

Topics in Microeconometrics

William Greene Department of Economics Stern School of Business

[Topic 7-Selection] 1/81









WILLIAM H. GREENE

Part 7: Sample Selection in Nonlinear and Panel Models







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. These may not be the same as the full population
- Can one generalize from the subpopulation to the full population?

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. ^(C))
- 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."

Heckman's Canonical Model

ALC IN COLUMN

A behavioral model:

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

LFP = 1 if $o^* \ge r^*$, 0 otherwise Prob(LFP=1) = $\Phi\left[(\mathbf{P} - \mathbf{v})/\sqrt{\sigma_v^2 + \sigma_u^2}\right]$ Desired Hours = $H^* = \gamma'\mathbf{W} + \varepsilon$ Actual Hours = H^* if LFP = 1 unobserved if LFP = 0 ε and u are correlated. ε and v might be correlated.

What is $E[H^* | \mathbf{w}, LFP = 1]$? Not $\gamma' \mathbf{w}$.

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>

Sample Selection Observations

Str. Mallin Str

- The selection 'problem' is caused by the correlation of the unobservables
 - Selection on observables is often manageable within the conventional model.
 - Selection on unobservables often requires a more detailed specification of the model – where does the effect come from?
- The 'bias' relates to the inconsistency of familiar estimators such as least squares
- The data are not biased; the (an) estimator is biased.

Standard Sample Selection Model

122 12

 $d_i^{\star} = \boldsymbol{\alpha}' \mathbf{Z}_i + U_i$ $d_{i} = 1(d_{i}^{*} > 0)$ $\gamma_i^{\star} = \boldsymbol{\beta}' \boldsymbol{X}_i + \varepsilon_i$ $y_i = y_i^*$ when $d_i = 1$, unobserved otherwise $(u_i,v_i) \sim \text{Bivariate Normal}[(0,0),(1,\rho\sigma,\sigma^2)]$ $E[y_i | y_i \text{ is observed}] = E[y_i | d_i = 1]$ $= \boldsymbol{\beta}' \mathbf{X}_{i} + \mathbf{E}[\varepsilon_{i} \mid \mathbf{d}_{i} = 1]$ $= \boldsymbol{\beta}' \mathbf{X}_{i} + \mathbf{E}[\varepsilon_{i} | \mathbf{u}_{i} > -\alpha' \mathbf{Z}_{i}]$ $= \boldsymbol{\beta}' \boldsymbol{x}_{i} + (\rho \sigma) \frac{\boldsymbol{\phi}(\boldsymbol{\alpha}' \boldsymbol{z}_{i})}{\boldsymbol{\Phi}(\boldsymbol{\alpha}' \boldsymbol{z}_{i})}$ $= \boldsymbol{\beta}' \mathbf{X}_{i} + \theta \lambda_{i}$

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



Selection as a Specification Error

- $E[y_i | \mathbf{x}_i, y_i \text{ observed}] = \mathbf{\beta}' \mathbf{x}_i + \Theta \lambda_i$
- 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.
 - \mathbf{z}_i and \mathbf{x}_i often have variables in common.
 - There is no specification error if $\theta = 0 \leftrightarrow \rho = 0$
- The "selection bias" is plim (b β)

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. Not more or less robust
- Simple to do 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).

Mroz 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	= wife's mother's education
WFED	= wife's father's education
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.

AND DECK

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 ; Lhs = WHRS ; Rhs = X,Lambda \$ REJECT ; LFP = 0 \$ REGRESS ; Lhs = WHRS ; Rhs = X \$

Participation Equation

Binomial Dependen Weighting Number o:	Probit Model t variable g variable f observations	LFP None 753	+		
+ Variable	++ Coefficient	Standard Error	+ b/St.Er.	+ P[Z >z]	++ Mean of X
	Index function f	or probability	•		· ·
Constant	1.00264501	.49994379	2.006	.0449	
KL6	90399802	.11434394	-7.906	.0000	.23771580
K618	05452607	.04021041	-1.356	.1751	1.35325365
WA	02602427	.01332588	-1.953	.0508	42.5378486
WE	.16038929	.02773622	5.783	.0000	12.2868526
HA	01642514	.01329110	-1.236	.2165	45.1208499
HE	05191039	.02040378	-2.544	.0110	12.4913679

Hours Equation

Sample Se	election Model					
Two stage	e least square	least squares regression				
LHS=WHRS	Mean	:	= 13	302.930		
Ì	Standard de	viation	= 7'	76.2744	Ì	
WTS=none	Number of o	bservs.	=	428	Ì	
Model siz	e Parameters	:	=	8	İ	
i	Degrees of	freedom	=	420	i	
 Residuals	Sum of squa:	res	= .:	2267214E+0	9 İ	
i	Standard er	ror of e	= 7.	34.7195	i	
 Correlati	on of disturban	ce in regr	essio	n	i	
and Selec	tion Criterion	(Rho)		8454	1	
+					+	
++	+			-+	+	.++
Variable	Coefficient	Standard 3	Error	b/St.Er.	P[Z >z]	Mean of X
Constant	2442.26665	1202.1	 1143	2.032	.0422	++
KL6	115.109657	282.00	3565	.408	.6831	.14018692
K618	-101.720762	38.283	3942	-2.657	.0079	1.35046729
AGE	14.6359451	53.191	5591	.275	.7832	41.9719626
AGESO	10078602	.6185	5252	163	.8706	1821.12150
WE	-102.203059	39.409	5323	-2.593	.0095	12.6588785
FAMINC	.01379467	.0034	5041	3.998	.0001	24130.4229
LAMBDA	-793.857053	494.54	1008	-1.605	.1084	.61466207

Selection "Bias" of OLS

_	L J					
	Variable	Coefficient	Standard Error	b/St.Er.	P[Z >z]	Mean of X
	Constant	2442.26665	1202.11143	2.032	.0422	-++
	KL6	115.109657	282.008565	.408	.6831	.14018692
	K618	-101.720762	38.2833942	-2.657	.0079	1.35046729
	AGE	14.6359451	53.1916591	.275	.7832	41.9719626
	AGESQ	10078602	.61856252	163	.8706	1821.12150
	WE	-102.203059	39.4096323	-2.593	.0095	12.6588785
	FAMINC	.01379467	.00345041	3.998	.0001	24130.4229
	LAMBDA	-793.857053	494.541008	-1.605	.1084	.61466207
-	+ ++ Variable	Coefficient	Standard Error	t-ratio	P[T >t]	Mean of X
-	Variable Constant	Coefficient 1812.12538	Standard Error 1144.33342	t-ratio 1.584	P[T >t] .1140	Mean of X
-	Variable Constant KL6	Coefficient 1812.12538 -299.128041	Standard Error 1144.33342 100.033124	t-ratio t-ratio 1.584 -2.990	P[T >t] .1140 .0030	.14018692
-	Variable Constant KL6 K618	Coefficient 1812.12538 -299.128041 -126.399697	Standard Error 1144.33342 100.033124 30.8728451	t-ratio t-s84 -2.990 -4.094	<pre>.1140 .0030 .0001</pre>	.14018692 1.35046729
-	Variable Constant KL6 K618 AGE	Coefficient 1812.12538 -299.128041 -126.399697 11.2795338	Standard Error 1144.33342 100.033124 30.8728451 53.8442084	t-ratio 1.584 -2.990 -4.094 .209	<pre>.1140 .0030 .0001 .8342</pre>	.14018692 1.35046729 41.9719626
-	Variable Constant KL6 K618 AGE AGESQ	Coefficient 1812.12538 -299.128041 -126.399697 11.2795338 26103541	Standard Error 1144.33342 100.033124 30.8728451 53.8442084 .62632815	t-ratio 1.584 -2.990 -4.094 .209 417	<pre> P[T >t] .1140 .0030 .0001 .8342 .6771</pre>	Mean of X .14018692 1.35046729 41.9719626 1821.12150
-	Variable Constant KL6 K618 AGE AGESQ WE	Coefficient 1812.12538 -299.128041 -126.399697 11.2795338 26103541 -47.3271780	Standard Error 1144.33342 100.033124 30.8728451 53.8442084 .62632815 17.2968137	t-ratio 1.584 -2.990 -4.094 .209 417 -2.736	<pre> P[T >t] .1140 .0030 .0001 .8342 .6771 .0065</pre>	Mean of X .14018692 1.35046729 41.9719626 1821.12150 12.6588785
-	Variable Constant KL6 K618 AGE AGESQ WE FAMINC	Coefficient 1812.12538 -299.128041 -126.399697 11.2795338 26103541 -47.3271780 .01261889	Standard Error 1144.33342 100.033124 30.8728451 53.8442084 .62632815 17.2968137 .00338906	l.584 -2.990 -4.094 .209 417 -2.736 3.723	<pre> P[T >t] .1140 .0030 .0001 .8342 .6771 .0065 .0002</pre>	.14018692 1.35046729 41.9719626 1821.12150 12.6588785 24130.4229

