

Mplus Short Courses  
Topic 2

**Regression Analysis, Exploratory Factor Analysis,  
Confirmatory Factor Analysis, And Structural  
Equation Modeling For Categorical, Censored,  
And Count Outcomes**

Linda K. Muthén  
Bengt Muthén

Copyright © 2009 Muthén & Muthén

[www.statmodel.com](http://www.statmodel.com)

03/09/2011

1

**Table Of Contents**

General Latent Variable Modeling Framework	7
Analysis With Categorical Observed And Latent Variables	11
Categorical Observed Variables	13
Logit And Probit Regression	18
British Coal Miner Example	25
Logistic Regression And Adjusted Odds Ratios	39
Latent Response Variable Formulation Versus Probability Curve Formulation	46
Ordered Polytomous Regression	49
Alcohol Consumption Example	55
Unordered Polytomous Regression	58
Censored Regression	65
Count Regression	67
Poisson Regression	68
Negative Binomial Regression	70
Path Analysis With Categorical Outcomes	73
Occupational Destination Example	81

2

## Table Of Contents (Continued)

Categorical Observed And Continuous Latent Variables	86
Item Response Theory	89
Exploratory Factor Analysis	113
Practical Issues	129
CFA With Covariates	142
Antisocial Behavior Example	147
Multiple Group Analysis With Categorical Outcomes	167
Exploratory Structural Equation Modeling	172
Multi-Group EFA Of Male And Female Aggressive Behavior	185
Technical Issues For Weighted Least Squares Estimation	199
References	206

3

## Mplus Background

- Inefficient dissemination of statistical methods:
  - Many good methods contributions from biostatistics, psychometrics, etc are underutilized in practice
- Fragmented presentation of methods:
  - Technical descriptions in many different journals
  - Many different pieces of limited software
- Mplus: Integration of methods in one framework
  - Easy to use: Simple, non-technical language, graphics
  - Powerful: General modeling capabilities
- Mplus versions
  - V1: November 1998
  - V2: February 2001
  - V3: March 2004
  - V4: February 2006
  - V5: November 2007
  - V5.2: November 2008
- Mplus team: Linda & Bengt Muthén, Thuy Nguyen, Tihomir Asparouhov, Michelle Conn, Jean Maninger

4

## Statistical Analysis With Latent Variables A General Modeling Framework

### Statistical Concepts Captured By Latent Variables

#### Continuous Latent Variables

- Measurement errors
- Factors
- Random effects
- Frailties, liabilities
- Variance components
- Missing data

#### Categorical Latent Variables

- Latent classes
- Clusters
- Finite mixtures
- Missing data

5

## Statistical Analysis With Latent Variables A General Modeling Framework (Continued)

### Models That Use Latent Variables

#### Continuous Latent Variables

- Factor analysis models
- Structural equation models
- Growth curve models
- Multilevel models

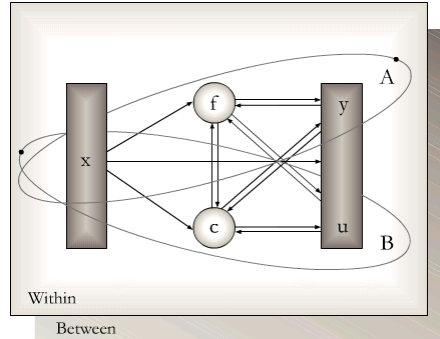
#### Categorical Latent Variables

- Latent class models
- Mixture models
- Discrete-time survival models
- Missing data models

Mplus integrates the statistical concepts captured by latent variables into a general modeling framework that includes not only all of the models listed above but also combinations and extensions of these models.

6

## General Latent Variable Modeling Framework



- Observed variables
  - x background variables (no model structure)
  - y continuous and censored outcome variables
  - u categorical (dichotomous, ordinal, nominal) and count outcome variables
- Latent variables
  - f continuous variables
    - interactions among f's
  - c categorical variables
    - multiple c's

7

## Mplus

Several programs in one

- Exploratory factor analysis
- Structural equation modeling
- Item response theory analysis
- Latent class analysis
- Latent transition analysis
- Survival analysis
- Growth modeling
- Multilevel analysis
- Complex survey data analysis
- Monte Carlo simulation

Fully integrated in the general latent variable framework

8

## Overview Of Mplus Courses

- **Topic 1.** August 20, 2009, Johns Hopkins University: Introductory - advanced factor analysis and structural equation modeling with continuous outcomes
- **Topic 2.** August 21, 2009, Johns Hopkins University: Introductory - advanced regression analysis, IRT, factor analysis and structural equation modeling with categorical, censored, and count outcomes
- **Topic 3.** March, 2010, Johns Hopkins University: Introductory and intermediate growth modeling
- **Topic 4.** March, 2010, Johns Hopkins University: Advanced growth modeling, survival analysis, and missing data analysis

9

## Overview Of Mplus Courses (Continued)

- **Topic 5.** August, 2010, Johns Hopkins University: Categorical latent variable modeling with cross-sectional data
- **Topic 6.** August 2010, Johns Hopkins University: Categorical latent variable modeling with longitudinal data
- **Topic 7.** March, 2011, Johns Hopkins University: Multilevel modeling of cross-sectional data
- **Topic 8.** March 2011, Johns Hopkins University: Multilevel modeling of longitudinal data

10

## **Analysis With Categorical Observed And Latent Variables**

11

## **Categorical Variable Modeling**

- Categorical observed variables
- Categorical observed variables, continuous latent variables
- Categorical observed variables, categorical latent variables

12

## Categorical Observed Variables

13

### Two Examples

**Alcohol Dependence And Gender In The NLSY**

	n	Not Dep	Dep	Prop	Odds (Prop/(1-Prop))
Female	4573	4317	256	0.056	0.059
Male	4603	3904	699	0.152	0.179
	9176	8221	955		

Odds Ratio = 0.179/0.059 = 3.019

Example wording: Males are three times more likely than females to be alcohol dependent.

