									
Data set Selected 	Data po 79 sample 42	oints Sum 53 28 	of weig] 753.0 428.0	hts 					
Sample Selection Model. Two step least squares regression LHS=WHRS Mean = Standard deviation = 776.27438 Number of observs. = 428 Model size Parameters = 8 Degrees of freedom = 420 Residuals Sum of squares = .226721E+09 Standard error of e = 734.71953 Fit R-squared = .10210 Adjusted R-squared = .08713 Model test F[7, 420] (prob) = 6.8(.0000) Not using OLS or no constant. Rsqrd & F may be < 0									
Correlati and Selec	on of disturband tion Criterion ((Rho) =	ion	84541					
WHRS	Coefficient	Standard Error	z	Prob. z >Z ≭	95% Co: Inte	nfidence erval			
Constant KL6 K618 AGE AGESQ WE FAMINC LAMBDA	2442.27** 115.110 -101.721*** 14.6359 10079 -102.203*** .01379*** -793.857	$\begin{array}{c} 1202.111\\ 282.0086\\ 38.28339\\ 53.19166\\ .61856\\ 39.40963\\ .00345\\ 494.5410\\ \end{array}$	2.03 .41 -2.66 .28 16 -2.59 4.00 -1.61	.0422 .6831 .0079 .7832 .8706 .0095 .0001 .1084	$\begin{array}{r} 86.17\\ -437.617\\ -176.755\\ -89.6178\\ -1.31315\\ -179.445\\ .00703\\ -1763.140\end{array}$	4798.36 667.836 -26.687 118.8897 1.11157 -24.962 .02056 175.426			

19-19/39

Part 19: Sample Selection

Selection "Bias"

WHRS	Coefficient	Standard Error	t	Prob. t >T*	95% Con Inte	fidence rval
Constant	1812.13	1144.333	$ \begin{array}{r} 1.58 \\ -2.99 \\ -4.09 \\ .21 \\42 \\ -2.74 \\ 3.72 \\ \end{array} $.1140	-430.73	4054.98
KL6	-299.128***	100.0331		.0030	-495.189	-103.067
K618	-126.400***	30.87285		.0001	-186.909	-65.890
AGE	11.2795	53.84421		.8342	-94.2532	116.8122
AGESQ	26104	.62633		.6771	-1.48862	.96655
WE	-47.3272***	17.29681		.0065	-81.2283	-13.4260
FAMINC	.01262***	.00339		.0002	.00598	.01926
WHRS	Coefficient	Standard Error	t	Prob. t >T*	95% Con Inte	fidence rval
Constant	2442.27**	1194.817	2.04	.0416	100.47	4784.06
KL6	115.110	252.7874	.46	.6491	-380.345	610.564
K618	-101.721***	33.75941	-3.01	.0027	-167.888	-35.554
AGE	14.6359	53.73825	.27	.7855	-90.6891	119.9610
AGESQ	10079	.63114	16	.8732	-1.33780	1.13623
WE	-102.203***	35.27561	-2.90	.0040	-171.342	-33.064
FAMINC	.01379***	.00344	4.01	.0001	.00704	.02054
LAMBDA	-793.857*	445.1168	-1.78	.0752	-1666.270	78.556

Heckman's corrected standard errors

WHRS	Coefficient	Standard Error	z	Prob. z >Z *	95% Confidence Interval	
Constant	2442.27**	1202.111	2.03	.0422	86.17	4798.36
KL6	115.110	282.0086	.41	.6831	-437.617	667.836
K618	-101.721***	38.28339	-2.66	.0079	-176.755	-26.687
AGE	14.6359	53.19166	.28	.7832	-89.6178	118.8897
AGESQ	10079	.61856	16	.8706	-1.31315	1.11157
WE	-102.203***	39.40963	-2.59	.0095	-179.445	-24.962
FAMINC	.01379***	.00345	4.00	.0001	.00703	.02056
LAMBDA	-793.857	494.5410	-1.61	.1084	-1763.140	175.426