Maximum Likelihood Estimation

$$\begin{split} &\log L = \sum_{d=1} \log \left[\frac{\exp\left(-\frac{1}{2} \left(\epsilon_{i} \mid \sigma\right)^{2}\right)}{\sigma \sqrt{2\pi}} \Phi\left(\frac{\rho(\epsilon_{i} \mid \sigma) + \boldsymbol{\alpha} \mid \boldsymbol{z}_{i}}{\sqrt{1 - \rho^{2}}}\right) \right] \\ &+ \sum_{d=0} \log \left[1 - \Phi(\boldsymbol{\alpha} \mid \boldsymbol{z}_{i}) \right] \\ &\text{Re parameterize this: let } \boldsymbol{q}_{i} = \boldsymbol{\alpha} \mid \boldsymbol{z}_{i} \\ &(1) \mid \theta = 1/\sigma \\ &(2) \mid \gamma = \beta/\sigma \text{ (Olsen transformation)} \\ &(3) \mid \tau = \rho/\sqrt{1 - \rho^{2}} \\ &(4) \text{ Constrain } \rho \text{ to be in (-1,1) by using} \\ &\psi = \frac{1}{2} ln\left(\frac{1 + \rho}{1 - \rho}\right) = a \tanh \rho, \text{ so } \rho = a \tanh^{-1}(\psi) = \frac{e x p(2\psi) - 1}{e x p(2\psi) + 1} \\ &\log L = \sum_{d=0} log \Phi(-\boldsymbol{q}_{i}) + \sum_{d=1}^{log} \frac{\log \theta - \frac{1}{2} log 2\pi - \frac{1}{2} (\theta y_{i} - \gamma \mid \boldsymbol{x}_{i})^{2}}{+ \log \Phi[\tau(\theta y_{i} - \gamma \mid \boldsymbol{x}_{i}) + \boldsymbol{q}_{i}\sqrt{1 + \tau^{2}}} \end{split}$$

[Topic 7-Selection] 18/81

ML Estima Maximum I Number of Iteratior Log like] Number of FIRST 7	ates of Selection Likelihood Estim observations as completed Lihood function parameters estimates are p	on Model nates 753 47 -3894.471 16 probit equation.		
Variable	Coefficient	Standard Error	b/St.Er.	P[Z >z]
۰۱ ٤	Selection (probi	it) equation for	LFP	·
Constant	1.01350651	.54823177	1.849	.0645
KL6	90129694	.11081111	-8.134	.0000
K618	05292375	.04137216	-1.279	.2008
WA	02491779	.01428642	-1.744	.0811
WE	.16396194	.02911763	5.631	.0000
HA	01763340	.01431873	-1.231	.2181
HE	05596671	.02133647	-2.623	.0087
C	Corrected regres	ssion, Regime 1		
Constant	1946.84517	1167.56008	1.667	.0954
KL6	-209.024866	222.027462	941	.3465
K618	-120.969192	35.4425577	-3.413	.0006
AGE	12.0375636	51.9850307	.232	.8169
AGESQ	22652298	.59912775	378	.7054
WE	-59.2166488	33.3802882	-1.774	.0761
FAMINC	.01289491	.00332219	3.881	.0001
SIGMA(1)	748.131644	59.7508375	12.521	.0000
RHO(1,2)	22965163	.50082203	459	.6466

MLE

100

MLE vs. Two Step

Two Stan

Service and the service of the servi

2442.26665	1202.11143	2.032	.0422	
115.109657	282.008565	.408	.6831	.14018692
-101.720762	38.2833942	-2.657	.0079	1.35046729
14.6359451	53.1916591	.275	.7832	41.9719626
10078602	.61856252	163	.8706	1821.12150
-102.203059	39.4096323	-2.593	.0095	12.6588785
.01379467	.00345041	3.998	.0001	24130.4229
-793.857053	494.541008	-1.605	.1084	.61466207
Standard erro	or of $e = 73$	4.7195		
on of disturbance	e in regression			
ion Criterion (H	Rho)	84541	L	
1946.84517	1167.56008	1.667	.0954	
-209.024866	222.027462	941	.3465	
-120.969192	35.4425577	-3.413	.0006	
12.0375636	51.9850307	.232	.8169	
22652298	.59912775	378	.7054	
-59.2166488	33.3802882	-1.774	.0761	
.01289491	.00332219	3.881	.0001	
748.131644	59.7508375	12.521	.0000	
22965163	.50082203	459	.6466	
	2442.26665 115.109657 -101.720762 14.6359451 10078602 -102.203059 .01379467 -793.857053 Standard erro on of disturbance ion Criterion (F 1946.84517 -209.024866 -120.969192 12.0375636 22652298 -59.2166488 .01289491 748.131644 22965163	2442.26665 1202.11143 115.109657 282.008565 -101.720762 38.2833942 14.6359451 53.1916591 10078602 .61856252 -102.203059 39.4096323 .01379467 .00345041 -793.857053 494.541008 Standard error of e = 73 on of disturbance in regression ion Criterion (Rho) 1946.84517 1167.56008 -209.024866 222.027462 -120.969192 35.4425577 12.0375636 51.9850307 22652298 .59912775 -59.2166488 33.3802882 .01289491 .00332219 748.131644 59.7508375 22965163 .50082203	2442.26665 1202.11143 2.032 115.109657 282.008565 .408 -101.720762 38.2833942 -2.657 14.6359451 53.1916591 .275 10078602 .61856252163 -102.203059 39.4096323 -2.593 .01379467 .00345041 3.998 -793.857053 494.541008 -1.605 Standard error of e = 734.7195 on of disturbance in regression ion Criterion (Rho)84541 1946.84517 1167.56008 1.667 -209.024866 222.027462941 -120.969192 35.4425577 -3.413 12.0375636 51.9850307 .232 22652298 .59912775378 -59.2166488 33.3802882 -1.774 .01289491 .00332219 3.881 748.131644 59.7508375 12.521 22965163 .50082203459	2442.26665 1202.11143 2.032 .0422 115.109657 282.008565 .408 .6831 -101.720762 38.2833942 -2.657 .0079 14.6359451 53.1916591 .275 .7832 10078602 .61856252163 .8706 -102.203059 39.4096323 -2.593 .0095 .01379467 .00345041 3.998 .0001 -793.857053 494.541008 -1.605 .1084 Standard error of e = 734.7195 on of disturbance in regression ion Criterion (Rho)84541 1946.84517 1167.56008 1.667 .0954 -209.024866 222.027462941 .3465 -120.969192 35.4425577 -3.413 .0006 12.0375636 51.9850307 .232 .8169 22652298 .59912775378 .7054 -59.2166488 33.3802882 -1.774 .0761 .01289491 .00332219 3.881 .0001 748.131644 59.7508375 12.521 .0000 22965163 .50082203459 .6466

Extension - Treatment Effect

1 38 TH

Str. Martin Sta

What is the value of an elite college education? $d_i^* = \mathbf{z} \mathbf{v} + u_i; d_i = 1[d_i^* > 0]$ (probit) $y_i^* = \mathbf{x} \mathbf{\beta} + \delta d_i + \varepsilon_i$ observed for everyone $[\varepsilon_i, u_i] \sim \text{Bivariate Normal}[0, 0, \sigma^2, \rho, 1]$ $E[\mathbf{y}_i^* | \mathbf{x}_i, \mathbf{d}_i = 1] = \mathbf{x}\mathbf{\beta} + \delta + E[\varepsilon_i | \mathbf{x}_i, \mathbf{d}_i = 1]$ $= \mathbf{x}\mathbf{\beta} + \delta + E[\varepsilon_i | \mathbf{x}_i, \mathbf{u}_i > -\mathbf{z} \mathbf{y}]$ $= \mathbf{X}\mathbf{\beta} + \delta + \rho\sigma\left(\frac{\phi(\mathbf{z}\mathbf{\dot{\gamma}})}{\Phi(\mathbf{z}\mathbf{\dot{\gamma}})}\right)$ $= \mathbf{X}\mathbf{\beta} + \delta + \rho\sigma\lambda_{i}$ $E[y_{i} * | x_{i}, d_{i} = 0] = \mathbf{x} \mathbf{\beta} + E[\varepsilon_{i} | x_{i}, d_{i} = 0]$ $= \mathbf{x}\mathbf{\beta} + \rho\sigma\left(\frac{-\phi(-\mathbf{z}\mathbf{\dot{\gamma}})}{\Phi(-\mathbf{z}\mathbf{\dot{\gamma}})}\right)$

Least squares is still biased and inconsistent. Left out variable [Topic 7-Selection] 21/81