**Colds And Vitamin C**

	n	No Cold	Cold	Prop	Odds
Placebo	140	109	31	0.221	0.284
Vitamin C	139	122	17	0.122	0.139

14

## Categorical Outcomes: Probability Concepts

- Probabilities:
 

	Alcohol Example			
– Joint: $P(u, x)$	Joint		Not Dep	Dep
– Marginal: $P(u)$				Conditional
– Conditional: $P(u   x)$				
	Female	.47	.03	.06
	Male	.43	.08	.15
	Marginal	.90	.11	
- Distributions:
  - Bernoulli:  $u = 0/1; E(u) = \pi$
  - Binomial: sum or prop. ( $u = 1$ ),  $E(prop.) = \pi$ ,  
 $V(prop.) = \pi(1 - \pi)/n, \hat{\pi} = prop$
  - Multinomial ( $\#parameters = \#cells - 1$ )
  - Independent multinomial (product multinomial)
  - Poisson

15

## Categorical Outcomes: Probability Concepts (Continued)

- Cross-product ratio (odds ratio):
 

		$u = 0$	$u = 1$
$x = 0$		$\pi_{00}$	$\pi_{01}$
$x = 1$		$\pi_{10}$	$\pi_{11}$

$$\pi_{00} \pi_{11} / (\pi_{01} \pi_{10}) = \frac{\pi_{11} / \pi_{10}}{\pi_{01} / \pi_{00}} =$$

$$P(u = 1, x = 1) / P(u = 0, x = 1) / P(u = 1, x = 0) / P(u = 0, x = 0)$$
- Tests:
  - Log odds ratio (approx. normal)
  - Test of proportions (approx. normal)
  - Pearson  $\chi^2 = \Sigma(O - E)^2 / E$  (e.g. independence)
  - Likelihood Ratio  $\chi^2 = 2 \Sigma O \log(O / E)$

16



## Further Readings On Categorical Variable Analysis

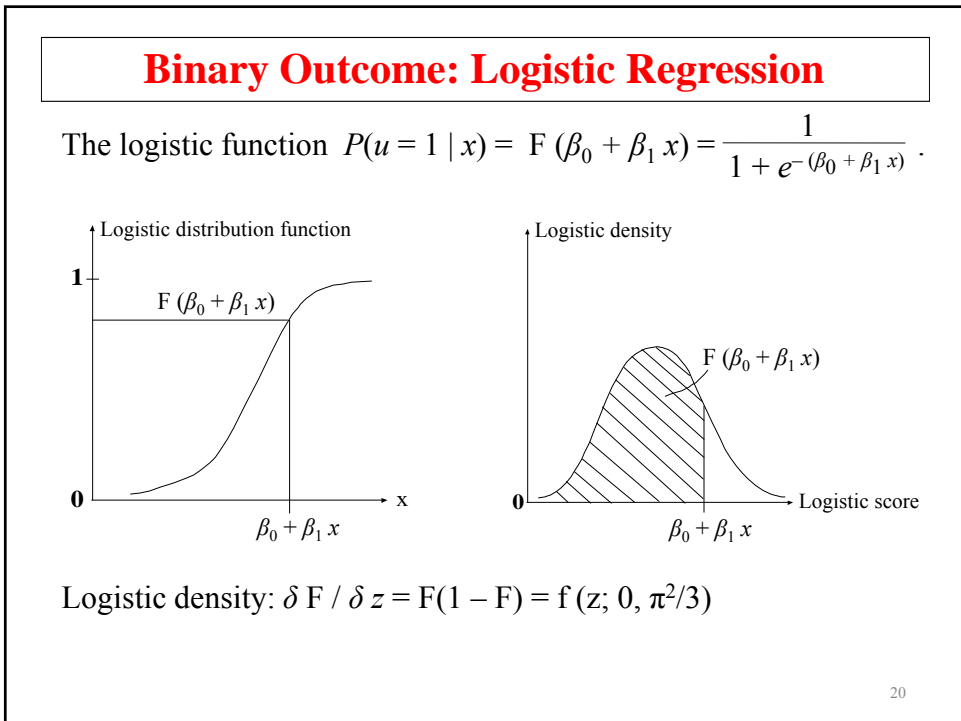
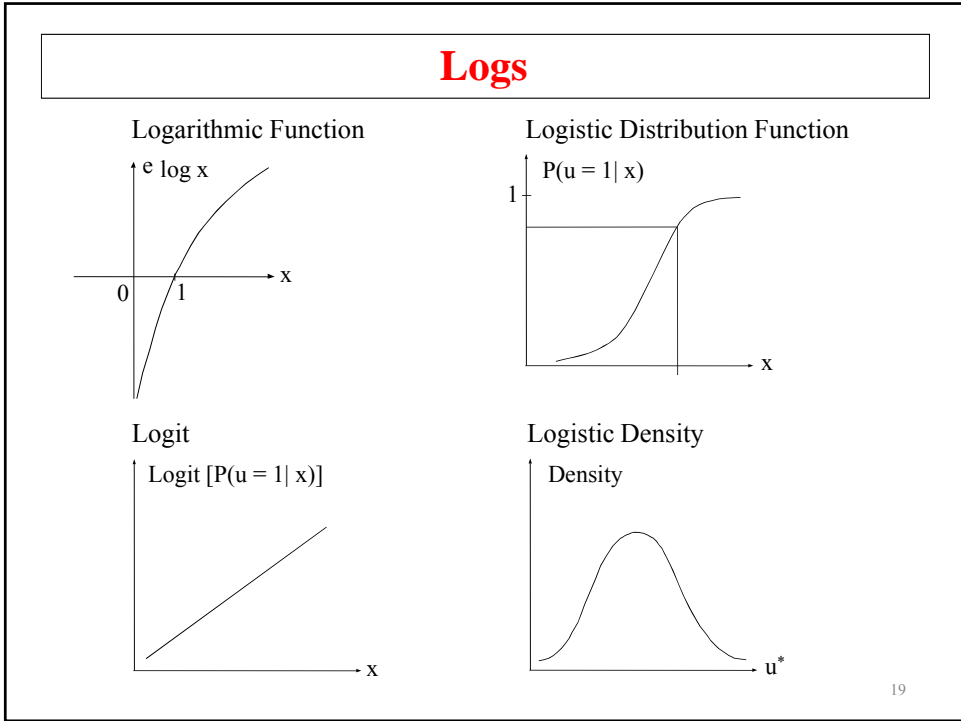
- Agresti, A. (2002). Categorical data analysis. Second edition. New York: John Wiley & Sons.
- Agresti, A. (1996). An introduction to categorical data analysis. New York: Wiley.
- Hosmer, D. W. & Lemeshow, S. (2000). Applied logistic regression. Second edition. New York: John Wiley & Sons.
- Long, S. (1997). Regression models for categorical and limited dependent variables. Thousand Oaks: Sage.

