### Statistical Methods, part 1 Module 2: Latent Class Analysis of Survey Error LCA

Dan Hedlin Stockholm University November 2012

Acknowledgement: Paul Biemer



### Latent Variable Analysis

	Latent Variable	Manifest Variable
Factor analysis	Continuous	Continuous
Latent trait analysis	Continuous	Discrete
Latent profile analysis	Discrete	Continuous
Latent class analysis	Discrete	Discrete

Bartholomew, D.J., and Knott, M. (1999). Latent Variable Models and Factor Analysis. London: Arnold.

# Latent Class Analysis (LCA)

• Developed in 1950's

– Lazarsfeld, Henry, Green

- Outcome variables can be dichotomous or polytomous (nominal or ordinal)
- Very extensive field
- Advanced use requires knowledge of loglinear models with and without latent variables

# Some applications of LCA

Nonsampling Error Analysis

- Identifying flawed questions and other questionnaire problems
- Estimating census undercount in a capturerecapture framework
- Characterizing respondents, interviewers, and questionnaire elements that contribute to survey error
- Adjusting for nonresponse and missing data in surveys
- We will only touch upon LCA applications

# Some Other Applications (cont'd)

**Social Science Analysis** 

- Causal modeling
- Log-linear analysis compensating for measurement error
- Cluster analysis
- Variable reduction and scale construction

# Why in this course?

- Very useful class of models
- Few people seem to know about LCA
- Course plan indicates a content of science and models. LCA does involve scientific issues.

### Use of LCA for Investigating Survey Error

- LCA methods and models are prone to misuse
  - E.g., adjustment vs. evaluation
  - Lack of attention to assumptions
  - Weak indicators
- LCA is best use in conjunction with other methods
- Avoid the temptation to believe the model is true
  - Consider the risks if the model is not true
  - Take steps to verify the model assumptions and to validate the estimates

# Maximum Likelihood Estimation

Recall that  $\hat{\pi}, \hat{\theta}$ , and  $\hat{\phi}$  contain all the information we need to estimate:

- $\pi$ , the true population proportion
- *SV*, the sampling variance
- *SRV*, the simple response variance,
- *R*, the reliability ratio
- Measurement bias of the sample proportion, Bias(*p*)
- Total variance, Var(*p*)
- Total MSE(*p*)

We can estimate  $\hat{\pi}, \hat{\theta}$ , and  $\hat{\phi}$  in some cases using only information from remeasurements of the characteristics of interest.

#### Latent Class Models for Measurement Error

Consider the following:

		Reinterview $(y_2)$	
		1	2
Interview $(y_1)$	1	$n_{11}$	<i>n</i> <sub>12</sub>
	2	$n_{21}$	<i>n</i> <sub>22</sub>

Or more generally: any second measurement that is similar (or next identical in type) to first measurement (i.e. not necessarily a better measurement)

# Slight Change in Notation

 $y_{tr}$  denotes observation *t* on unit *r*  $y_{tr} = 1$  denotes positive response and = 2 denotes a negative response

#### Assume:

- SRS (simple random sample)
- Two parallel measurements are available
- $\gamma_{\theta\phi} = 0$  (homogeneity, ie constant probabilities for misclassification)

Pr  $(n_{11}, n_{12}, n_{21}, n_{22})$  is a multinomial

### **Multinomial Distribution**

# Consider the distribution of the $2 \times 2$ table with cells (11, 12, 21, 22)

Let  $\pi_{11} = P(y_i \text{ falls in cell 11})$  and define similarly.  $\pi_{12}, \pi_{21}, \pi_{22}$ 

Then,

$$\mathbf{P}(n_{11}, n_{12}, n_{21}, n_{22}) = \begin{pmatrix} n \\ n_{11} & n_{12} & n_{21} \\ n_{11} & n_{12} & n_{21} & n_{22} \end{pmatrix} \pi_{11}^{n_{11}} \pi_{12}^{n_{21}} \pi_{21}^{n_{22}} \pi_{22}^{n_{22}}$$