SPIRIT PROPERTY

$$E[\mathbf{y}_{i}^{*}|\mathbf{x}_{i},\mathbf{d}_{i}=1] - E[\mathbf{y}_{i}^{*}|\mathbf{x}_{i},\mathbf{d}_{i}=0]$$

$$= \mathbf{x}\mathbf{\beta} + \delta + \rho\sigma\left(\frac{\phi(\mathbf{z}\mathbf{\dot{\gamma}})}{\Phi(\mathbf{z}\mathbf{\dot{\gamma}})}\right)$$

$$- \mathbf{x}\mathbf{\beta} - \rho\sigma\left(\frac{-\phi(-\mathbf{z}\mathbf{\dot{\gamma}})}{\Phi(-\mathbf{z}\mathbf{\dot{\gamma}})}\right)$$

$$= \delta + \rho\sigma\left[\left(\frac{\phi(\mathbf{z}\mathbf{\dot{\gamma}})}{\Phi(\mathbf{z}\mathbf{\dot{\gamma}})}\right) - \left(\frac{-\phi(-\mathbf{z}\mathbf{\dot{\gamma}})}{\Phi(-\mathbf{z}\mathbf{\dot{\gamma}})}\right)\right]$$

$$= \text{Treatment} + \text{Selection Effect}$$

Sample Selection in Exponential Regression

An approach modeled on Heckman's model

Regression Equation:
$$Prob[y=j|x,u]=P(\lambda)$$
;
 $\lambda=exp(x\beta + \theta u)$
Selection Equation: $d=1[z\delta > 0]$ (The usual probit)

 $[u,\varepsilon] \sim n[0,0,1,1,\rho]$ (Var[u] is absorbed in θ)

Estimation: Nonlinear Least Squares: [Terza (1998).] $E[y|x,d=1]=exp(x\beta\rho+)^{-2} \frac{\Phi(z\delta+\rho)}{\Phi(z\delta)}$

Panel Data and Selection

S.S.M. Store

Selection equation with time invariant individual effect

$$\mathbf{d}_{it} = \mathbf{1}[\mathbf{z}\mathbf{\dot{y}} + \theta_i + \eta_{it} > 0]$$

Observation mechanism: $(y_{it}, \mathbf{x}_{it})$ observed when $d_{it} = 1$

Primary equation of interest

Common effects linear regression model

 $y_{it} \mid (d_{it} = 1) = \textbf{x} \textbf{\beta} + \alpha_i + \epsilon_{it}$

"Selectivity" as usual arises as a problem when the unobservables are correlated; Corr(ϵ_{it} , η_{it}) $\neq 0$.

The common effects, θ_i and α_i make matters worse.

Panel Data and Sample Selection Models: A Nonlinear Time Series

- I. 1990-1992: Fixed and Random Effects Extensions
- II. 1995 and 2005: Model Identification through Conditional Mean Assumptions
- III. 1997-2005: Semiparametric Approaches based on Differences and Kernel Weights
- IV. 2007: Return to Conventional Estimators, with Bias Corrections

Panel Data Sample Selection Models

No. Martin

Verbeek, Economics Letters, 1990. $d_{it} = 1[z_{ij} + w_i + \eta_{it} > 0]$ (Random effects probit) $y_{it} | (d_{it} = 1) = x\beta + \alpha_i + \varepsilon_{it}$; (Fixed effects regression) Proposed "marginal likelihood" based on joint normality

$$\begin{split} \log \mathbf{L}_{i} &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \prod_{t=1}^{\mathsf{T}_{i}} \Phi \left[(2\mathsf{d}_{it} - 1) \frac{\mathbf{z} \mathbf{\hat{y}} + \Delta_{it} + \mathsf{u}_{i,1} + \mathsf{d}_{it} \mathsf{u}_{i,2}}{\sqrt{\sigma_{\eta}^{2} (1 - \mathsf{d}_{it} \rho^{2})}} \right] \mathsf{f}(\mathsf{u}_{i,1}, \mathsf{u}_{i,2}) \mathsf{d}\mathsf{u}_{i,1} \mathsf{d}\mathsf{u}_{i,2} \\ \Delta_{it} &= (\rho / \sigma_{\varepsilon}) \mathsf{d}_{it} \left[(\mathsf{y}_{it} - \overline{\mathsf{y}}_{i}) - (\mathsf{x}_{it} \mathbf{\beta} \overline{\mathsf{x}}_{i})^{\mathsf{r}} \right] \end{split}$$

(Integrate out the random effects; difference out the fixed effects.) u_{i,1}, u_{i,2} are time invariant uncorrelated standard normal variables How to do the integration? Natural candidate for simulation. (Not mentioned in the paper. Too early.) [Verbeek and Nijman: Selectivity "test" based on this model, International Economic Review, 1992.]

Zabel – Economics Letters

AL MALLIN

- Inappropriate to have a mix of FE and RE models
- Two part solution
 - Treat both effects as "fixed"
 - Project both effects onto the group means of the variables in the equations (Mundlak approach)
 - Resulting model is two random effects equations
- Use both random effects

Selection with Fixed Effects

$$y_{it}^{*} = \eta_{i} + \mathbf{x}_{it}^{'} \boldsymbol{\beta} + \varepsilon_{it}, \ \eta_{i} = \overline{\mathbf{x}}_{i}^{'} \boldsymbol{\pi} + \tau w_{i}, w_{i} \sim N[0,1]$$

$$d_{it}^{*} = \theta_{i} + \mathbf{z}_{it}^{'} \boldsymbol{\alpha} + u_{it}, \ \theta_{i} = \overline{\mathbf{z}}_{i}^{'} \boldsymbol{\delta} + \omega v_{i}, v_{i} \sim N[0,1]$$

$$(\varepsilon_{it}, u_{it}) \sim N_{2}[(0,0), (\sigma^{2}, 1, \rho\sigma)].$$

$$L_{i} = \int_{-\infty}^{\infty} \prod_{d_{it}=0} \Phi[-\mathbf{z}_{it}^{'} \boldsymbol{\alpha} - \overline{\mathbf{z}}_{i}^{'} \boldsymbol{\delta} - \omega v_{i}] \phi(v_{i}) dv_{i}$$

$$\times \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \prod_{d_{it}=1} \left(\Phi\left[\frac{\mathbf{z}_{it}^{'} \boldsymbol{\alpha} + \overline{\mathbf{z}}_{i}^{'} \boldsymbol{\delta} + \omega v_{i} + (\rho / \sigma) \varepsilon_{it}}{\sqrt{1 - \rho^{2}}}\right] dv_{i} dw_{i}$$

$$\varepsilon_{it} = y_{it} - \mathbf{x}_{it}^{'} \boldsymbol{\beta} - \overline{\mathbf{x}}_{i}^{'} \boldsymbol{\pi} - \tau w_{i}$$

Practical Complications

The bivariate normal integration is actually the product of two univariate normals, because in the specification above, v_i and w_i are assumed to be uncorrelated. Vella notes, however, "... given the computational demands of estimating by maximum likelihood induced by the requirement to evaluate multiple integrals, we consider the applicability of available simple, or two step procedures."



Simulation

The first line in the log likelihood is of the form $E_v[\Pi_{d=0}\Phi(...)]$ and the second line is of the form $E_w[E_v[\Phi(...)\phi(...)/\sigma]]$. Using simulation instead, the simulated likelihood is

$$L_{i}^{S} = \frac{1}{R} \sum_{r=1}^{R} \prod_{d_{it}=0} \Phi \left[-\mathbf{z}_{it}' \boldsymbol{\alpha} - \overline{\mathbf{z}}_{i}' \boldsymbol{\delta} - \omega v_{i,r} \right]$$

$$\times \frac{1}{R} \sum_{r=1}^{R} \prod_{d_{it}=1} \Phi \left[\frac{\mathbf{z}_{it}' \boldsymbol{\alpha} + \overline{\mathbf{z}}_{i}' \boldsymbol{\delta} + \omega v_{i,r} + (\rho / \sigma) \varepsilon_{it,r}}{\sqrt{1 - \rho^{2}}} \right] \frac{1}{\sigma} \phi \left(\frac{\varepsilon_{it,r}}{\sigma} \right)$$

$$\varepsilon_{it,r} = y_{it} - \mathbf{x}_{it}' \boldsymbol{\beta} - \overline{\mathbf{x}}_{i}' \boldsymbol{\pi} - \tau w_{i,r}$$

Correlated Effects

Suppose that w_i and v_i are bivariate standard normal with correlation ρ_{vw} . We can project w_i on v_i and write $w_i = \rho_{vw} v_i + (1 - \rho_{vw}^2)^{1/2} h_i$

where h_i has a standard normal distribution. To allow the correlation, we now simply substitute this expression for wi in the simulated (or original) log likelihood, and add ρ_{vw} to the list of parameters to be estimated. The simulation is then over still independent normal variates, v_i and h_i .