Uncorrected standard errors - OLS

WHRS	Coefficient	Standard Error	t	Prob. t >T*	95% Cc Int	nfidence erval
Constant	2442.27**	1194.817	2.04	.0416	100.47	4784.06
KL6	115.110	252.7874	.46	.6491	-380.345	610.564
K618	-101.721***	33.75941	-3.01	.0027	-167.888	-35.554
AGE	14.6359	53.73825	.27	.7855	-90.6891	119.9610
AGESQ	10079	.63114	16	.8732	-1.33780	1.13623
WE	-102.203***	35.27561	-2.90	.0040	-171.342	-33.064
FAMINC	.01379***	.00344	4.01	.0001	.00704	.02054
LAMBDA	-793.857*	445.1168	-1.78	.0752	-1666.270	78.556

Heteroscedasticity robust standard errors (cluster = 1)

WHRS	Coefficient	Clustered Std.Error	t	Prob. t >T *	95% Co Int	95% Confidence Interval	
Constant	2442.27**	1232.878	1.98	.0482	25.87	4858.66	
KL6	115.110	302.2937	.38	.7036	-477.375	707.594	
K618	-101.721***	33.77584	-3.01	.0028	-167.920	-35.521	
AGE	14.6359	57.43322	.25	.7990	-97.9311	127.2030	
AGESQ	10079	.67949	15	.8822	-1.43257	1.23100	
WE	-102.203***	36.79001	-2.78	.0057	-174.310	-30.096	
FAMINC	.01379***	.00409	3.37	.0008	.00578	.02181	
LAMBDA	-793.857*	470.2006	-1.69	.0921	-1715.433	127.719	

19-21/39

Part 19: Sample Selection

$$\begin{split} & \text{Maximum Likelihood Estimation} \\ & \text{logL} = \sum_{d=1} \ \log \left[\frac{\exp\left(-\frac{1}{2} (\epsilon_i \ / \ \sigma)^2\right)}{\sigma \sqrt{2\pi}} \Phi\left(\frac{\rho(\epsilon_i \ / \ \sigma) + \alpha' \mathbf{z}_i}{\sqrt{1 - \rho^2}}\right) \right] \\ & + \sum_{d=0} \ \log [1 - \Phi(\alpha' \mathbf{z}_i)] \\ & \text{Re parameterize this: let } q_i = \alpha' \mathbf{z}_i \\ & (1) \ \theta = 1/\sigma \\ & (2) \ \gamma = \beta/\sigma \text{ (Olsen transformation)} \\ & (3) \ \tau = \rho/\sqrt{1 - \rho^2} \\ & (4) \text{ Constrain } \rho \text{ to be in (-1,1) by using} \\ & \psi = \frac{1}{2} \ln\left(\frac{1 + \rho}{1 - \rho}\right) = \text{atanh}\rho, \text{ so } \rho = \text{atanh}^{-1}(\psi) = \frac{\exp(2\psi) - 1}{\exp(2\psi) + 1} \\ & \text{logL} = \sum_{d=0} \log \Phi(-q_i) + \sum_{d=1}^{\log \theta - \frac{1}{2} \log 2\pi - \frac{1}{2} (\theta y_i - \gamma' \mathbf{x}_i)^2}{+\log \Phi[\tau(\theta y_i - \gamma' \mathbf{x}) + q_i \sqrt{1 + \tau^2}]} \end{split}$$