For the interview-reinterview table,

$$\pi_{11} = P(y_{1r} = 1, y_{2r} = 1)$$
  

$$\pi_{21} = P(y_{1r} = 2, y_{2r} = 1)$$
  

$$\pi_{12} = P(y_{1r} = 1, y_{2r} = 2)$$
  

$$\pi_{22} = P(y_{1r} = 2, y_{2r} = 2)$$

$$\pi_{11} = P(y_{1r} = 1, y_{2r} = 1 | \mu_r = 1) \times P(\mu_r = 1)$$
  
+  $P(y_{1r} = 1, y_{2r} = 1 | \mu_r = 2) \times P(\mu_r = 2)$   
=  $(1 - \theta)(1 - \theta)\pi + \phi\phi(1 - \pi)$   
$$\boxed{= (1 - \theta)^2 \pi + \phi^2 (1 - \pi)}$$

$$\pi_{21} = P(y_{1r} = 2, y_{2r} = 1 | \mu_r = 1) \times P(\mu_r = 1)$$
  
+  $P(y_{1r} = 2, y_{2r} = 1 | \mu_r = 2) \times P(\mu_r = 2)$   
=  $\theta(1 - \theta)\pi + (1 - \phi)\phi(1 - \pi)$ 

By assumption: 
$$\pi_{12} = \pi_{21}$$

$$\pi_{22} = P(y_{1r} = 2, y_{2r} = 2 | \mu_r = 1) \times P(\mu_r = 1)$$
  
+  $P(y_{1r} = 2, y_{2r} = 2 | \mu_r = 2) \times P(\mu_r = 2)$   
=  $\theta^2 \pi + (1 - \phi)^2 (1 - \pi)$ 

### Summarising:

Likelihood = L( $\pi$ ,  $\theta$ ,  $\phi$  |  $n_{11}$ ,  $n_{12}$ ,  $n_{21}$ ,  $n_{22}$ )

$$\mathbf{P}(n_{11}, n_{12}, n_{21}, n_{22}) = \begin{pmatrix} n \\ n_{11} & n_{12} & n_{21} \\ n_{11} & n_{12} & n_{21} & n_{22} \end{pmatrix} \pi_{11}^{n_{11}} \pi_{12}^{n_{21}} \pi_{21}^{n_{22}} \pi_{22}^{n_{22}}$$

$$= \begin{pmatrix} n \\ n_{11} & n_{12} & n_{21} & n_{22} \end{pmatrix} \begin{bmatrix} [(1-\theta)^2 \pi + \phi^2 (1-\pi)]^{n_{11}} \\ \times [\theta(1-\theta)\pi + \phi(1-\phi)(1-\pi)]^{n_{12}+n_{21}} \\ \times [\theta^2 \pi + (1-\theta)^2 (1-\pi)]^{n_{22}} \end{bmatrix}$$

# Maximum Likelihood Estimation

The best (ML) estimates of

π, θ, φ

are those that maximize *L* as a function of  $\pi$ ,  $\theta$ ,  $\phi$ Unfortunately, no unique solution exists (ie model is not identifiable).

3 parameters2 degrees of freedom (d.f.)

Why only 2 d.f. and not 3 since it is a 2x2 table?

# param > # d.f. and model is not identifiable.

Notation for common probability

Suppose we let

 $\theta = \phi = \varepsilon$ , say

[This is not realistic. Why not? Nevertheless...] Then, plugging into the previous formula,

L( $n_{11}, n_{12}, n_{21}, n_{22} | \pi, \varepsilon$ ) = constant [ $(1-\varepsilon)^2 \pi + \varepsilon^2 (1-\pi)$ ] $^{n_{11}} \times [\varepsilon(1-\varepsilon)]^{n_{12}+n_{21}}$ 

$$\times \left[\varepsilon^2 \pi + (1 - \varepsilon)^2 (1 - \pi)\right]^{n_{22}}$$

### Two parameters $\pi, \epsilon$ and two d.f., so now $\Rightarrow \pi, \epsilon$ estimable

But, no d.f. left over to estimate model fit (no chi-square test). Solution fits the data perfectly

==> model is <u>saturated</u>.

No d.f. test model validity or fit.



Find MLE of  $\pi, \varepsilon$ 

Numerical methods to solve equations, e.g.: grid search Newton-Raphson Iterative Proportional Fitting (IPF) EM-algorithm

### More General Model (Hui-Walter, 1980)

Interview  $\pi$ ,  $\theta(\text{int.})$ ,  $\phi(\text{int.})$  $\theta(\text{reint.})$ ,  $\theta(\text{reint.})$ 

Assumes that error distributions for interview and reinterview are not equal.