17

## Logit And Probit Regression

- Dichotomous outcome
- Adjusted log odds
- Ordered, polytomous outcome
- Unordered, polytomous outcome
- Multivariate categorical outcomes

18



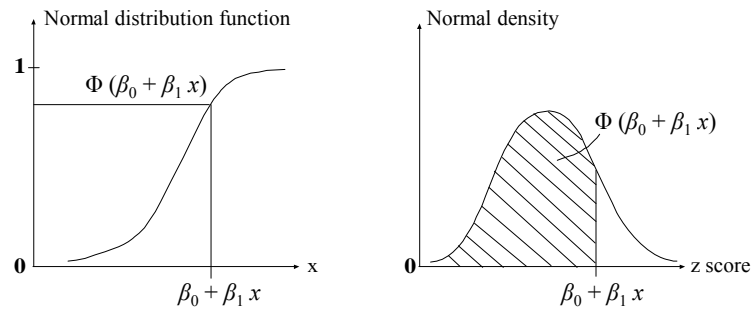
## Binary Outcome: Probit Regression

Probit regression considers

$$P(u = 1 | x) = \Phi(\beta_0 + \beta_1 x), \quad (60)$$

where  $\Phi$  is the standard normal distribution function. Using the inverse normal function  $\Phi^{-1}$ , gives a linear probit equation

$$\Phi^{-1}[P(u = 1 | x)] = \beta_0 + \beta_1 x. \quad (61)$$



21

## Interpreting Logit And Probit Coefficients

- Sign and significance
- Odds and odds ratios
- Probabilities

22

## Logistic Regression And Log Odds

$$\begin{aligned} \text{Odds } (u = 1 | x) &= P(u = 1 | x) / P(u = 0 | x) \\ &= P(u = 1 | x) / (1 - P(u = 1 | x)). \end{aligned}$$

The logistic function

$$P(u = 1 | x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

gives a log odds linear in  $x$ ,

$$\text{logit} = \log [\text{odds } (u = 1 | x)] = \log [P(u = 1 | x) / (1 - P(u = 1 | x))]$$

$$= \log \left[ \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}} / \left( 1 - \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}} \right) \right]$$

$$= \log \left[ \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}} * \frac{1 + e^{-(\beta_0 + \beta_1 x)}}{e^{-(\beta_0 + \beta_1 x)}} \right]$$

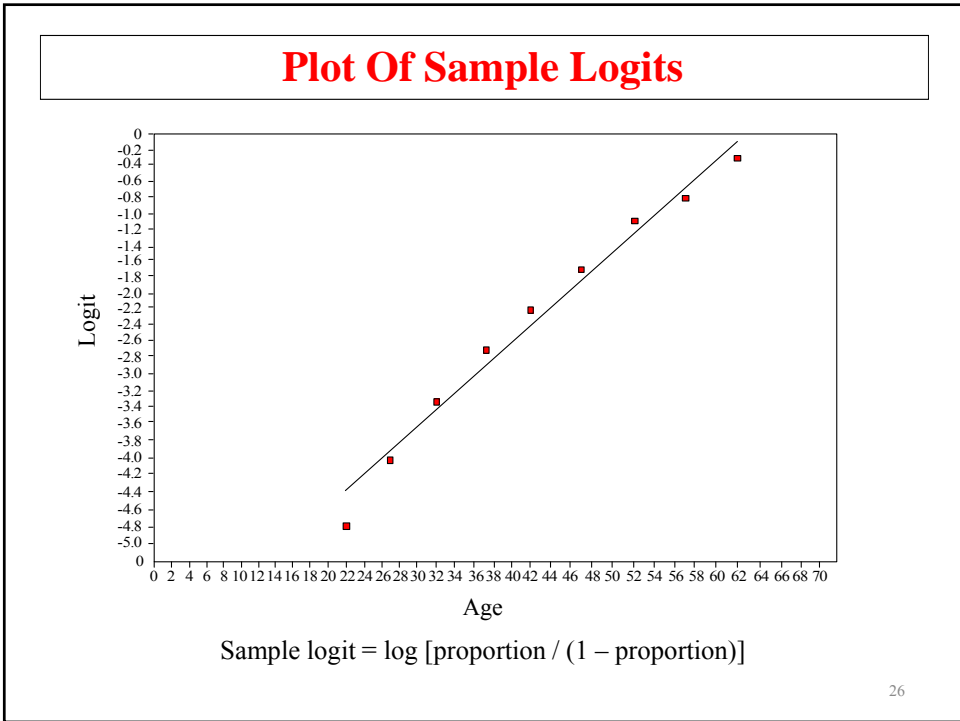
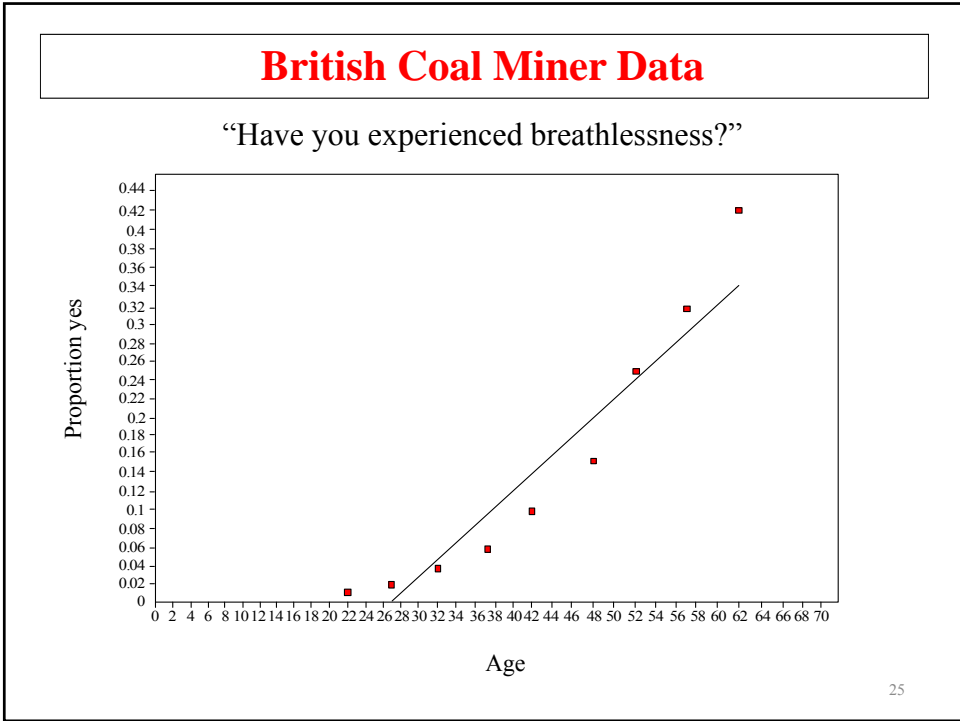
$$= \log \left[ e^{(\beta_0 + \beta_1 x)} \right] = \beta_0 + \beta_1 x$$

23

## Logistic Regression And Log Odds (Continued)

- $\text{logit} = \log \text{odds} = \beta_0 + \beta_1 x$
- When  $x$  changes one unit, the *logit* (*log odds*) changes  $\beta_1$  units
- When  $x$  changes one unit, the *odds* changes  $e^{\beta_1}$  units

24



**British Coal Miner Data (Continued)**

<i>Age (x)</i>	<i>N</i>	<i>N Yes</i>	<i>Proportion Yes</i>	<i>OLS Estimated Probability</i>	<i>Logit Estimated Probability</i>	<i>Probit Estimated Probability</i>
22	1,952	16	0.008	-0.053	0.013	0.009
27	1,791	32	0.018	-0.004	0.022	0.018
32	2,113	73	0.035	0.045	0.036	0.034
37	2,783	169	0.061	0.094	0.059	0.060
42	2,274	223	0.098	0.143	0.095	0.100
47	2,393	357	0.149	0.192	0.148	0.156
52	2,090	521	0.249	0.241	0.225	0.231
57	1,750	558	0.319	0.290	0.327	0.322
62	1,136	478	0.421	0.339	0.448	0.425
	18,282	2,427	0.130			