Conditional Means

Wooldridge (1995) proposes an estimator that can be based on straightforward applications of conventional, everyday methods.

$$y_{it}^{*} = \eta_{i} + \mathbf{x}_{it}^{\prime} \boldsymbol{\beta} + \varepsilon_{it},$$

$$d_{it}^{*} = \theta_{i} + \mathbf{z}_{it}^{\prime} \boldsymbol{\alpha} + u_{it},$$

$$(\varepsilon_{it}, u_{it}) \sim N_{2}[(0, 0), (\sigma^{2}, 1, \rho\sigma)].$$

Under the mean independence assumption

$$\mathbf{E}[\mathbf{\varepsilon}_{it} \mid \mathbf{\eta}_i, \mathbf{\theta}_i, \mathbf{z}_{i1}, \dots, \mathbf{z}_{iT}, \mathbf{v}_{i1}, \dots, \mathbf{v}_{iT}, \mathbf{d}_{i1}, \dots, \mathbf{d}_{iT}] = \rho u_{it}, \quad \mathbf{v}_{it} = \mathbf{\theta}_i + u_{it}$$

 $\mathbf{E}[v_{it}|\mathbf{x}_{i1},\dots,\mathbf{x}_{iT}, \mathbf{\eta}_i, \mathbf{\theta}_i, \mathbf{z}_{i1},\dots,\mathbf{z}_{iT}, \mathbf{v}_{i1},\dots,\mathbf{v}_{iT}, \mathbf{d}_{i1},\dots,\mathbf{d}_{iT}] = \mathbf{\eta}_i + \mathbf{x}_{it}'\mathbf{\beta} + \rho u_{it}.$

This suggests an approach to estimating the model parameters, however it requires computation of u_{it} . That would require estimation of θ_i which cannot be done, at least not consistently – and that precludes simple estimation of u_{it} .



A Feasible Estimator

To escape the dilemma, Wooldridge suggests Chamberlain's approach to the fixed effects model,

$$\Theta_i = f_0 + \mathbf{z}_{i1}'\mathbf{f}_1 + \mathbf{z}_{i2}'\mathbf{f}_2 + \ldots + \mathbf{z}_{iT}'\mathbf{f}_T + h_i.$$

With this substitution,

$$d_{it}^* = \mathbf{z}_{it}' \boldsymbol{\alpha} + \mathbf{f}_0 + \mathbf{z}_{i1}' \mathbf{f}_1 + \mathbf{z}_{i2}' \mathbf{f}_2 + \dots + \mathbf{z}_{iT}' \mathbf{f}_T + h_i + u_{it}$$

= $\mathbf{z}_{it}' \boldsymbol{\alpha} + \mathbf{f}_0 + \mathbf{z}_{i1}' \mathbf{f}_1 + \mathbf{z}_{i2}' \mathbf{f}_2 + \dots + \mathbf{z}_{iT}' \mathbf{f}_T + w_{it}$

where w_{it} is independent of \mathbf{z}_{it} , t = 1, ..., T. This now implies that

$$\mathbf{E}[\mathbf{y}_{it}|\mathbf{x}_{i1},\dots,\mathbf{x}_{iT}, \mathbf{\eta}_i, \mathbf{\theta}_i, \mathbf{z}_{i1},\dots,\mathbf{z}_{iT}, \mathbf{v}_{i1},\dots,\mathbf{v}_{iT}, d_{i1},\dots,d_{iT}] = \mathbf{\eta}_i + \mathbf{x}_{it}'\mathbf{\beta} + \mathbf{\rho}(w_{it} - h_i)$$

= $(\mathbf{\eta}_i - \mathbf{\rho}\mathbf{h}_i) + \mathbf{x}_{it}'\mathbf{\beta} + \mathbf{\rho}w_{it}.$

[Topic 7-Selection] 33/81



Estimation

To complete the estimation procedure, we now compute *T* cross sectional probit models (reestimating $f_0, f_1,...$ each time), and compute $\hat{\lambda}_{it}$ from each one. The resulting equation,

$$y_{it} = a_i + \mathbf{x}_{it}' \mathbf{\beta} + \rho \,\hat{\lambda}_{it} + \mathbf{v}_{it}$$

now forms the basis for estimation of β and ρ by using a conventional fixed effects linear regression with the observed data.

Kyriazidou - Semiparametrics

Assume 2 periods

Estimate selection equation by FE logit

Use first differences and weighted least squares:

$$\hat{\boldsymbol{\beta}} = \left[\boldsymbol{\Sigma}_{i=1}^{N} \boldsymbol{d}_{i1} \boldsymbol{d}_{i2} \boldsymbol{\Psi}_{i} \Delta \boldsymbol{x}_{i} \Delta \boldsymbol{x}_{i}' \right]^{-1} \left[\boldsymbol{\Sigma}_{i=1}^{N} \boldsymbol{d}_{i1} \boldsymbol{d}_{i2} \hat{\boldsymbol{\Psi}}_{i} \Delta \boldsymbol{x}_{i} \Delta \boldsymbol{y}_{i} \right]$$
$$\hat{\boldsymbol{\Psi}}_{i} = \frac{1}{h} \mathsf{K} \left[\frac{\Delta \boldsymbol{w}_{i}' \hat{\boldsymbol{\alpha}}}{h} \right] \text{ kernel function.}$$

Use with longer panels - any pairwise differences Extensions based on pairwise differences by Rochina-Barrachina and Dustman/Rochina-Barrachina (1999)

Bias Corrections

- Val and Vella, 2007 (Working paper)
- Assume fixed effects
 - Bias corrected probit estimator at the first step
 - Use fixed probit model to set up second step Heckman style regression treatment.


Postscript

- What selection process is at work?
 - All of the work examined here (and in the literature) assumes the selection operates anew in each period
 - An alternative scenario: Selection into the panel, once, at baseline.
 - Alternative: Sequential selection = endogenous attrition (Wooldridge 2002, inverse probability weighting)
- Why aren't the time invariant components correlated?
- Other models
 - All of the work on panel data selection assumes the main equation is a linear model.
 - Any others? Discrete choice? Counts?

Attrition

- In a panel, t=1,...,T individual I leaves the sample at time K_i and does not return.
- If the determinants of attrition (especially the unobservables) are correlated with the variables in the equation of interest, then the now familiar problem of sample selection arises.

Dealing with Attrition in a QOL Study

St. Martin

- The attrition issue: Appearance for the second interview was low for people with initial low QOL (death or depression) or with initial high QOL (don't need the treatment). Thus, missing data at exit were clearly related to values of the dependent variable.
- Solutions to the attrition problem
 - Heckman selection model (used in the study)
 - **D** Prob[Present at exit|covariates] = $\Phi(z'\theta)$ (Probit model)
 - Additional variable added to difference model $\lambda_i = \Phi(\mathbf{z}_i \boldsymbol{\theta}) / \Phi(\mathbf{z}_i \boldsymbol{\theta})$
 - The FDA solution: fill with zeros. (!)

An Early Attrition Model

Hausman, J. and Wise, D., "Attrition Bias in Experimental and Panel Data: The Gary Income Maintenance Experiment," Econometrica, 1979. A two period model:

Structural response model (Random Effects Regression)

$$y_{i1} = \mathbf{x}\mathbf{\beta} + \varepsilon_{i1} + u_i$$

$$y_{i2} = \mathbf{x}\mathbf{\beta} + \varepsilon_{i2} + u_i$$

Attrition model for observation in the second period (Probit)

$$z_{i2}^{*} = \delta y_{i2} + x \Theta + v_{i2}$$

 $z_{i2} = 1(z_{i2}^{*} > 0)$

Endogeneity "problem"

$$\begin{aligned} \rho_{12} &= Corr[\epsilon_{i1} + u_{i}, \epsilon_{i2} + u_{i}] = \sigma_{u}^{2} / (\sigma_{\epsilon}^{2} + \sigma_{u}^{2}) \\ \tau &= Corr[v_{i2}, \epsilon_{i2} + u_{i}] = Corr[v_{i2} + \delta(\epsilon_{i2} + u_{i}), \epsilon_{i2} + u_{i}) \end{aligned}$$

Methods of Estimating the Attrition Model

- Heckman style "selection" model
- Two step maximum likelihood
- Full information maximum likelihood
- Two step method of moments estimators
- Weighting schemes that account for the "survivor bias"

Selection Model

Reduced form probit model for second period observation equation

$$Z_{i2}^{*} = X_{i2}^{\prime}(\theta + \delta\beta) + W_{i}^{\prime}\alpha + \delta(\varepsilon_{i2} + u_{i} + v_{i})$$
$$= r_{i2}^{\prime}\gamma + h_{i2}$$
$$Z_{i2} = 1(Z_{i2}^{*} > 0)$$

Conditional means for observations observed in the second period

$$\mathsf{E}[\mathsf{y}_{i2} \mid \mathsf{x}_{i2}, \mathsf{z}_{i2} = 1] = \mathsf{x}_{i2}'\beta + (\rho_{12}\sigma_{\varepsilon})\frac{\phi(\mathsf{r}_{i2}'\gamma)}{\Phi(\mathsf{r}_{i2}'\gamma)}$$

First period conditional means for observations observed in the second period

$$\mathsf{E}[\mathsf{y}_{i1} \mid \mathsf{x}_{i1}, \mathsf{z}_{i2} = 1] = \mathsf{x}_{i1}'\beta + (\rho_{12}\tau\sigma_{\varepsilon})\frac{\phi(\mathsf{r}_{i2}'\gamma)}{\Phi(\mathsf{r}_{i2}'\gamma)}$$

(1) Estimate probit equation

(2) Combine these two equations with a period dummy variable, use

OLS with a constructed regressor in the second period

THE TWO DISTURBANCES ARE CORRELATED.