			M	ILE				
ML Estima Dependent Log like Estimatic Inf.Cr.A	ates of Selection variable inbood function based on N = C = 7820.9 Å	n Model WH -3894.470 753, K = IC∕N = 10.3	RS 94 16 86				Two Step Es	stimates
WHRS	Coefficient	Standard Error	z	Prob. z >Z*	95% Co Int	nfidence erval	Coefficient	Standard Error
Constant KL6 K618 WA WE HA HE	Selection (prob: 1.01351** 90130*** 05292 02492* .16396*** 01763 05597***	it) equation .51518 .11218 .03999 .01382 .02783 .01379 .01379 .02020	for LFP. 1.97 -8.03 -1.32 -1.80 5.89 -1.28 -2.77	.0492 .0000 .1857 .0714 .0000 .2011 .0056	.00376 -1.12116 13130 05201 .10942 04467 09555	2.02325 68143 .02545 .00217 .21850 .00940 01638	Index function f 1.00265** 90400*** 05453 02602* .16039*** 01643 05191**	or probabil .49994 .11434 .04021 .01333 .02774 .01329 .02040
Constant KL6 K618 AGESQ WE FAMINC SIGMA(1) RH0(1,2)	1946.85* -209.025 -120.969*** 12.0376 22652 -59.2166* 01289*** 748.132*** 22965	1167.225 221.3726 35.45852 51.99025 .59914 33.33120 .00332 59.64630 .49962	1.67 94 -3.41 .23 38 -1.78 3.88 12.54 46	.0953 .3451 .0006 .8169 .7054 .0756 .0001 .0000 .6458	$\begin{array}{r} -340.87\\ -642.907\\ -190.467\\ -89.8615\\ -1.40082\\ -124.5446\\ .00639\\ 631.227\\ -1.20890\end{array}$	4234.56 224.857 -51.472 113.9366 .94777 6.1113 .01940 865.036 .74959	2442.27** 115.110 -101.721*** 14.6359 10079 -102.203*** .01379*** -793.857	1202.111 282.0086 38.28339 53.19166 .61856 39.40963 .00345 494.5410

Standard error corrected for selection \$939.01825 Correlation of disturbance in regression and Selection Criterion (Rho) = -.84541

How to Handle Selectivity

- The 'Mills Ratio' approach just add a 'lambda' to whatever model is being estimated?
 - The Heckman model applies to a probit model with a linear regression.
 - The conditional mean in a nonlinear model is not something "+lambda"
- The model can sometimes be built up from first principles

Received Sunday, April 27, 2014

I have a paper regarding strategic alliances between firms, and their impact on firm risk. While observing how a firm's strategic alliance formation impacts its risk, I need to correct for two types of selection biases. The reviews at Journal of Marketing asked us to correct for the propensity of firms to enter into alliances, and also the propensity to select a specific partner, before we examine how the partnership itself impacts risk.

Our approach involved conducting a probit of alliance formation propensity, take the inverse mills and include it in the second selection equation which is also a probit of partner selection. Then, we include inverse mills from the second selection into the main model. The review team states that this is not correct, and we need an MLE estimation in order to correctly model the set of three equations. The Associate Editor's point is given below. Can you please provide any guidance on whether this is a valid criticism of our approach. Is there a procedure in LIMDEP that can handle this set of three equations with two selection probit models?

AE's comment:

"Please note that the procedure of using an inverse mills ratio is only consistent when the main equation where the ratio is being used is linear. In non-linear cases (like the second probit used by the authors), this is not correct. Please see any standard econometric treatment like Greene or Wooldridge. A MLE estimator is needed which will be far from trivial to specify and estimate given error correlations between all three equations."

19-25/39

A Bivariate Probit Model

Labor Force Participation Equation

 $d^* = \alpha' z + u$ $d = 1(d^* > 0)$ Full Time or Part Time? $f^* = \beta' x + \varepsilon$ $f = 1(f^* > 0)$ Probability Model: Nonparticipant: Prob[d=0] = $\Phi(-\alpha' z)$ Participant and Full Time Prob[f=1,d=1] = Prob[f=1|d=1]Prob[d=1] $= Bivariate Normal(\beta'x,\alpha'z,\rho)$

Participant and Part Time

Prob[f=0,d=1]= Prob[f=0|d=1]Prob[d=1]
= Bivariate Normal(
$$\beta$$
'x,- α 'z,- ρ)