As before, we assume

- Single latent variable
- Local independence (we do not assume parallel measurements) Need at least 5 d.f. to estimate these five parameters:

 $\pi$ ,  $\theta$ (int.),  $\phi$ (int.),  $\theta$ (reint.),  $\phi$ (reint.)

### Trick: Introduce a Bivariate Grouping Variable, G

Group 1 (G=1)	Reinter	view	
		1	2
Interview	1	$n_{111}$	<i>n</i> <sub>112</sub>
	2	<i>n</i> <sub>121</sub>	<i>n</i> <sub>122</sub>
Group 2 ( <i>G</i> =2)	Reinter	view	
		1	2
Interview	1	<i>n</i> <sub>211</sub>	<i>n</i> <sub>212</sub>
	2	$n_{221}$	<i>n</i> <sub>222</sub>

How many d.f. are there?

- 3 from G=1 table
- 3 from G=2 table

==> up to 6 parameters can be estimated

10 Parameters in the unrestricted model

Group 1 
$$\pi_1$$
,  
 $\theta_1(\text{int.}) \theta_1(\text{reint.})$   
 $\phi_1(\text{int.}) \phi_1(\text{reint.})$   
 $\theta_2(\text{int.}) \theta_2(\text{reint.})$   
 $\phi_2(\text{int.}) \phi_2(\text{reint.})$ 

Need to constrain the model to save 4 df's

#### **Hui-Walter constraints**

$$\begin{aligned} \theta_1(\text{int.}) &= \theta_2(\text{int.}) = \theta_{\text{int}} \\ \phi_1(\text{int.}) &= \phi_2(\text{int.}) = \phi_{\text{int}} \\ \theta_1(\text{reint.}) &= \theta_2(\text{reint.}) = \theta_{\text{reint}} \\ \phi_1(\text{reint.}) &= \phi_2(\text{reint.}) = \phi_{\text{reint}} \end{aligned}$$

In words, assume classification probabilities do not differ by group. Now, 6 parameters and 6 d.f. so model can be estimated Completely saturated, i.e., 0 d.f. for error

### Likelihood

Likelihood = P( $\pi_1, \pi_2 \theta_{int}, \phi_{int}, \theta_{reint}, \phi_{reint} | n_{gij}, g = 1, 2, i = 1, 2, j = 1, 2$ )

 $\propto (\pi_{111})^{n_{111}} (\pi_{121})^{n_{121}} (\pi_{112})^{n_{112}} (\pi_{122})^{n_{122}} (\pi_{211})^{n_{211}} (\pi_{221})^{n_{221}} (\pi_{212})^{n_{212}} (\pi_{222})^{n_{222}}$ 

$$=\prod_{g}\prod_{i}\prod_{j}(\pi_{gij})^{n_{gij}}$$

where

$$\pi_{gij} = P(G = g, y_1 = i, y_2 = j)$$

# Likelihood

These probabilities can be rewritten in terms of  $\pi_1, \pi_2 \theta_{int}, \phi_{int}, \theta_{reint}, \phi_{reint}, \phi_{reint}$ 

$$\pi_{111} = \pi_{g=1} \left[ (1 - \theta_{int})(1 - \theta_{reint})\pi_1 + \phi_{int} \phi_{reint} (1 - \pi_1) \right]$$
  
$$\pi_{211} = \pi_{g=2} \left[ (1 - \theta_{int})(1 - \theta_{reint})\pi_2 + \phi_{int} \phi_{reint} (1 - \pi_2) \right]$$

and so on.

Plug these into L and maximize to obtain ML estimates

 $\hat{\pi}_1, \hat{\pi}_2 \hat{\theta}_{int}, \hat{\phi}_{int}, \hat{\theta}_{reint}, \hat{\phi}_{reint}$ 

### Results for 1996 US Current Population Survey Reinterview

#### **Input Data**

	Male	S			Fema	les	
Interview	Reinterview Response			Interview	Reinterview Response		
Response	EMP	UNEMP	NLF	Response	EMP	UNEMP	NLF
EMP	2372	14	29	EMP	2087	6	60
UNEMP	10	90	27	UNEMP	10	75	41
NLF	75	18	974	NLF	87	33	1639

Are the Hui-Walter assumptions plausible for these data?