TREAT THIS IS A SUR MODEL. (EQUIVALENT TO MDE)

Maximum Likelihood

$$\begin{split} \text{LogL}_{i} &= \frac{-\log 2\pi}{2} - \log \sigma_{\epsilon} - \frac{(y_{i1} - x_{i1}'\beta)^{2}}{2\sigma_{\epsilon}^{2}} \\ &+ z_{i2} \begin{bmatrix} \log \sigma_{\epsilon} + \log \sqrt{1 - \rho_{12}^{2}} - \frac{[(y_{i2} - \rho_{12}y_{i1}) - (x_{12} - \rho_{12}x_{i1})'\beta)]^{2}}{2\sigma_{\epsilon}^{2}(1 - \rho_{12}^{2})} \\ &+ \log \Phi \bigg(\frac{r_{i2}'\gamma + (\tau / \sigma_{\epsilon})(y_{i2} - x_{i2}'\beta)}{\sqrt{1 - \tau^{2}}} \bigg) \\ &+ (1 - z_{i2}) \bigg[\log \Phi \bigg(- \frac{r_{i2}'\gamma + (\rho_{12}\tau / \sigma_{\epsilon})(y_{i1} - x_{i1}'\beta)}{\sqrt{1 - \rho_{12}\tau^{2}}} \bigg) \bigg] \end{split}$$

- (1) See H&W for FIML estimation
- (2) Use the invariance principle to reparameterize
- (3) Estimate γ separately and use a two step ML with Murphy and Topel correction of asymptotic covariance matrix.

TABLE IV^a

PARAMETER ESTIMATES OF THE EARNINGS FUNCTION STRUCTURAL MODEL WITH AND WITHOUT A CORRECTION FOR ATTRITION

	With attrition maximum likeli (standard	n correction: hood estimates d errors)	Without attrition correction: generalized least squares estimates (standard errors)		
Variables	Earnings function parameters	Attrition parameters	Earnings function parameters		
Constant	5.8539	6347	5.8911		
	(0.0903)	(.3351)	(.0829)		
Experimental effect	0822	.2414	0793		
	(.0402)	(.1211)	(.0390)		
Time effect	.0940	_	.0841		
	(.0520)	_	(.0358)		
Education	.0209	0204	.0136		
	(.0052)	(.0244)	(.0050)		
Experience	.0037	0038	.0020		
	(.0013)	(.0061)	(.0013)		
Income	0131	.1752	0115		
	(.0050)	(.0470)	(.0044)		
Union	.2159	1.4290	.2853		
	(.0362)	(0.1252)	(.0330)		
Poor health	0601	.2480	0578		
	(.0330)	(.1237)	(.0326)		
	$\hat{\sigma}_{\eta}^2 = .1832$ (.0057)	<i>l</i> * = 64.35	$\hat{\sigma}_{\eta}^2 = .1236$		
	$\hat{\rho}_{12} = .2596$ (.0391)	$ \rho_{23} =1089 $ (1.0429)	$\hat{\rho}_{12} = .2003$		

A Model of Attrition

- Nijman and Verbeek, Journal of Applied Econometrics, 1992
- Consumption survey (Holland, 1984 1986)
 - Exogenous selection for participation (rotating panel)
 - Voluntary participation (missing not at random attrition)

Attrition Model

The main equation

$$\begin{split} y_{i,t} &= \beta_0 + \bm{x}'_{i,t} \bm{\beta} + \tilde{\alpha}_i + \epsilon_{i,t}, \text{ Random effects consumption function} \\ \tilde{\alpha}_i &= \overline{\bm{x}}'_i \bm{\theta} + \bm{u}_i, & \text{Mundlak device; } \bm{u}_i \text{ uncorrelated with } \bm{X}_i \\ y_{i,t} &= \beta_0 + \bm{x}'_{i,t} \bm{\beta} + \overline{\bm{x}}'_i \bm{\theta} + \bm{u}_i + \epsilon_{i,t}, \text{ Reduced form random effects model} \\ \text{The selection mechanism} \end{split}$$

a_{it} = 1[individual i asked to participate in period t] Purely exogenous

a_{it} may depend on observables, but does not depend on unobservables

- $r_{it} = 1$ [individual i chooses to participate if asked] Endogenous.
 - r_{it} is the endogenous participation dummy variable
 - $\boldsymbol{a}_{_{it}}=0 \Longrightarrow \boldsymbol{r}_{_{it}}=0$
 - $a_{it} = 1 \Rightarrow$ the selection mechanism operates

Selection Equation

The main equation

 $y_{i,t} = \beta_0 + \mathbf{x}'_{i,t}\mathbf{\beta} + \overline{\mathbf{x}}'_i\mathbf{\theta} + \mathbf{u}_i + \varepsilon_{i,t}$, Reduced form random effects model The selection mechanism.

38 38

Str. Malling

 $\begin{array}{ll} r_{it} &= 1 [\text{individual i chooses to participate if asked}] \quad \text{Endogenous.} \\ &r_{it} \text{ is the endogenous participation dummy variable} \\ &a_{it} = 0 \Rightarrow r_{it} = 0 \\ &a_{it} = 1 \Rightarrow \quad \text{the selection mechanism operates} \\ &r_{it} &= 1 [\gamma_0 + \textbf{x}'_{i,t}\gamma + \overline{\textbf{x}}'_{i}\mu + \textbf{z}'_{i,t}\delta + \textbf{v}_i + \textbf{w}_{i,t} > 0] \text{ all observed if } a_{it} = 1 \\ & \quad \text{State dependence: } \textbf{z} \text{ may include } r_{i,t-1} \end{array}$

Latent persistent unobserved heterogeneity: $\sigma_v^2 > 0$. "Selection" arises if $Cov[\epsilon_{i,t}, w_{i,t}] \neq 0$ or $Cov[u_i, v_i] \neq 0$

Estimation Using One Wave

ND YOUR

- Use any single wave as a cross section with observed lagged values.
- Advantage: Familiar sample selection model
- Disadvantages
 - Loss of efficiency
 - "One can no longer distinguish between state dependence and unobserved heterogeneity."

One Wave Model

A standard sample selection model.

$$\begin{split} y_{it} &= \beta_0 + \bm{x}_{it}' \bm{\beta} + \overline{\bm{x}}_i' \bm{\theta} + (\bm{u}_i + \bm{\epsilon}_{it}) \\ r_{it} &= \bm{1} [\gamma_0 + \bm{x}_{it}' \bm{\gamma} + \overline{\bm{x}}_i' \bm{\mu} + \delta_1 r_{i,t-1} + \delta_2 \bm{a}_{i,t-1} + (\bm{v}_i + \bm{w}_{it}) > 0] \\ \text{With only one period of data and } r_{i,t-1} \text{ exogenous,} \\ \text{this is the Heckman sample selection model.} \\ \text{If } > 0, \text{ then } r_{i,t-1} \text{ is correlated with } \bm{v}_i \text{ and the Heckman approach fails.} \end{split}$$

An assumption is required:

(1) Include $r_{i,t-1}$ and assume no unobserved heterogeneity (2) Exclude $r_{i,t-1}$ and assume there is no state dependence. In either case, now if $Cov[(u_i + \varepsilon_{it}), (v_i + w_{it})]$ we can use OLS. Otherwise, use the maximum likelihood estimator.

Maximum Likelihood Estimation

- Because numerical integration is required in one or two dimensions for every individual in the sample at each iteration of a high dimensional numerical optimization problem, this is, though feasible, not computationally attractive.
 - The dimensionality of the optimization is irrelevant
 - This is much easier in 2008 than it was in 1992 (especially with simulation) The authors did the computations with Hermite quadrature.

Testing for Selection?

States and and

- Selectivity is parameterized in these models coefficients are correlations or covariances.
- Maximum Likelihood Results
 - Covariances were highly insignificant.
 - \square LR statistic=0.46.
- Two step results produced the same conclusion based on a Hausman test
- ML Estimation results looked like the two step results.