FT/PT Selection Model

+			+		
FIML Est:	imates of Bivari	ate Probit Model	LI		
Dependent	t variable	FULLFP	I		
Weighting	g variable	None	Ι		
Number of observations		753	1	Full Tin	ne = Hours > 1000
Log like	lihood function	-723.9798	I		
Number of	f parameters	16	I		
Selection	n model based on	LFP	I		
+			+		
+	++		-+	+	-++
Variable	Coefficient	Standard Error	b/St.Er.	P[Z >z]	Mean of X
+	tt Trdex equatio	n for FULLTTME	-+	+	-++
Constant	94532822	1 61674948	585	5587	
WW	- 02764944	01941006	-1 424	1543	4 17768154
кт.6	.04098432	26250878	156	8759	14018692
к618	- 13640024	.05930081	-2.300	.0214	1,35046729
AGE	.03543435	07530788	471	6380	41,9719626
AGESO	00043848	.00088406	496	.6199	1821.12150
WE	08622974	.02808185	-3.071	.0021	12.6588785
FAMINC	.210971D-04	.503746D-05	4.188	.0000	24130.4229
:	Index equatio	n for LFP			
Constant	.98337341	.50679582	1.940	.0523	
KL6	88485756	.11251971	-7.864	.0000	.23771580
K618	04101187	.04020437	-1.020	.3077	1.35325365
WA	02462108	.01308154	-1.882	.0598	42.5378486
WE	.16636047	.02738447	6.075	.0000	12.2868526
HA	01652335	.01287662	-1.283	.1994	45.1208499
HE	06276470	.01912877	-3.281	.0010	12.4913679
1	Disturbance corr	elation			
RHO(1,2)	84102682	.25122229	-3.348	.0008	

Part 19: Sample Selection

Building a Likelihood for a Poisson Regression Model with Selection

Poisson Probability Functions

 $P(y_i | x_i) = exp(-\lambda_i)\lambda_i^{y} / y_i!$

Covariates and Unobserved Heterogeneity

 $\lambda(\mathbf{x}_{i'} \varepsilon_i) = \exp(\mathbf{x}_i' \mathbf{\beta} + \varepsilon_i)$

Conditional Contribution to the Log Likelihood

 $\log L_i \mid \varepsilon_i = -\lambda(\mathbf{x}_i, \varepsilon_i) + y_i \log \lambda(\mathbf{x}_i, \varepsilon_i) - \log y_i!$

Probit Selection Mechanism

 $\begin{aligned} d_{i}^{*} &= \mathbf{z}_{i}^{\prime} \mathbf{Y} + u_{i}, \ d_{i} = \mathbf{1}[d_{i}^{*} > 0] \\ [\varepsilon_{i}, u_{i}] &\sim \mathsf{BVN} \begin{bmatrix} 0 \\ 0 \end{pmatrix}, \begin{bmatrix} \sigma^{2} & \rho \sigma \\ \rho \sigma & 1 \end{bmatrix} \\ y_{i}, \mathbf{x}_{i} \text{ observed only when } d_{i} = \mathbf{1}. \end{aligned}$

Building the Likelihood

The Conditional Probit Probability

 $\begin{aligned} \mathbf{u}_{i} \mid \boldsymbol{\varepsilon}_{i} \sim \mathsf{N}[(\rho / \sigma)\boldsymbol{\varepsilon}_{i}, (1 - \rho^{2})] \\ \mathsf{Prob}[\mathsf{d}_{i} = 1 \mid \mathbf{z}_{i}, \boldsymbol{\varepsilon}_{i}] &= \Phi\left[\frac{\mathbf{z}_{i}'\gamma + (\rho / \sigma)\boldsymbol{\varepsilon}_{i}}{\sqrt{1 - \rho^{2}}}\right] \\ \mathsf{Prob}[\mathsf{d}_{i} = 0 \mid \mathbf{z}_{i}, \boldsymbol{\varepsilon}_{i}] &= \Phi\left[\frac{-\mathbf{z}_{i}'\gamma - (\rho / \sigma)\boldsymbol{\varepsilon}_{i}}{\sqrt{1 - \rho^{2}}}\right] \end{aligned}$