### Hui-Walter Model Estimates of Classification Probabilities, fictitious numbers

Interview					
True Status	Observed Status				
	EMP	UNEMP	NLF		
EMP	98.6	1.4	0.0		
	(0.1)	(0.1)	(n/a)		
UNEMP	5.6	61.6	27.9		
	(15.2)	(11.1)	(5.3)		
NLF	2.6	0.0	97.4		
	(1.5)	(n/a)	(1.1)		

Standard errors in brackets

# Now: Three Indicator LC Model

#### DATA

	C=1		C=2	
	B=1	B=2	B=1	B=2
A=1	<i>n</i> <sub>111</sub>	<i>n</i> <sub>121</sub>	<i>n</i> <sub>112</sub>	<i>n</i> <sub>122</sub>
A=2	<i>n</i> <sub>211</sub>	<i>n</i> <sub>221</sub>	<i>n</i> <sub>212</sub>	<i>n</i> <sub>222</sub>

## **Classical Latent Class Model**

 Requires a minimum of three indicators of X for unrestricted model

– Say, A, B, C

Assumes local independence
 – i.e., P(A and B and C|X) = P(A|X)P(B|X)P(C|X)

• Notation  $\pi_{X=x} \text{ or } \pi_x^X = P(X = x),$   $\pi_{A=a|X=x} \text{ or } \pi_{a|x}^{A|X} = P(A = a \mid X = x)$ Also,  $\pi_x^X$ ,  $\pi_{abc}^{ABC}$ 

### Interpretations of Local Independence

$$\pi_{abc|x}^{ABC|X} = \pi_{a|x}^{A|X} \pi_{b|x}^{B|X} \pi_{c|x}^{C|X}$$

 $\Rightarrow$  classification errors are independent i.e., A, B, C represent independent selections from an individual's response distribution

We can refer to this model as  $\{A|XB|XC|X\}$  or  $\{AXBXCX\}$  (borrowed from log-linear modeling)

### Latent Class Model for Three Indicators

$$L(\boldsymbol{\pi} \mid \mathbf{n}) \approx \prod_{a} \prod_{b} \prod_{c} (\pi_{abc}^{ABC})^{n_{abc}}$$



### Path Diagram



What is not shown is as important as what is shown

### MLE Methods of Estimation

•  $\pi$ -probability model

– Our focus in this course

- log-linear model with latent variable
- modified path model

The latter two methods require knowledge of log-linear models which is not assumed in this course

# LCA Software include:

- Free software
  - *− ℓ*EM
  - "SAS" by PennState University
- Commercial software
  - Mplus
  - Latent Gold

Others: http://www.johnuebersax.com/stat/soft.htm

# Design of the NHSDA

- National, multistage, household survey
- 1994, 1995, 1996 data
- 43,825 total interviews
- data are weighted
- drug questions are repeated
- See Biemer and Wiesen (2002)
Three Indicators of Past-Year Marijuana Use

The Recency Question (Indicator A)

## How long has it been since you last used marijuana or hashish?

- A = 1 if either "Within the past 30 days" or "More than 30 days but within past 12 months"
- A = 2 if otherwise

Three Indicators of Past-Year Marijuana Use

The Frequency Question (Indicator B)

Now think about the past 12 months from your 12-month reference date through today. On how many days in the past 12 months did you use marijuana or hashish?

B = 1 if response is 1 or more days;

B = 2 if otherwise

Three Indicators of Past-Year Marijuana Use

#### The Composite Question (Indicator C ) 7 questions such as

- used in last 12 months?
- spent a great deal of time getting it, using it, or getting over its effects?
- used drug much more often or in larger amounts than intended?
- C = 1 if response is positive to any question;
- C = 2 if otherwise

#### Evidence of Reporting Error: Inconsistency in A, B, and C

Comparison	1994	1995	1996
A vs. B	1.35	1.48	1.61
A vs. C	4.80	2.14	2.48
B vs. C	4.96	2.31	2.69
A vs. B vs. C	5.55	2.96	3.39

#### 1995 Marijuana Data

	C=	=1	C=2		
	B=1	B=2	B=1	B=2	
A=1	1158	1158 11		3	
A=2	114	191	135	16064	

#### Introduction to *e*EM

<u> </u>	
man 3	
lat 1	
dim 2 2 2 2	
lab X A B C	
mod XAIXBIXCIX	
dat [1158 73 11 3 114 135 191 16064]	