Selectivity in Nonlinear Models

A.Matha

- 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

A Bivariate Probit Model

Labor Force Participation Equation

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

Participant and Part Time

Prob[f=0,d=1] = Prob[f=0|d=1]Prob[d=1]
= Bivariate Normal(
$$\beta' \mathbf{x}, -\alpha' \mathbf{z}, -\rho$$
)

[Topic 7-Selection] 53/81

FT/PT Selection Model

FIML Estimates of Bivariate	Probit Model
Dependent variable	FULLFP
Weighting variable	None
Number of observations	753
Log likelihood function	-723.9798
Number of parameters	16
Selection model based on LFE)

Full Time = Hours > 1000

+----+ |Variable | Coefficient | Standard Error |b/St.Er.|P[|Z|>z] | Mean of X| +----+

	Index	equation	for F	ULLTIME							
Constant	.94	532822	1.6	51674948		.585	.5	587			
WW	02	764944	.0	1941006	-1	.424	.1	543	4	.17768	154
KL6	.04	098432	.2	26250878		.156	.8	759		.14018	692
K618	13	640024	.0	5930081	-2	.300	.0	214	1	.35046	729
AGE	.03	543435	.0	7530788		.471	.6	380	4	1.9719	626
AGESQ	00	043848	.0	0088406	-	.496	.6	199	1	821.12	150
WE	08	622974	.0	2808185	- 3	.071	.0	021	1	2.6588	785
FAMINC	.2109	71D-04	.503	3746D-05	4	.188	.0	000	2	4130.42	229
	Index	equation	for I	JFP							
Constant	.98	337341	.5	0679582	1	.940	.0	523			
KL6	88	485756	.1	1251971	-7	.864	.0	000		.23771	580
К618	04	101187	.0	4020437	-1	.020	.3	077	1	.35325	365
WA	02	462108	.0	1308154	-1	.882	.0	598	4	2.53784	486
WE	.16	636047	.0	2738447	б	.075	.0	000	1	2.2868	526
HA	01	652335	.0	1287662	-1	.283	.1	994	4	5.12084	499
HE	06	276470	.0)1912877	- 3	.281	.0	010	1	2.4913	679
	Disturban	ce correl	ation	ı							
RHO(1,2)	84	102682	.2	25122229	-3	.348	.0	800			

Building a Likelihood for a Poisson Regression Model with Selection

Poisson Probability Functions

 $\mathsf{P}(\mathsf{y}_{i} \mid \mathbf{x}_{i}) = \exp(-\lambda_{i})\lambda_{i}^{y} / y_{i}!$

Covariates and Unobserved Heterogeneity

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

Conditional Contribution to the Log Likelihood

$$logL_{i} | \epsilon_{i} = -\lambda(\mathbf{x}_{i}, \epsilon_{i}) + y_{i} \log \lambda(\mathbf{x}_{i}, \epsilon_{i}) - \log y_{i}!$$

Probit Selection Mechanism

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

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 $L_i(y_i, d_i = 1) | \epsilon_{i} = [f(y_i | \mathbf{x}_i, \epsilon_i, d_i = 1) Prob[d_i = 1 | \mathbf{z}_i, \epsilon_i]$ $L_i(d_i = 0) = Prob[d_i = 0 | \mathbf{z}_i, \epsilon_i]$

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

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

 $logL_i = logf(y_i, d_i)$

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

Poisson Model with Sample Selection. Dependent variable DOCVIS Log likelihood function -3592.42064 Restricted log likelihood -6076.83457 Chi squared [2 d.f.] 4968.82786 Significance level .00000 McFadden Pseudo R-squared .4088336 Estimation based on N = 27326, K = 12 Inf.Cr.AIC = 7208.8 AIC/N = .264 Mean of LHS Variable = 3.12451 Restr. Log-L is Poisson+Probit (indep). LogL for initial probit = -2442.4091 LogL for initial Poisson= -3634.4254 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 EDUC HHNINC MARRIED HHKIDS Constant AGE EDUC HHNINC Sigma Rho	Parameters of P 1.22286 .01528** 06984** 38472 22330 18695 Parameters of P: -3.40612*** .00634*** .06403*** .75095*** Standard Deviat 1.07226*** Correlation of 2 02597	oisson/Neg. E 1.03286 .00654 .02745 .36241 .15519 .14693 robit Selecti .12282 .00178 .00757 .09452 ion of Hetero .05985 Heterogeneity .25188	Binomial 1.18 2.34 -2.54 -1.06 -1.44 -1.27 on Model -27.73 3.57 8.45 7.95 ogeneity 17.92 & Selec 10	Probabil .2364 .0194 .0110 .2884 .1502 .2032 .0000 .0004 .0000 .0000 .0000 .0000 .0000	ity 80151 .00247 12364 -1.09504 52748 47493 -3.64684 .00286 .04919 .56570 .95495 51964	3.24722 .02809 01603 .32559 .08087 .10103 -3.16540 .00982 .07887 .93620 1.18957 .46770	

THE P

100

Poisson Regression Dependent variable DOCVIS Log likelihood function -1659.65094 Restricted log likelihood -1751.09275 Chi squared [6 d.f.] 182.88363 Significance level .00000 McFadden Pseudo R-squared .0522199 Estimation based on N = 27326, K = 7 Inf.Cr.AIC = 3333.3 AIC/N = .122 Chi- squared = 2887.14468 RsqP= .1766 G - squared = 2219.26013 RsqD= .0761 Overdispersion tests: g=mu(i) : 5.121 Overdispersion tests: g=mu(i)^2: 5.728 Cov. matrix corrected for 2 step estimation							
DOCVIS	Coefficient	Standard Error	z	Prob. z >Z *	95% Co Int	Confidence Interval	
Constant AGE EDUC HHNINC MARRIED HHKIDS MlsRatio	1.13394 .01440 05828 19541 37640 19118 .21226	4082.032 6.37321 63.44852 742.8324 .62847 .27342 1130.177	.00 .00 .00 60 70 .00	.9998 .9982 .9993 .9998 .5492 .4844 .9999	-7999.50166 -12.47686 -124.41510 -1456.12012 -1.60818 72707 -2214.89378	8001.76954 12.50566 124.29854 1455.72930 .85539 .34472 2215.31831	
Note: ***	; **, * ==> Si	gnificance at	1%, 5%,	10% le	vel.		

a

Stochastic Frontier Model: ML

AL TH

Statute and

 $y_i = \boldsymbol{\beta' x}_i + v_i - u_i$

where $u_i = |\sigma_u U_i| = \sigma_u |U_i|, U_i \sim N[0, 1^2],$

 $v_i = \sigma_v V_i, V_i \sim N[0, 1^2].$

$$\log L(\boldsymbol{\beta}, \boldsymbol{\sigma}, \boldsymbol{\lambda}) = \sum_{i=1}^{N} \left[\frac{1}{2} \log \left(\frac{2}{\pi} \right) - \log \boldsymbol{\sigma} - \frac{1}{2} (\varepsilon_i / \boldsymbol{\sigma})^2 + \log \Phi(-\lambda \varepsilon_i / \boldsymbol{\sigma}) \right]$$

where $\varepsilon_i = y_i - \beta' \mathbf{x}_i = v_i - u_i$,

$$\lambda \qquad = \sigma_u / \sigma_v,$$

$$\sigma = \sqrt{\sigma_v^2 + \sigma_u^2}$$

Sample Selected SF Model

$$d_{i} = \mathbf{1}[\boldsymbol{\alpha}'\mathbf{z}_{i} + w_{i} > 0], \ w_{i} \sim \mathbf{N}[0, 1^{2}]$$

$$y_{i} = \boldsymbol{\beta}'\mathbf{x}_{i} + \varepsilon_{i}, \ \varepsilon_{i} \sim \mathbf{N}[0, \sigma_{\varepsilon}^{2}]$$

$$(y_{i}, \mathbf{x}_{i}) \text{ observed only when } d_{i} = 1.$$

$$\varepsilon_{i} = v_{i} - u_{i}$$

$$u_{i} = |\sigma_{u}U_{i}| = \sigma_{u} |U_{i}| \text{ where } U_{i} \sim \mathbf{N}[0, 1^{2}]$$

$$v_{i} = \sigma_{v}V_{i} \text{ where } V_{i} \sim \mathbf{N}[0, 1^{2}].$$

$$(w_{i}, v_{i}) \sim \mathbf{N}_{2}[(0, 1), (1, \rho\sigma_{v}, \sigma_{v}^{2})]$$

Same in the local

$$f(y_{i} | \mathbf{x}_{i}, | U_{i} |, d_{i}, \mathbf{z}_{i}) = \begin{bmatrix} d_{i} \frac{\exp\left(-\frac{1}{2}(y_{i} - \beta' x_{i} + \sigma_{u} | U_{i} |)^{2} / \sigma_{v}^{2}\right)}{\sigma_{v} \sqrt{2\pi}} \Phi\left(\frac{\rho(y_{i} - \beta' x_{i} + \sigma_{u} | U_{i} |) / \sigma_{\varepsilon} + \boldsymbol{\alpha}' \mathbf{z}_{i}}{\sqrt{1 - \rho^{2}}}\right) + (1 - d_{i}) \Phi(-\boldsymbol{\alpha}' \mathbf{z}_{i}) \end{bmatrix}$$