SOURCE: Ashford & Sowden (1970), Muthén (1993)

Logit model:  $\chi^2_{LRT}(7) = 17.13$  ( $p > 0.01$ )

Probit model:  $\chi^2_{LRT}(7) = 5.19$

27

**Coal Miner Data**

<i>x</i>	<i>u</i>	<i>w</i>
22	0	1936
22	1	16
27	0	1759
27	1	32
32	0	2040
32	1	73
37	0	2614
37	1	169
42	0	2051
42	1	223
47	0	2036
47	1	357
52	0	1569
52	1	521
57	0	1192
57	1	558
62	0	658
62	1	478

28

## Mplus Input For Categorical Outcomes

- Specifying dependent variables as categorical – use the CATEGORICAL option

CATEGORICAL ARE u1 u2 u3;

- Thresholds used instead of intercepts – only different in sign
- Referring to thresholds in the model – use \$ number added to a variable name – the number of thresholds is equal to the number of categories minus 1

u1\$1 refers to threshold 1 of u1

u1\$2 refers to threshold 2 of u1

29

## Mplus Input For Categorical Outcomes (Continued)

u2\$1 refers to threshold 1 of u2

u2\$2 refers to threshold 2 of u2

u2\$3 refers to threshold 3 of u2

u3\$1 refers to threshold 1 of u3

- Referring to scale factors – use { } to refer to scale factors

{u1@1 u2 u3};

30

## Input For Logistic Regression Of Coal Miner Data

```

TITLE:      Logistic regression of coal miner data
DATA:      FILE = coalminer.dat;
VARIABLE:  NAMES = x u w;
           CATEGORICAL = u;
           FREQWEIGHT = w;
DEFINE:    x = x/10;
ANALYSIS:  ESTIMATOR = ML;
MODEL:     u ON x;
OUTPUT:    TECH1 SAMPSTAT STANDARDIZED;

```

31

## Input For Probit Regression Of Coal Miner Data

```

TITLE:      Probit regression of coal miner data
DATA:      FILE = coalminer.dat;
VARIABLE:  NAMES = x u w;
           CATEGORICAL = u;
           FREQWEIGHT = w;
DEFINE:    x = x/10;
MODEL:     u ON x;
OUTPUT:    TECH1 SAMPSTAT STANDARDIZED;

```

32



## Output Excerpts Logistic Regression Of Coal Miner Data

### Model Results

	Estimates	S.E.	Est./S.E.	Std	StdYX
U ON X	1.025	0.025	41.758	1.025	0.556
Thresholds U\$1	6.564	0.124	52.873		

$$\text{Odds: } e^{1.025} = 2.79$$

As  $x$  increases 1 unit (10 years), the odds of breathlessness increases 2.79

33

## Estimated Logistic Regression Probabilities For Coal Miner Data

$$P(u=1|x) = \frac{1}{1+e^{-L}},$$

$$\text{where } L = -6.564 + 1.025 \times x$$

$$\text{For } x = 6.2 \text{ (age 62)}$$

$$L = -6.564 + 1.025 \times 6.2 = -0.209$$

$$P(u=1|\text{age 62}) = \frac{1}{1+e^{0.209}} = 0.448$$

34

## Output Excerpts Probit Regression Of Coal Miner Data

### Model Results

	Estimates	S.E.	Est./S.E.	Std	StdYX
U					
ON					
X	0.548	0.013	43.075	0.548	0.545
Thresholds					
U\$1	3.581	0.062	57.866	3.581	3.581

### R-Square

Observed Variable	Residual Variance	R-Square
U	1.000	0.297

35

## Estimated Probit Regression Probabilities For Coal Miner Data

$$\begin{aligned}
 P(u = 1 | x = 62) &= \Phi(\hat{\beta}_0 + \hat{\beta}_1 x) \\
 &= 1 - \Phi(\hat{\tau} - \hat{\beta}_1 x) \\
 &= \Phi(-\hat{\tau} + \hat{\beta}_1 x).
 \end{aligned}$$

$$\Phi(-3.581 + 0.548 * 6.2) = \Phi(-0.1834) \approx 0.427$$

Note:  $\text{logit } \hat{\beta} \approx \text{probit } \hat{\beta} * c$   
 where  $c = \sqrt{\pi^2 / 3} = 1.81$

36

**Categorical Outcomes: Logit And Probit Regression  
With One Binary And One Continuous X**

$$P(u = 1 | x_1, x_2) = F[\beta_0 + \beta_1 x_1 + \beta_2 x_2], \quad (22)$$

$P(u = 0 | x_1, x_2) = 1 - P[u = 1 | x_1, x_2]$ , where  $F[z]$  is either the standard normal ( $\Phi[z]$ ) or logistic ( $1/[1 + e^{-z}]$ ) distribution function.

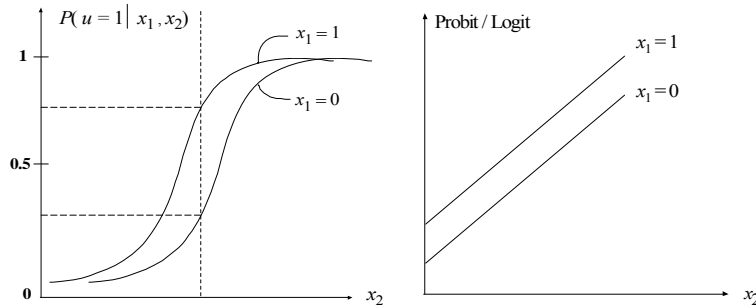
**Example:** Lung cancer and smoking among coal miners

- $u$  lung cancer ( $u = 1$ ) or not ( $u = 0$ )
- $x_1$  smoker ( $x_1 = 1$ ), non-smoker ( $x_1 = 0$ )
- $x_2$  years spent in coal mine

37

**Categorical Outcomes: Logit And Probit Regression  
With One Binary And One Continuous X**

$$P(u = 1 | x_1, x_2) = F[\beta_0 + \beta_1 x_1 + \beta_2 x_2], \quad (22)$$



38

## Logistic Regression And Adjusted Odds Ratios

Binary  $u$  variable regression on a binary  $x_1$  variable and a continuous  $x_2$  variable:

$$P(u = 1 | x_1, x_2) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}}, \quad (62)$$

which implies

$$\log \text{odds} = \text{logit} [P(u = 1 | x_1, x_2)] = \beta_0 + \beta_1 x_1 + \beta_2 x_2. \quad (63)$$

This gives

$$\log \text{odds}_{\{x_1=0\}} = \text{logit} [P(u = 1 | x_1 = 0, x_2)] = \beta_0 + \beta_2 x_2, \quad (64)$$

and

$$\log \text{odds}_{\{x_1=1\}} = \text{logit} [P(u = 1 | x_1 = 1, x_2)] = \beta_0 + \beta_1 + \beta_2 x_2. \quad (65)$$

39

## Logistic Regression And Adjusted Odds Ratios (Continued)

The log odds ratio for  $u$  and  $x_1$  adjusted for  $x_2$  is

$$\log OR = \log \left[ \frac{\text{odds}_1}{\text{odds}_0} \right] = \log \text{odds}_1 - \log \text{odds}_0 = \beta_1 \quad (66)$$

so that  $OR = \exp(\beta_1)$ , constant for all values of  $x_2$ . If an interaction term for  $x_1$  and  $x_2$  is introduced, the constancy of the OR no longer holds.