#### Introduction to *et EM* (cont'd)

- man 3
- lat 1
- dim 2 2 2 2
- lab X A B C

\*Number of manifest variables \*Number of latent variables \*Dimensions \*Labels mod X A | X B | X C | \* Model

dat [1158 73 11 3 114 135 191 \*Data 16064]

#### Introduction to $\ell EM$ (cont'd)

🔏 Input - marij example.inp	
man 3	
lat 1	
dim 2 2 2 2	
lab X A B C	
mod XAIXBIXCIX	
sta A X [.9 .1 .1 .9]	
sta B X [.9 .1 .1 .9]	
sta C X [.9 .1 .1 .9]	
dat [1158 73 11 3 114 135 191 16064	]
npa	
nR2	
nla	
	-

#### Introduction to $\ell EM$ (cont'd)

```
*Number of manifest variables
man 3
                              *Number of latent variables
lat 1
                              *Dimensions
dim 2 2 2 2
                              *Labels
lab X A B C
                              *Model
mod X A|X B|X C|X
                              *Starting values for A|X
sta A|X [.9 .1 .1 .9]
                              *Starting values for B|X
sta B|X [.9 .1 .1 .9]
                              *Starting values for C|X
sta C|X [.9 .1 .1 .9]
                                                   *Data
dat [1158 73 11 3 114 135 191 16064]
                              *Output control
npa
                              *Output control
nR2
                              *Output control
nla
```

```
LEM: log-linear and event history analysis with missing data.
Developed by Jeroen Vermunt (c), Tilburg University, The
Netherlands.
Version 1.0b (September 18, 1997).
*** INPUT ***
  man 3
  lat 1
  dim 2 2 2 2
  lab X A B C
  mod X A | X B | X C | X
  sta A|X [.9 .1 .1 .9]
  sta B|X [.9 .1 .1 .9]
  sta C|X [.9 .1 .1 .9]
  dat [1158 73 11 3 114 135 191 16064]
  npa
  nR2
 46
hla
```

\*\*\* STATISTICS \*\*\*

Number of iterations = 8 Converge criterion = 0.0000008181 Seed random values = 5307

= 0.0000 (0.000)X-squared = 0.0000 (0.000)L-squared Cressie-Read = 0.0000 (0.000)Dissimilarity index = 0.0000Degrees of freedom = 0= -7371.28926Log-likelihood Number of parameters = 7 (+1) = 17749.0Sample size BIC(L-squared) = 0.0000AIC(L-squared) = 0.0000BIC(log-likelihood) = 14811.0671AIC(log-likelihood) = 14756.5785

WARNING: no information is provided on identification of parameters

#### \*\*\* FREQUENCIES \*\*\*

А	В	С	observed	estimated	std. res.
1	1	1	1158.000	1157.999	0.000
1	1	2	73.000	73.001	-0.000
1	2	1	11.000	11.001	-0.000
1	2	2	3.000	2.999	0.000
2	1	1	114.000	114.001	-0.000
2	1	2	135.000	134.999	0.000
2	2	1	191.000	191.000	0.000
2	2	2	16064.000	16064.000	-0.000

- \*\*\* (CONDITIONAL) PROBABILITIES \*\*\*
- \* P(X) \*
  - 1 0.0768 2 0.9232
- \* P(A|X) \*

1	1	0.9115
2	1	0.0885
1	2	0.0001
2	2	0.9999

\* P(B|X) \*

1	1	0.9906
2	1	0.0094
1	2	0.0079
2	2	0.9921

\* P(C|X) \*

1	1	0.9407
2	1	0.0593
1	2	0.0117
2	2	0.9883

#### **Fit Statistics**

- X-squared
- L-squared
- Dissimilarity index
- Degrees of freedom
- Log-likelihood
- Number of parameters
- Sample size
- BIC(L-squared)
- AIC(L-squared)

#### **Pearson Statistic**



 $\hat{m}_{abc}$  = model estimated frequency in cell (*a*, *b*, *c*)

Distributed approximately as a chisquare random variable if the model is true. Poor approximation if average cell size is less than 5.

#### Likelihood Ratio Chi-square Statistic



Distributed approximately as a chisquare random variable if the model is true. Poor approximation if average cell size is less than 5.