Conditional Contribution to Likelihood $L_i(y_i, d_i = 1) | \epsilon_i = [f(y_i | \mathbf{x}_i, \epsilon_i, d_i = 1) \operatorname{Prob}[d_i = 1 | \mathbf{z}_i, \epsilon_i]$ $L_i(d_i = 0) = \operatorname{Prob}[d_i = 0 | \mathbf{z}_i, \epsilon_i]$

19-29/39

Dear Professor Greene,

I am doing a project investigating the impact of hedge fund manager's coinvestment on the survival probability of the fund. As fund managers' coinvestment decision is self-selection which might cause endogeneity issue, I jointly estimate the co-investment decision (Probit model) and the survival probability (Hazard model) to account for endogeneity of co-investment decision. I received one comment saying that I should use Heckman's two procedure to correct for endogeneity. My understanding is the Heckman's approach applies to a Probit and a LINEAR model. Since hazard model is nonlinear, simply adding inverse Mill's ration in the hazard model is wrong.

What I am asking is if my understanding of this is correct? If so, why can we not simply add Mill's ratio in a nonlinear model?

Conditional Likelihood

Conditional Density (not the log) $f(y_i, d_i = 1 | \epsilon_i) = [f(y_i | \epsilon_i, d_i = 1)]Prob[d_i = 1 | \epsilon_i]$ $f(y_i, d_i = 0 | \epsilon_i) = Prob[d_i = 0 | \epsilon_i]$ Unconditional Densities

$$f(y_{i}, d_{i} = 1) = \int_{-\infty}^{\infty} [f(y_{i} | \varepsilon_{i}, d_{i} = 1)] \operatorname{Prob}[d_{i} = 1 | \varepsilon_{i}] \frac{1}{\sigma} \phi\left(\frac{\varepsilon}{\sigma}\right) d\varepsilon$$
$$f(y_{i}, d_{i} = 0) = \int_{-\infty}^{\infty} \operatorname{Prob}[d_{i} = 0 | \varepsilon_{i}] \frac{1}{\sigma} \phi\left(\frac{\varepsilon}{\sigma}\right) d\varepsilon$$
$$Log Likelihoods$$

 $logL_i = logt(y_i, d_i)$

19-31/39

Poisson Model with Selection

Strategy:

- Hermite quadrature or maximum simulated likelihood.
- Not by throwing a 'lambda' into the unconditional likelihood
- Could this be done without joint normality?
 - How robust is the model?
 - Is there any other approach available?
 - Not easily. The subject of ongoing research

Nonnormality Issue

- How robust is the Heckman model to nonnormality of the unobserved effects?
- Are there other techniques
 - Parametric: Copula methods
 - Semiparametric: Klein/Spady and Series methods
- Other forms of the selection equation e.g., multinomial logit
- Other forms of the primary model: e.g., as above.

Application: Health Care Usage

German Health Care Usage Data, 7,293 Individuals, Varying Numbers of Periods

This is an unbalanced panel with 7,293 individuals. There are altogether 27,326 observations. The number of observations ranges from 1 to 7.

(Frequencies are: 1=1525, 2=2158, 3=825, 4=926, 5=1051, 6=1000, 7=987).

(Downloaded from the JAE Archive)

Variables in the file are

DOCTOR	= 1(Number of doctor visits > 0)	
HOSPITAL	= 1(Number of hospital visits > 0)	
HSAT	= health satisfaction, coded 0 (low) - 10 (high)	
DOCVIS	= number of doctor visits in last three months	
HOSPVIS	= number of hospital visits in last calendar year	
PUBLIC	= insured in public health insurance = 1; otherwise = 0	
ADDON	= insured by add-on insurance = 1; otherswise = 0	
HHNINC	= household nominal monthly net income in German marks /	10000.
	(4 observations with income=0 were dropped)	
HHKIDS	= children under age 16 in the household = 1; otherwise = 0	
EDUC	= years of schooling	
AGE	= age in years	
MARRIED	= marital status	

```
SAMPLE : All $
NAMELIST : z = one,age,educ,married,hhkids,hhninc $
NAMELIST : x = one,age,hsat $
PROBIT : Lhs = public : Rhs = z : Hold $
SELECT : Lhs = docvis : Rhs = x : mle$
POISSON : Lhs = docvis : Rhs = x : Selection : MLE $
BIVARIATE: Lhs = doctor,public : Rh1=x : Rh2=z : Selection $
```