Example wording:

“The odds of lung cancer adjusted for years is OR times higher for smokers than for nonsmokers”

“The odds ratio adjusted for years is OR”

40

## Analysis Of NLSY Data: Odds Ratios For Alcohol Dependence And Gender

### Adjusting for Age First Started Drinking (n=9176)

Observed Frequencies, Proportions, and Odds Ratios					
Age 1st	Frequency		Proportion Dependent		
	Female	Male	Female	Male	OR
12 or <	85	223	.071	.233	3.98
13	105	180	.133	.256	2.24
14	198	308	.086	.253	3.60
15	331	534	.106	.185	1.91
16	800	990	.079	.152	2.09
17	725	777	.070	.170	2.72
18 or >	2329	1591	.030	.089	3.16

41

## Analysis Of NLSY Data: Odds Ratios For Alcohol Dependence And Gender (Continued)

Estimated Probabilities and Odds Ratios						
Age 1st	Logit			Probit		
	Female	Male	OR	Female	Male	OR
12 or <	.141	.304	2.66	.152	.298	2.37
13	.117	.260	2.66	.125	.257	2.42
14	.096	.220	2.66	.102	.220	2.48
15	.078	.185	2.66	.082	.186	2.55
16	.064	.154	2.66	.065	.155	2.63
17	.052	.127	2.66	.051	.128	2.72
18 or >	.042	.105	2.66	.040	.104	2.82

Logit model:  $\chi^2_p(12) = 54.2$

Probit model:  $\chi^2_p(12) = 46.8$

42

## Analysis Of NLSY Data: Odds Ratios For Alcohol Dependence And Gender (Continued)

### Dependence on Gender and Age First Started Drinking

	Logit Regression				Probit Regression				Unstd. Coeff Rescaled To Logit
	Unstd. Coeff.	s.e.	t	Std.	Unstd. Coeff.	s.e.	t	Std.	
Intercept	0.84	.32	2.6		-0.42	.18	-2.4		
Male	0.98	.08	12.7	0.51	0.50	.04	13.1	0.48	0.91
Age 1st	-0.22	.02	-11.6	-0.19	-0.12	.01	-11.0	-0.19	-0.22
R <sup>2</sup>	0.12				0.08				

$$OR = e^{0.98} = 2.66$$

$$\text{logit } \beta \approx \text{probit } \beta * c$$

$$\text{where } c = \sqrt{\pi^2 / 3} = 1.81$$

43

## NELS 88

**Table 2.2** – Odds ratios of eighth-grade students in 1988 performing below basic levels of reading and mathematics in 1988 and dropping out of school, 1988 to 1990, by basic demographics

Variable	Below basic mathematics	Below basic reading	Dropped out
<b>Sex</b>			
Female vs. male	0.81*	0.73**	0.92
<b>Race — ethnicity</b>			
Asian vs. white	0.82	1.42**	0.59
Hispanic vs. white	2.09**	2.29**	2.01**
Black vs. white	2.23**	2.64**	2.23**
Native American vs. white	2.43**	3.50**	2.50**
<b>Socioeconomic status</b>			
Low vs. middle	1.90**	1.91**	3.95**
High vs. middle	0.46**	0.41**	0.39*

SOURCE: U.S. Department of Education, National Center for Education Statistics, National Education Longitudinal Study of 1988 (NELS:88), "Base Year and First Follow-Up surveys.

44

## NELS 88

**Table 2.3** – Adjusted odds ratios of eighth-grade students in 1988 performing below basic levels of reading and mathematics in 1988 and dropping out of school, 1988 to 1990, by basic demographics

Variable	Below basic mathematics	Below basic reading	Dropped out
Sex			
Female vs. male	0.77**	0.70**	0.86
Race — ethnicity			
Asian vs. white	0.84	1.46**	0.60
Hispanic vs. white	1.60**	1.74**	1.12
Black vs. white	1.77**	2.09**	1.45
Native American vs. white	2.02**	2.87**	1.64
Socioeconomic status			
Low vs. middle	1.68**	1.66**	3.74**
High vs. middle	0.49**	0.44**	0.41*

45

## Latent Response Variable Formulation Versus Probability Curve Formulation

Probability curve formulation in the binary  $u$  case:

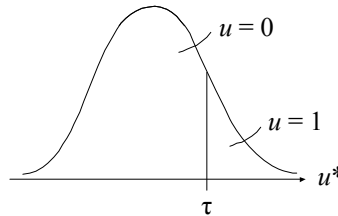
$$P(u = 1 | x) = F(\beta_0 + \beta_1 x), \tag{67}$$

where  $F$  is the standard normal or logistic distribution function.

Latent response variable formulation defines a threshold  $\tau$  on a continuous  $u^*$  variable so that  $u = 1$  is observed when  $u^*$  exceeds  $\tau$  while otherwise  $u = 0$  is observed,

$$u^* = \gamma x + \delta, \tag{68}$$

where  $\delta \sim N(0, V(\delta))$ .



46

**Latent Response Variable Formulation Versus Probability Curve Formulation (Continued)**

$$P(u = 1 | x) = P(u^* > \tau | x) = 1 - P(u^* \leq \tau | x) = \tag{69}$$

$$= 1 - \Phi[(\tau - \gamma x) V(\delta)^{-1/2}] = \Phi[-\tau + \gamma x) V(\delta)^{-1/2}]. \tag{70}$$

Standardizing to  $V(\delta) = 1$  this defines a probit model with intercept  $(\beta_0) = -\tau$  and slope  $(\beta_1) = \gamma$ .

Alternatively, a logistic density may be assumed for  $\delta$ ,

$$f[\delta ; 0, \pi^2/3] = dF/d\delta = F(1 - F), \tag{71}$$

where in this case  $F$  is the logistic distribution function  $1/(1 + e^{-\delta})$ .