#### **Dissimilarity Index**

$$D = 2\sum_{abc} \frac{|n_{abc} - \hat{m}_{abc}|}{2n}$$

Smallest proportion of observations that would need to be reallocated to other cells to make the model fit perfectly. Should be less than 0.05 for a well-fitting model.

#### **Degrees of Freedom**

### df = number of cells number of estimated model parameters

#### Log-Likelihood

 $\log(L) = \sum n_{abc} \log(\hat{\pi}_{abc})$ abc

#### *BIC, AIC, BIC-L*<sup>2</sup> and *AIC-L*<sup>2</sup>

 $BIC = -2\log L + (\log n) \times npar$  $AIC = -2\log L + 2npar$  $BIC(L^{2}) = L^{2} - df \log(n)$  $AIC(L^{2}) = L^{2} - df 2$ 





#### ℓEM Input Code

```
* 4 manifest variables
   - 4
man
            * 1 latent variable
   1
lat
dim 2 2 2 2 2
              * S=1(Males)
=2(Females)
lab X S A B C
mod SX A|X B|X C|X
         * 16 records in data set
rec 16
        * last column is a count
rco
sta A|X [.9 .1 .1 .9]
sta B|X [.9 .1 .1 .9]
sta C|X [.9 .1 .1 .9]
```

#### *ℓ*EM Input Code

dat 1 1	[ 1 1	1	1	698	* *	Data Sex	orde: A	cing B	is C	Count
11111112222222222222222222222222222222	11222211112222	122112211221122	212121212121212	5 2 65 83 121 7495 460 30 6 1 49 52 70 8569]						

nla

Three Indicator Model with a Grouping Variable – Fully Saturated



#### Model: GX A|XG B|XG C|XG

# Three Indicator Model with a Grouping Variable – Fully Saturated

$$\pi_{gabc}^{GABC} = \sum_{x=1}^{2} \pi_{gx}^{GX} \pi_{a|xg}^{A|XG} \pi_{b|xg}^{B|XG} \pi_{c|xg}^{C|XG}$$

#### *ℓ*EM Input Code

man	4
lat	1
dim	2 2 2 2 2
lab	X S A B C
mod	SX A   SX B   SX C   SX Starting values for each
rec	16 group. Order consistent
rec rco	16 group. Order consistent with lab statement
rec rco sta	16       group. Order consistent with lab statement         A SX [.9 .1 .9 .1 .1 .9 .1 .1 .9 .1 .9]
rec rco sta sta	16       group. Order consistent with lab statement         A SX [.9 .1 .9 .1 .1 .9 .1 .9 .1 .9]         B SX [.9 .1 .9 .1 .1 .9 .1 .9]

#### ℓEM Results for Grouping Variable Models

Model	d.f	L <sup>2</sup>	р	D	BIC-L <sup>2</sup>
X A X B X C X	7	49.0	0.0	0.024	117.4799
SX A X B X C X	6	33.2	0.0	0.006	-25.4787
SX A XS B XS C XS	0	0	n/a	0	0

#### *ℓ*EM Estimates for Model 2

* 7	** (CONDITIONAL)	PROBABILITIES	* * *	
*	P(XS) *			
	1 1	0.0461	* P(B X) *	
	1 2 2 1	0.0308	1   1	0.9905
	2 2	0.4896		0.0095
*	P(A X) *		2   2	0.9921
	1   1	0.9110	* P(C X) *	
	2   1 1   2 2   2	0.0890 0.0001 0.9999	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	0.9402 0.0598 0.0117
				0.5005

#### Specifying Hui-Walter Model in ℓEM

#### Model: GX AX BX



#### Two Parameter, Two Indicator Model



#### A Model for Estimating Mode Effects

#### Study Design (Biemer, 2001)

Draw a sample and randomly divide. Assign one subsample to F-F; the other to CATI. Reinterview both by CATI.		Reinterview by CATI	
		B=1	B=2
Interview by Face to face	A=1	<i>п</i> <sub>F11</sub>	<b>n</b> <sub>F12</sub>
	A=2	<i>п</i> <sub>F21</sub>	п <sub>F22</sub>
Interview by CATI	A=1	<i>п</i> <sub>Т11</sub>	<i>n</i> <sub>T12</sub>
	A=2	<i>n</i> <sub>T21</sub>	n <sub>T22</sub>