Binomial Dependent Log like Restricte Chi squar Significa McFadden Estimatic Inf.Cr.Al Results r	Probit Model t variable lihood function ed log likelihood red [5 d.f.] ance level Pseudo R-squared on based on N = 2 IC = 16652.6 AIC retained for SELEC	PUBL -8320.323 -9711.251 2781.855 .000 .14322 7326, K = ∕N = .6 TION model.	IC 99 53 08 00 85 6 09				
PUBLIC	Coefficient	Standard Error	z	Prob. z >Z *	95% Con Inte	nfidence erval	
Constant AGE EDUC MARRIED HHKIDS HHNINC	Index function fo 3.63081*** .00115 17193*** 02762 06940*** 97958***	r probabili .07341 .00111 .00406 .02903 .02503 .05581	ty 49.46 1.03 -42.30 95 -2.77 -17.55	.0000 .3011 .0000 .3413 .0056 .0000	3.48693 00103 17990 08452 11845 -1.08895	3.77469 .00333 16397 .02927 02035 87020	

Sample Probit Selecti Results Data se Selecte	Selection Model selection equati ion rule is: Obse s of selection: Data et 27 ed sample 24	on based on rvations wit points S 326 203	PUBLIC h PUBLIC Sum of we: 27326 24203	= 1 ights .0 .0	1	-	
Sample Se Two step LHS=DOCVI Model siz Residuals Fit Model tes Not using Standard Correlati and Selec	election Model least square IS Mean Standard dev Number of ob te Parameters Degrees of f s Sum of squar Standard err R-squared Adjusted R-s st F[3, 24199 g OLS or no const error corrected ion of disturbanc ction Criterion (s regression = iation = servs. = reedom = es = or of e = quared =] (prob) = ant. Rsqrd & for selectic e in regress Rho) =	3.3 5.8 719 5.4 1320.7(.0 F may be on 5.4 sion1	 31207 38224 24203 4 24199 9470. 45265 14069 14069 14059 0000) ≥ < 0 47621 18222			
DOCVIS	Coefficient	Standard Error	z	Prob z >Z*	. 95 *	% Confidence Interval	
Constant AGE HSAT LAMBDA	8.22501*** .02939*** 89510*** 99788***	.20329 .00319 .01564 .23717	40.46 9.22 -57.23 -4.21	.0000 .0000 .0000 .0000	7.82 .02 92 -1.46	2657 8.6234 314 .0356 25758644 2735330	15 53 15 13