47

**Latent Response Variable Formulation: R<sup>2</sup>, Standardization, And Effects On Probabilities**

$$u^* = \gamma x + \delta$$

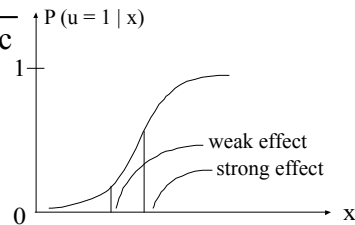
- $R^2(u^*) = \gamma^2 V(x) / (\gamma^2 V(x) + c)$ , where  $c = 1$  for probit and  $\pi^2 / 3$  for logit (McKelvey & Zavoina, 1975)

- Standardized  $\gamma$  refers to the effect of  $x$  on  $u^*$ ,

$$\hat{\gamma}_s = \hat{\gamma} SD(x) / SD(u^*),$$

$$SD(u^*) = \sqrt{\hat{\gamma}^2 V(x) + c}$$

- Effect of  $x$  on  $P(u = 1)$  depends on  $x$  value

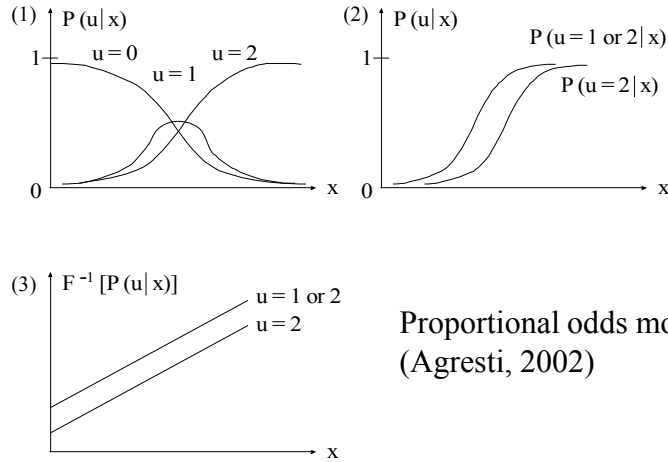


48



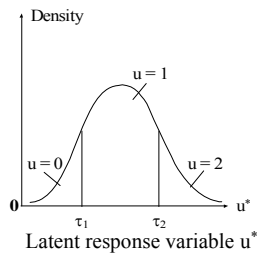
## Modeling With An Ordered Polytomous u Outcome

u polytomous with 3 categories



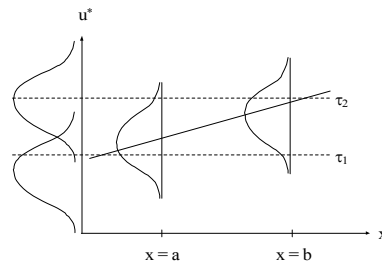
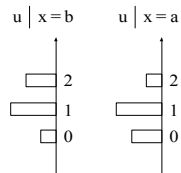
49

## Ordered Polytomous Outcome Using A Latent Response Variable Formulation



Latent response variable regression:

$$u_i^* = \gamma x_i + \delta_i$$



50

## Ordered Polytomous Outcome Using A Latent Response Variable Formulation (Continued)

A categorical variable  $u$  with  $C$  ordered categories,

$$u = c, \text{ if } \tau_{j,c} < u^* \leq \tau_{j,c+1} \quad (72)$$

for categories  $c = 0, 1, 2, \dots, C - 1$  and  $\tau_0 = -\infty, \tau_C = \infty$ .

Example: a single  $x$  variable and a  $u$  variable with three categories. Two threshold parameters,  $\tau_1$  and  $\tau_2$ .

Probit:

$$u^* = \gamma x + \delta, \text{ with } \delta \text{ normal} \quad (73)$$

$$P(u = 0 | x) = \Phi(\tau_1 - \gamma x), \quad (74)$$

$$P(u = 1 | x) = \Phi(\tau_2 - \gamma x) - \Phi(\tau_1 - \gamma x), \quad (75)$$

$$P(u = 2 | x) = 1 - \Phi(\tau_2 - \gamma x) = \Phi(-\tau_2 + \gamma x). \quad (76)$$

51

## Ordered Polytomous Outcome Using A Latent Response Variable Formulation (Continued)

$$P(u = 1 \text{ or } 2 | x) = P(u = 1 | x) + P(u = 2 | x) \quad (77)$$

$$= 1 - \Phi(\tau_1 - \gamma x) \quad (78)$$

$$= \Phi(-\tau_1 + \gamma x) \quad (79)$$

$$= 1 - P(u = 0 | x), \quad (80)$$

that is, a linear probit for,

$$P(u = 2 | x) = \Phi(-\tau_2 + \gamma x), \quad (81)$$

$$P(u = 1 \text{ or } 2 | x) = \Phi(-\tau_1 + \gamma x). \quad (82)$$

Note: same slope  $\gamma$ , so parallel probability curves

52

### Logit For Ordered Categorical Outcome

$$P(u = 2 | x) = \frac{1}{1 + e^{-(\beta_2 + \beta x)}}, \quad (83)$$

$$P(u = 1 \text{ or } 2 | x) = \frac{1}{1 + e^{-(\beta_1 + \beta x)}}. \quad (84)$$

Log odds for each of these two events is a linear expression,

$$\text{logit} [P(u = 2 | x)] = \quad (85)$$

$$= \log[P(u = 2 | x) / (1 - P(u = 2 | x))] = \beta_2 + \beta x, \quad (86)$$

$$\text{logit} [P(u = 1 \text{ or } 2 | x)] = \quad (87)$$

$$= \log[P(u = 1 \text{ or } 2 | x) / (1 - P(u = 1 \text{ or } 2 | x))] = \beta_1 + \beta x. \quad (88)$$

Note: same slope  $\beta$ , so parallel probability curves

53

### Logit For Ordered Categorical Outcome (Continued)

When  $x$  is a 0/1 variable,

$$\text{logit} [P(u = 2 | x = 1)] - \text{logit} [P(u = 2 | x = 0)] = \beta \quad (89)$$

$$\text{logit} [P(u = 1 \text{ or } 2 | x = 1)] - \text{logit} [P(u = 1 \text{ or } 2 | x = 0)] = \beta \quad (90)$$

showing that the ordered polytomous logistic regression model has constant odds ratios for these different outcomes.

54

## Alcohol Consumption: Ordered Polytomous Regression

$u$ : “On the days that you drink, how many drinks do you have per day, on the average?”