#### Model Assumptions

- X denotes (latent) true characteristic
- *G*=1 denotes F-F sample
- G=2 denotes CATI sample  $\pi_{x|1}^{X|G} \neq \pi_{x|2}^{X|G}$
- CATI interview and CATI reinterview share a common mode effect; i.e.,

$$\pi_{2|12}^{A|XG} = \pi_{2|12}^{B|XG} = \pi_{2|11}^{B|XG} = \theta_T \text{ and}$$
$$\pi_{1|22}^{A|XG} = \pi_{1|22}^{B|XG} = \pi_{1|21}^{B|XG} = \phi_T$$

#### Summary of Model Parameters

 $\pi_{\varrho}^{G}, \pi_{111}^{X|G}, \pi_{12}^{X|G}, \theta_{T}, \theta_{F}, \phi_{T}, \phi_{F}$ 

Leaving 0 d.f. for testing fit

#### Does Anyone Smoke Inside the Home?

		Reinterview by CATI	
		R=1	R=2
Interview by Face to face	F=1	334	70
	F=2	29	1233
Interview by CATI	T=1	282	20
	T=2	9	931
## *ℓ*EM Input Code

\* Example: Does Anyone Smoke Inside the Home? lat 1 man 3 dim 2 2 2 2 lab X G A B mod G X | G A|GX eq2 B|GX eq2 sta X|G [.3 .4 .7 .6] des [0 1 0 2 3 0 4 0 0 2 0 2 4 0 4 0] dat [ 334 70 29 1233 282 20 9 931]

## *ℓ*EM Input Code

\* Example: Does Anyone Smoke Inside the Home? lat 1 man 3 dim 2 2 2 2 lab X G A B G denotes the split sample \* mod G X|G \* P(X|G) varies across split \* A|GX eq2 For G=1 A = ff and B=CATI \* B|GX eq2 For G=2 A & B are both CATI sta X|G [.3 .4 .7 .6] des [0 1 0 2 \*X=1 Int: Theta\_F and Theta\_T 3 0 4 0 \*X=2 Int: Phi\_F and Phi\_T 0 2 0 2 \*X=1 Reint: Theta T 4 0 4 0] \*X=2 Reint: Phi T dat [ 334 70 29 1233 282 20 9 931]

## *l***EM Fit Statistics**

=	4.1720	(0.0000)
=	4.2787	(0.0000)
=	4.1952	(0.0000)
=	0.0019	
=	0	
=	-4047.3	3498
=	7 (+1)	
=	2908.0	
=	4.2787	
=	4.2787	
=	8150.49	65
=	8108.67	00
		<pre>= 4.1720 = 4.2787 = 4.1952 = 0.0019 = 0 = -4047.3 = 7 (+1) = 2908.0 = 4.2787 = 4.2787 = 8150.49 = 8108.67</pre>

## **ℓ**EM Estimates

* P(G) * 1 2	0.5729 0.4271	* P(A XG) * 1   1 1 0.920 2   1 1 0.079 1   1 2 0.951 2   1 2 0.048 1   2 1 0.048	* .9207 .0793 .9516 .0484 .0413	P(B XG) 1   1 1 2   1 1 1 2	* 0.9516 0.0484 0.9516
* P(X G)	*		.9587	2   1 2	0.0484
1   1	0.2288	1 2 2 0	.0002	1   2 1	0.0002
1   2	0.2507	2   2 2 0	.9998	2   2 1	0.9998
2   1	0.7712			1   2 2	0.0002
2   2	0.7493			2   2 2	0.9998

Importance of Model Validity Depends on the Application

- In some applications, validity can be supported by ability to identify real questionnaire problems.
- In other applications, this type of validation may be quite difficult
- Further, LCA methodology is being pushed in the US to adjust the reported survey estimates for misclassification bias.
  - Unemployment rate
  - Expenditures
  - Total population size in a census

Analyzing Unequal Weighting and Clustered Samples ("Analysis of survey data")

- Option 1: Ignore the weights and clustering
  Pros: Often these do not change estimates of
  - classification error
  - Cons: can lead to false inference about classification error (see, for example, Patterson, Dayton and Graubard (2002))

Analyzing Unequal Weighting and Clustered Samples (cont'd)

- Option 2: Use software (Latent Gold, Mplus) that properly account for the weights
  - Pros: Unbiased estimates and asymptotically unbiased standard errors
  - Cons: Convergence problems in some cases