Sample Selection in a Nonlinear Model

$$d_i = 1(\alpha' \mathbf{z}_i + w_i > 0) \ w_i \sim N[0,1],$$

$$g_i | \varepsilon_i = g(\beta' \mathbf{x}_i, \sigma_{\varepsilon} \varepsilon_i) \varepsilon_i \sim N[0, 1]$$

$$y_i | \mathbf{x}_i, \varepsilon_i \sim f[y_i | g(\boldsymbol{\beta}' \mathbf{x}_i, \sigma_{\varepsilon} \varepsilon_i)]$$

$$[w_i, \varepsilon_i] \sim N[(0,1), (1,\rho,1)]$$

 y_i , \mathbf{x}_i are observed only when $z_i = 1$.

$$f(y_i, d_i | \mathbf{x}_i, \mathbf{z}_i) = \int_{-\infty}^{\infty} \{ (1 - d_i) + d_i f[y_i] g(\boldsymbol{\beta}' \mathbf{x}_i, \boldsymbol{\sigma}_{\varepsilon} \varepsilon_i)] \} \times \\ \Phi \Big((2d_i - 1) [\boldsymbol{\alpha}' \mathbf{z}_i + \rho \varepsilon_i] / \sqrt{1 - \rho^2} \Big) \phi(\varepsilon_i) d\varepsilon_i ,$$

Simulated Log Likelihood for a Simpler Model

$$\log L_{S}(\boldsymbol{\beta}, \boldsymbol{\sigma}_{u}, \boldsymbol{\sigma}_{v}, \boldsymbol{\alpha}, \boldsymbol{\rho}) = \sum_{i=1}^{N} \log \frac{1}{R} \sum_{r=1}^{R} \left[\begin{array}{c} d_{i} \left[\frac{\exp\left(-\frac{1}{2}(y_{i} - \boldsymbol{\beta}' \mathbf{x}_{i} + \boldsymbol{\sigma}_{u} | \boldsymbol{U}_{ir} |)^{2} / \boldsymbol{\sigma}_{v}^{2}\right) \right] \\ \Phi\left(\frac{\rho(y_{i} - \boldsymbol{\beta}' \mathbf{x}_{i} + \boldsymbol{\sigma}_{u} | \boldsymbol{U}_{ir} |) / \boldsymbol{\sigma}_{\varepsilon} + \boldsymbol{\alpha}' \mathbf{z}_{i}}{\sqrt{1 - \boldsymbol{\rho}^{2}}} \right) \right] \\ + (1 - d_{i}) \Phi(-\boldsymbol{\alpha}' \mathbf{z}_{i}) \end{array}$$

[Topic 7-Selection] 64/81

A 2 Step MSL Approach

MALTIN

$$\log L_{S,C}(\beta,\sigma_u,\sigma_v,\rho) = \sum_{i=1}^{N} \log \frac{1}{R} \sum_{r=1}^{R} \left[\frac{d_i \frac{\exp\left(-\frac{1}{2}(y_i - \beta' \mathbf{x}_i + \sigma_u |U_{ir}|)^2 / \sigma_v^2)\right)}{\sigma_v \sqrt{2\pi}} \times \right] + (1 - d_i)\Phi(-a_i)$$

where $a_i = \hat{\boldsymbol{\alpha}}' \mathbf{z}_i$

$$\log L_{S,C}(\beta,\sigma_u,\sigma_v,\rho) = \sum_{d_i=1} \log \frac{1}{R} \sum_{r=1}^{R} \left[\frac{\exp\left(-\frac{1}{2}(y_i - \beta' \mathbf{x}_i + \sigma_u |U_{ir}|)^2 / \sigma_v^2)\right)}{\sigma_v \sqrt{2\pi}} \times \right] \Phi\left(\frac{\rho(y_i - \beta' \mathbf{x}_i + \sigma_u |U_{ir}|) / \sigma_\varepsilon + a_i}{\sqrt{1 - \rho^2}}\right) \right]$$

[Topic 7-Selection] 65/81

Simulated ML for the SF Model

Same in

$$f(y_{i} | \mathbf{x}_{i}, |U_{i}|) = \frac{\exp[-\frac{1}{2}(y_{i} - \boldsymbol{\beta}'\mathbf{x}_{i} + \sigma_{u} / U_{i} |)^{2} / \sigma_{v}^{2}]}{\sigma_{v} \sqrt{2\pi}}$$

$$f(y_{i} | \mathbf{x}_{i}) = \int_{|U_{i}|} \frac{\exp[-\frac{1}{2}(y_{i} - \boldsymbol{\beta}'\mathbf{x}_{i} + \sigma_{u} / U_{i} |)^{2} / \sigma_{v}^{2}]}{\sigma_{v} \sqrt{2\pi}} p(|U_{i}|) d |U_{i}|$$

$$p(|U_{i}|) = \frac{2\exp[-\frac{1}{2}|U_{i}|^{2}]}{\sqrt{2\pi}}, |U_{i}| \ge 0.$$

$$f(y | \mathbf{x}_{i}) \approx \frac{1}{R} \sum_{r=1}^{R} \frac{\exp[-\frac{1}{2}(y_{i} - \boldsymbol{\beta}'\mathbf{x}_{i} + \sigma_{u} / U_{ir} |)^{2} / \sigma_{v}^{2}]}{\sigma_{v} \sqrt{2\pi}}$$

$$\log L_{S}(\boldsymbol{\beta},\boldsymbol{\sigma}_{u},\boldsymbol{\sigma}_{v}) = \sum_{i=1}^{N} \log \left\{ \frac{1}{R} \sum_{r=1}^{R} \frac{\exp[-\frac{1}{2}(y_{i} - \boldsymbol{\beta}' \mathbf{x}_{i} + \boldsymbol{\sigma}_{u} / U_{ir} |)^{2} / \boldsymbol{\sigma}_{v}^{2}]}{\boldsymbol{\sigma}_{v} \sqrt{2\pi}} \right\}$$

This is simply a linear regression with a random constant term, $\alpha_i = \alpha - \sigma_u |U_i|$

Nonnormality Issue

St. Mallin .

- 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.

A Study of Health Status in the Presence of Attrition



Research Article

The dynamics of health in the British Household Panel Survey



Model for Self Assessed Health

- British Household Panel Survey (BHPS)
 - Waves 1-8, 1991-1998
 - Self assessed health on 0,1,2,3,4 scale
 - Sociological and demographic covariates
 - Dynamics inertia in reporting of top scale
- Dynamic ordered probit model
 - Balanced panel analyze dynamics
 - Unbalanced panel examine attrition

Dynamic Ordered Probit Model

Latent Regression - Random Utility $h_{it}^* = \beta' \mathbf{x}_{it} + \gamma' \mathbf{H}_{i,t-1} + \alpha_i + \varepsilon_{it}$ $\mathbf{x}_{it} = \text{relevant covariates and control variables}$ $\mathbf{H}_{i,t-1} = 0/1 \text{ indicators of reported health status in previous period}$ $\mathbf{H}_{i,t-1}(j) = 1[\text{Individual i reported h}_{it} = j \text{ in previous period}], j=0,...,4$ Ordered Choice Observation Mechanism $h_{it} = j \text{ if } \mu_{j-1} < h_{it}^* \leq \mu_j, j = 0,1,2,3,4$

Ordered Probit Model - $\varepsilon_{ii} \sim N[0,1]$ Random Effects with Mundlak Correction and Initial Conditions $\alpha_i = \alpha_0 + [\alpha'_1 \mathbf{H}_{i,1} + \alpha'_2 \overline{\mathbf{x}}_i] + u_i, \ u_i \sim N[0,\sigma^2]$

Random Effects Dynamic Ordered Probit Model

A MALTIN ST

Random Effects Dynamic Ordered Probit Model $\begin{aligned} h_{it} ^{*} &= \mathbf{x}_{it}^{\prime} \boldsymbol{\beta} + \boldsymbol{\Sigma}_{j=1}^{J} \boldsymbol{\gamma}_{j} h_{i,t-1}(j) + \boldsymbol{\alpha}_{i} + \boldsymbol{\epsilon}_{i,t} \\ h_{i,t} &= j \text{ if } \boldsymbol{\mu}_{j-1} < h_{it} ^{*} < \boldsymbol{\mu}_{j} \\ h_{i,t}(j) &= 1 \text{ if } h_{i,t} = j \\ P_{it,j} &= P[h_{it} = j] = \Phi(\boldsymbol{\mu}_{j} - \mathbf{x}_{it}^{\prime} \boldsymbol{\beta} - \boldsymbol{\Sigma}_{j=1}^{J} \boldsymbol{\gamma}_{j} h_{i,t-1}(j) - \boldsymbol{\alpha}_{i}) \\ &- \Phi(\boldsymbol{\mu}_{j-1} - \mathbf{x}_{it}^{\prime} \boldsymbol{\beta} - \boldsymbol{\Sigma}_{j=1}^{J} \boldsymbol{\gamma}_{j} h_{i,t-1}(j) - \boldsymbol{\alpha}_{i}) \end{aligned}$