Ordinal $u$ : (“Alameda Scoring”)	$x$ 's: Age: whole years 20 – 64
0 non-drinker	Income: 1 $\leq$ \$4,999
1 1-2 drinks per day	2 \$5,000 – \$9,999
2 3-4 drinks per day	3 \$10,000 – \$14,999
3 5 or more drinks per day	4 \$15,000 – \$24,999
	5 $\geq$ \$25,000

$N = 713$  Males with regular physical activity levels

Source: Golden (1982), Muthén (1993)

55

## Alcohol Consumption: Ordered Polytomous Regression (Continued)

$$P(u = 0 | x) = \Phi(\tau_1 - \gamma'x) \quad (11)$$

$$P(u = 1 | x) = \Phi(\tau_2 - \gamma'x) - \Phi(\tau_1 - \gamma'x),$$

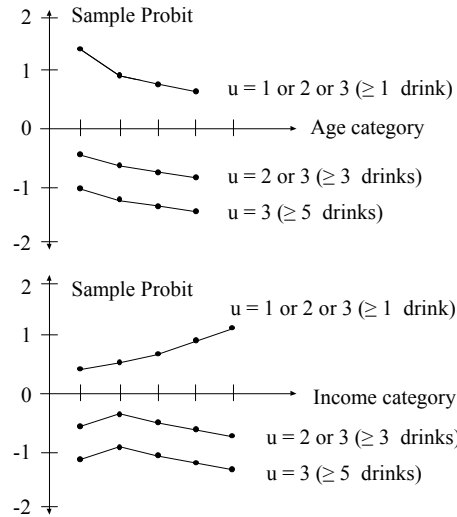
$$P(u = 2 | x) = \Phi(\tau_3 - \gamma'x) - \Phi(\tau_2 - \gamma'x),$$

$$P(u = 3 | x) = \Phi(-\tau_3 + \gamma'x).$$

Ordered  $u$  gives a single slope

56

## Alcohol Consumption: Ordered Polytomous Regression (Continued)



57

## Polytomous Outcome: Unordered Case

Multinomial logistic regression:

$$P(u_i = c | x_i) = \frac{e^{\beta_{0c} + \beta_{1c} x_i}}{\sum_{k=1}^K e^{\beta_{0k} + \beta_{1k} x_i}}, \quad (91)$$

for  $c = 1, 2, \dots, K$ , where we standardize to

$$\beta_{0K} = 0, \quad (92)$$

$$\beta_{1K} = 0, \quad (93)$$

which gives the log odds

$$\log[P(u_i = c | x_i) / P(u_i = K | x_i)] = \beta_{0c} + \beta_{1c} x_i, \quad (94)$$

for  $c = 1, 2, \dots, K - 1$ .

58

## Multinomial Logistic Regression Special Case Of K = 2

$$\begin{aligned}
 P(u_i = 1 | x_i) &= \frac{e^{\beta_{01} + \beta_{11} x_i}}{e^{\beta_{01} + \beta_{11} x_i} + 1} \\
 &= \frac{e^{-(\beta_{01} + \beta_{11} x_i)}}{e^{-(\beta_{01} + \beta_{11} x_i)}} * \frac{e^{\beta_{01} + \beta_{11} x_i}}{e^{\beta_{01} + \beta_{11} x_i} + 1} \\
 &= \frac{1}{1 + e^{-(\beta_{01} + \beta_{11} x_i)}}
 \end{aligned}$$

which is the standard logistic regression for a binary outcome.

59

## Input For Multinomial Logistic Regression

```

TITLE:          multinomial logistic regression
DATA:          FILE = nlsy.dat;
VARIABLE:      NAMES = u x1-x3;
               NOMINAL = u;
MODEL:        u ON x1-x3;

```

60

**Output Excerpts**  
**Multinomial Logistic Regression:**  
**4 Categories Of ASB In The NLSY**

		Estimates	S.E.	Est./S.E.
U#1	ON			
	AGE94	-.285	.028	-10.045
	MALE	2.578	.151	17.086
	BLACK	.158	.139	1.141
U#2	ON			
	AGE94	.069	.022	3.182
	MALE	.187	.110	1.702
	BLACK	-.606	.139	-4.357
U#3	ON			
	AGE94	-.317	.028	-11.311
	MALE	1.459	.101	14.431
	BLACK	.999	.117	8.513
Intercepts				
	U#1	-1.822	.174	-10.485
	U#2	-.748	.103	-7.258
	U#3	-.324	.125	-2.600

61

**Estimated Probabilities**  
**For Multinomial Logistic Regression:**  
**4 Categories Of ASB In The NLSY**

**Example 1: x's = 0**

	exp	probability = exp/sum
log odds (u=1) = -1.822	0.162	0.069
log odds (u=2) = -0.748	0.473	0.201
log odds (u=3) = -0.324	0.723	0.307
log odds (u=4) = 0	1.0	0.424
sum	2.358	1.001

62

## Estimated Probabilities For Multinomial Logistic Regression: 4 Categories Of ASB In The NLSY (Continued)

**Example 2: x = 1, 1, 1**

$$\text{log odds (u=1)} = -1.822 + (-0.285*1) + (2.578*1) + (0.158*1) = 0.629$$

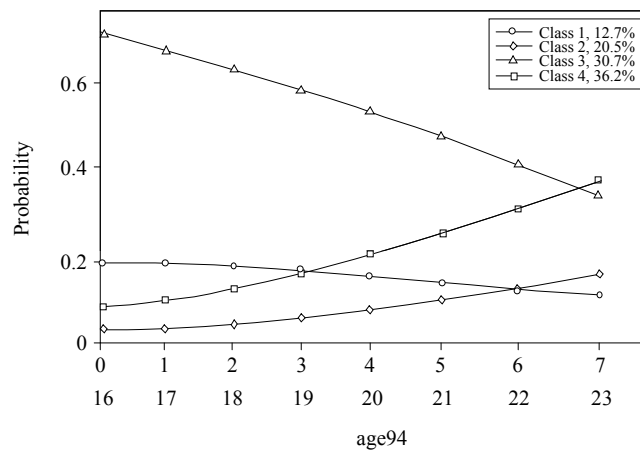
$$\text{log odds (u=2)} = -0.748 + 0.069*1 + 0.187*1 + (-0.606*1) = -1.098$$

$$\text{log odds (u=3)} = -0.324 + (-0.317*1) + 1.459*1 + 0.999*1 = 1.817$$

	exp	probability = exp/sum
log odds (u=1) = 0.629	1.876	0.200
log odds (u=2) = -1.098	0.334	0.036
log odds (u=3) = 1.817	6.153	0.657
log odds (u=4) = 0	1.0	0.107
sum	9.363	1.000

63

## Estimated Probabilities For Multinomial Logistic Regression: 4 Categories Of ASB In The NLSY (Continued)



64



## Censored-Normal (Tobit) Regression

$$y^* = \pi_0 + \pi x + \delta \quad V(\delta) \text{ identifiable}$$

Continuous – unlimited:  $y = y^*$

$$\text{Continuous-censored: } y = \begin{cases} c_L, & \text{if } y^* \leq c_L \\ y^*, & \text{if } c_L < y^* < c_U \\ c_U, & \text{if } y^* \geq c_U \end{cases}$$

Censoring from below,  $c_L = 0, c_U = \infty$ :