Sample Selection Model Two step least squares regression Standard error corrected for selection 5.47621 Correlation of disturbance in regression and Selection Criterion (Rho) = - 18222								
	· · · · · · · · · · · · · · · · · · ·	Standard Prob 95% Confidence						
DOCVIS	Coefficient	Error	z	z >Z*	Interval			
Constant AGE HSAT LAMBDA	8.22501*** .02939*** 89510*** 99788***	.20329 .00319 .01564 .23717	40.46 9.22 -57.23 -4.21	.0000 .0000 .0000 .0000	7.82657 .02314 92575 -1.46273	8.62345 .03563 86445 53303		
ML Estimates of Selection Model Dependent variable DOCVIS Log likelihood function -83716.13744 Estimation based on N = 27326, K = 11 Inf.Cr.AIC = 167454.3 AIC/N = 6.128 Model estimated: Aug 01, 2012, 23:04:48								
DOCVIS	9 Coefficient	Standard Error	z	Prob. z >Z *	95% Confidence Interval			
Constant AGE EDUC MARRIED HHKIDS HHNINC	Selection (probit) 3.63881*** .00117 17218*** 02901 07355*** 98832*** Corrected regressio	equation .07361 .00115 .00412 .02925 .02512 .05164 on, Regime .0202	for PUBL: 49.44 1.02 -41.83 99 -2.93 -19.14 1 41.70	IC .0000 .3073 .0000 .3212 .0034 .0034	3.49455 00108 18024 08633 12279 -1.08953	3.78308 .00342 16411 .02831 02431 88712		
Constant AGE HSAT SIGMA(1) RHO(1,2)	8.10465*** .03040*** 89792*** 5.45984*** 08970*	.00325 .01389 .01019 .04810	41.79 9.35 -64.65 535.78 -1.86	.0000 .0000 .0000 .0000 .0622	7.72457 .02403 92514 5.43987 18398	8.48473 .03678 87070 5.47982 .00458		

Poisson Model with Sample Selection. Dependent variable DOCVIS Log likelihood function -60829.53023 Restricted log likelihood -205953.60785 Chi squared [2 d.f.] 290248.15524 Significance level .00000 McFadden Pseudo R-squared .7046445 Estimation based on N = 27326, K = 11 Inf.Cr.AIC = 121681.1 AIC/N = 4.453 Restr. Log-L is Poisson+Probit (indep). LogL for initial probit = -8320.3674 LogL for initial Poisson= -197633.2404 Means for Psn/Neg.Bin. use selected data. Means for Probit based on all observations.							
DOCVIS	Coefficient	Standard Error	z	Prob. z >Z *	95% Confidence Interval		
Constant AGE HSAT	Parameters of Pois 1.55949*** .01053*** 25228*** Parameters of Prob	sson/Neg. B .04592 .00078 .00339 bit Selectio	inomial H 33.96 13.50 -74.50 on Model	Probabil: .0000 .0000 .0000	ity 1.46948 .00900 25892	1.64949 .01206 24564	
Constant AGE EDUC MARRIED HHKIDS HHNINC	3.55923*** .00175 16992*** 02874 06369** 93895***	.07381 .00115 .00412 .02937 .02518 .05262	48.22 1.52 -41.21 98 -2.53 -17.84	.0000 .1285 .0000 .3277 .0114 .0000	3.41456 00051 17800 08631 11305 -1.04208	3.70390 .00401 16184 .02882 01434 83582	
Sigma Rho	Standard Deviation 1.16922*** Correlation of He 02853	n of Hetero .00768 terogeneity .05534	geneity 152.19 & Select 52	.0000 tion .6062	1.15416 13698	1.18427 .07993	

FIML Estimates of Bivariate Probit Model Dependent variable DOCPUB Log likelihood function -22945.59406 Estimation based on N = 27326, K = 10 Inf.Cr.AIC = 45911.2 AIC/N = 1.680 Selection model based on PUBLIC Selected obs. 24203, Nonselected: 3123							
DOCTOR PUBLIC	Coefficient	Standard Error	z	Prob. z >Z *	95% Confidence Interval		
Constant AGE HSAT	Index equation 1.18664*** .00892*** 17226*** Index equation	for DOCTOR .05011 .00079 .00416 for PUPLIC	23.68 11.29 -41.43	. 0000 . 0000 . 0000	1.08843 .00737 18041	1.28485 .01047 16411	
Constant AGE EDUC MARRIED HHKIDS HHNINC RHO(1,2)	3.63944*** .00098 17210*** 02446 07540*** 97675*** Disturbance correl 13237**	.07354 .00115 .00411 .02927 .02511 .05162 lation .05753	49.49 .85 -41.86 84 -3.00 -18.92 -2.30	.0000 .3949 .0000 .4033 .0027 .0027 .0000	3.49530 00127 18016 08182 12462 -1.07792 24513	3.78357 .00323 16405 .03290 02619 87558 01960	