Parameterize Random Effects

$$\alpha_{i} = \alpha_{0} + \Sigma_{j=1}^{J} \alpha_{1,j} \mathbf{h}_{i,1}(j) + \boldsymbol{\alpha}' \overline{\mathbf{x}}_{i} + \mathbf{u}_{i}$$

Simulation or Quadrature Based Estimation

$$InL = \sum_{i=1}^{N} In \int_{\alpha_i} \prod_{t=1}^{T_i} P_{it,j} f(\alpha_j) d\alpha_j$$

Data

TTP:

Service and

Table I. Variable definitions

SAHSelf-Assessed Health: 5 if excellent, 4 if good, 3 if fair, 2 if poor, 1 if veryWIDOW1 if widowed, 0 otherwiseSINGLE1 if never married, 0 otherwiseDIV/SEP1 if divorced or separated, 0 otherwiseNON-WHITE1 if a member of ethnic group other than white, 0 otherwiseDEGREE1 if highest academic qualification is a degree or higher degree, 0 otherwiseHND/A1 if highest academic qualification is O level or CSE, 0 otherwiseO/CSE1 if highest academic qualification is O level or CSE, 0 otherwiseHHSIZENumber of people in household including respondentNCH04Number of children in household aged 0-4NCH1218Number of children in household aged 12-18INCOMEEquivalized annual real household income in poundsAGEAge in years at 1st December of current wave	poor
--	------
Variable of Interest



Figure 1. Self-assessed health status by wave



Dynamics

Table II. Transition matrices, balanced panel

(a) Men

SAH	EX	GOOD	FAIR	POOR	VERY POOR	Ν
EX	0.600	0.342	0.046	0.010	0.002	5485
GOOD	0.184	0.651	0.142	0.019	0.004	9263
FAIR	0.055	0.361	0.471	0.100	0.012	3433
POOR	0.029	0.120	0.340	0.418	0.093	1031
VERY POOR	0.032	0.073	0.133	0.423	0.339	248
Ν	5231	9287	3565	1111	266	19460

(b) Women

SAH	EX	GOOD	FAIR	POOR	VERY POOR	Ν
EX	0.572	0.353	0.059	0.013	0.004	5164
GOOD	0.150	0.657	0.162	0.026	0.005	11 306
FAIR	0.040	0.362	0.465	0.116	0.017	4928
POOR	0.021	0.156	0.360	0.365	0.098	1587
VERY POOR	0.014	0.106	0.192	0.326	0.362	423
Ν	4884	11 329	5082	1649	464	23 408



Attrition

Table V. Sample size, drop-outs and attrition rates by wave

(a) All data

FULL SAMPLE					EX at <i>t</i> – 1	GOOD at $t - 1$	FAIR at $t - 1$	POOR at $t - 1$	VPOOR at $t - 1$
Wave	No. individuals	Survival rate	Drop-outs	Attrition rate	Attrition rate	Attrition rate	Attrition rate	Attrition rate	Attrition rate
1	10 2 5 6								
2	8957	87.33%	1299	12.67%	11.54%	12.57%	13.01%	13.73%	23.74%
3	8162	79.58%	795	8.88%	8.08%	8.13%	9.65%	12.62%	19.46%
4	7825	76.30%	337	4.13%	6.67%	6.54%	6.73%	10.35%	14.74%
5	7430	72.45%	395	5.05%	6.21%	6.18%	7.87%	9.11%	16.34%
6	7238	70.57%	192	2.58%	3.11%	3.24%	5.06%	10.47%	18.83%
7	7102	69.25%	136	1.88%	3.15%	3.85%	4.79%	8.83%	8.75%
8	6839	66.68%	263	3.70%	3.43%	3.82%	5.30%	5.88%	17.01%

Testing for Attrition Bias

Table 9: Verbeek and Nijman tests for attrition: based on dynamic ordered probit models with Wooldridge specification of correlated effects and initial conditions

		MEN				WOMEN		
	β	Std.err.	t-test	p-value	β	Std.err.	t-test	p-value
NEXT WAVE	.199	.035	5.67	.000	.060	.034	1.77	.077
ALL WAVES	.139	.031	4.46	.000	.071	.029	2.45	.014
NUMBER OF	.031	.009	3.54	.000	.016	.008	1.88	.060
WAVES								

Three dummy variables added to full model with unbalanced panel suggest presence of attrition effects.

Probability Weighting Estimators

- A Patch for Attrition
- (1) Fit a participation probit equation for each wave.
- (2) Compute p(i,t) = predictions of participation for each individual in each period.
 - Special assumptions needed to make this work
- Ignore common effects and fit a weighted pooled log likelihood: Σ_i Σ_t [d_{it}/p(i,t)]logLP_{it}.

Attrition Model with IP Weights



Assumes (1) Prob(attrition|all data) = Prob(attrition|selected variables) (ignorability) (2) Attrition is an 'absorbing state.' No reentry. Obviously not true for the GSOEP data above. Can deal with point (2) by isolating a subsample of those present at wave 1 and the monotonically shrinking subsample as the waves progress.

Inverse Probability Weighting

Str. Marine Local

Panel is based on those present at WAVE 1, N1 individuals Attrition is an absorbing state. No reentry, so $N1 \ge N2 \ge ... \ge N8$. Sample is restricted at each wave to individuals who were present at the previous wave.

- $d_{it} = 1$ [Individual is present at wave t].
- $\mathbf{d}_{i1} = 1 \quad \forall \quad \mathbf{i}, \, \mathbf{d}_{it} = \mathbf{0} \implies \mathbf{d}_{i,t+1} = \mathbf{0}.$
- $\tilde{\mathbf{x}}_{i1}$ = covariates observed for all i at entry that relate to likelihood of being present at subsequent waves.

(health problems, disability, psychological well being, self employment, unemployment, maternity leave, student, caring for family member, ...)

Probit model for $d_{it} = 1[\delta' \tilde{\mathbf{x}}_{i1} + w_{it}]$, t = 2,...,8. $\hat{\pi}_{it} =$ fitted probability.

Assuming attrition decisions are independent, $\hat{P}_{it} = \prod_{s=1}^{t} \hat{\pi}_{is}$

Inverse probability weight $\hat{W}_{it} = \frac{d_{it}}{\hat{P}_{it}}$ Weighted log likelihood $\log L_W = \sum_{i=1}^{N} \sum_{t=1}^{8} \log L_{it}$ (No common effects.)

[Topic 7-Selection] 79/81

Estimated Partial Effects by Model

Та	able 12: Average p	artial effects on prol	bability of reporting	excellent health for	selected variables	
a) Men						
	(1) Pooled model, balanced sample	(2) Pooled model, unbalanced sample	(3) Pooled model, IPW-1	(4) Pooled model, IPW-2	(5) Random effects, balanced sample	(6) Random effects, unbalanced sample
Ln(INCOME)	.009 (.004)	.009 (.004)	.009 (.004)	.011 (.005)	.013 (.006)	.012 (.005)
Mean Ln(INCOME)	.049 (.024)	.043 (.022)	.042 (.021)	.045 (.022)	.066 (.028)	.056 (.025)
DEGREE	.010 (.005)	.017 (.009)	.018 (.009)	.018 (.009)	.015 (.006)	.027 (.012)
HND/A	.019 (.009)	.021 (.011)	.021 (.010)	.022 (.011)	.028 (.011)	.030 (.013)
O/CSE	.016 (.008)	.020 (.010)	.020 (.010)	.020 (.010)	.024 (.010)	.028 (.012)
SAHEX(t-1)	.234 (.087)	.231 (.090)	.231 (.090)	.230 (.089)	.082 (.031)	.085 (.035)
SAHFAIR(t-1)	170 (.085)	163 (.084)	162 (.084)	162 (.083)	080 (.034)	077 (.036)
SAHPOOR(t-1)	242 (.167)	233 (.163)	232 (.162)	232 (.162)	151 (.077)	145 (.078)
SAHVPOOR(t-1)	260 (.198)	253 (.197)	255 (.199)	255 (.200)	184 (.104)	179 (.106)

Partial Effect for a Category



These are 4 dummy variables for state in the previous period. Using first differences, the 0.234 estimated for SAHEX means transition from EXCELLENT in the previous period to GOOD in the previous period, where GOOD is the omitted category. Likewise for the other 3 previous state variables. The margin from 'POOR' to 'GOOD' was not interesting in the paper. The better margin would have been from EXCELLENT to POOR, which would have (EX,POOR) change from (1,0) to (0,1).

[Topic 7-Selection] 81/81