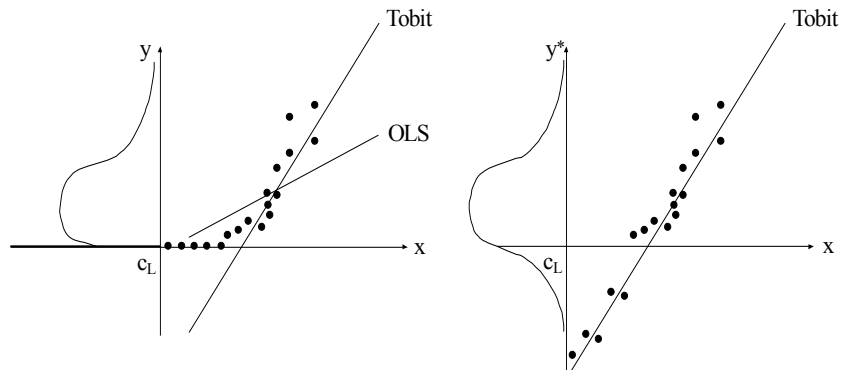
$$P(y > 0 | x) = F \left( \frac{\pi_0 + \pi x}{\sqrt{V(\delta)}} \right) \quad (\text{Probit Regression})$$

$$E(y | y > 0, x) = \pi_0 + \pi x + f/F \sqrt{V(\delta)}$$

**Classical Tobit**

65

## OLS v. Tobit Regression For Censored y But Normal y\*



66

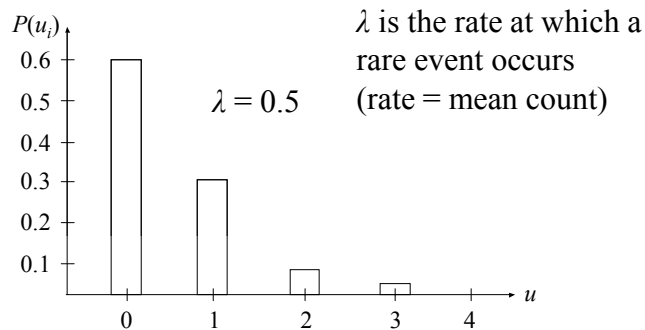
## Regression With A Count Dependent Variable

67

## Poisson Regression

A Poisson distribution for a count variable  $u_i$  has

$$P(u_i = r) = \frac{\lambda_i^r e^{-\lambda_i}}{r!}, \text{ where } u_i = 0, 1, 2, \dots$$



Regression equation for the log rate:

$$e^{\log \lambda_i} = \ln \lambda_i = \beta_0 + \beta_1 x_i$$

68

## Zero-Inflated Poisson (ZIP) Regression

A Poisson variable has mean = variance.

Data often have variance > mean due to preponderance of zeros.

$\pi = P$  (being in the zero class where only  $u = 0$  is seen)

$1 - \pi = P$  (not being in the zero class with  $u$  following a Poisson distribution)

A mixture at zero:      ZIP mean count:      ZIP variance of count:

$$P(u = 0) = \pi + \underbrace{(1 - \pi) e^{-\lambda}}_{\text{Poisson part}} \quad \lambda (1 - \pi) \quad \lambda (1 - \pi) (1 + \lambda * \pi)$$

The ZIP model implies two regressions:

$$\text{logit}(\pi_i) = \gamma_0 + \gamma_1 x_i,$$

$$\ln \lambda_i = \beta_0 + \beta_1 x_i$$

69

## Negative Binomial Regression

Unobserved heterogeneity  $\varepsilon_i$  is added to the Poisson model

$$\ln \lambda_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \text{ where } \exp(\varepsilon) \sim \Gamma$$

Poisson assumes

Negative binomial assumes

$$E(u_i | x_i) = \lambda_i$$

$$E(u_i | x_i) = \lambda_i$$

$$V(u_i | x_i) = \lambda_i$$

$$V(u_i | x_i) = \lambda_i (1 + \lambda_i \alpha)$$

NB with  $\alpha = 0$  gives Poisson. When the dispersion parameter  $\alpha > 0$ , the NB model gives substantially higher probability for low counts and somewhat higher probability for high counts than Poisson.

Further variations are zero-inflated NB and zero-truncated NB (hurdle model or two-part model).

70

## Mplus Specifications

Variable command	Type of dependent variable	Variance/residual variance
CATEGORICAL = u;	Binary, ordered polytomous	No
NOMINAL = u;	Unordered, polytomous (nominal)	No
CENSORED = y (b); = y (a);	Censored normal (Tobit) Censored from below or above	Yes
COUNT = u; u (p);	Poisson	No
= u (i); u (pi);	Zero-inflated Poisson	No
= u (nb);	Negative binomial	
= u (nbi);	Zero-inflated negative binomial	
= u (nbt);	Zero-truncated negative binomial	
= u (nbh);	Negative binomial hurdle	

71

## Further Readings On Censored and Count Regressions

- Hilbe, J. M. (2007). Negative binomial regression. Cambridge, UK: Cambridge University Press.
- Lambert, D. (1992). Zero-inflated Poisson regression, with an application to defects in manufacturing. Technometrics, 34, 1-13.
- Long, S. (1997). Regression models for categorical and limited dependent variables. Thousand Oaks: Sage.
- Maddala, G.S. (1983). Limited-dependent and qualitative variables in econometrics. Cambridge: Cambridge University Press.
- Tobin, J (1958). Estimation of relationships for limited dependent variables. Econometrica, 26, 24-36.

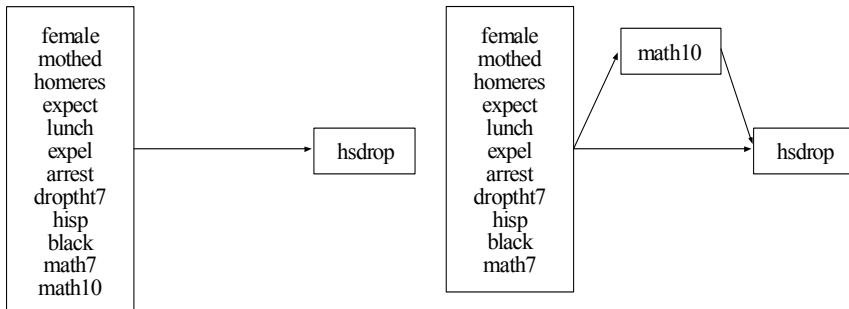
72

**Path Analysis With Categorical Outcomes**

**Path Analysis With A Binary Outcome And A Continuous Mediator With Missing Data**

**Logistic Regression**

**Path Model**



## Input For A Path Analysis With A Binary Outcome And A Continuous Mediator With Missing Data Using Monte Carlo Integration

```

TITLE:      Path analysis with a binary outcome and a continuous
            mediator with missing data using Monte Carlo integration
DATA:      FILE = lsaydropout.dat;
VARIABLE:  NAMES ARE female mothed homeres math7 math10 expel arrest
            hisp black hsdrop expect lunch droptht7;
            MISSING = ALL(9999);
            CATEGORICAL = hsdrop;
ANALYSIS:  ESTIMATOR = ML;
            INTEGRATION = MONTECARLO(500);
MODEL:     hsdrop ON female mothed homeres expect math7 math10 lunch
            expel arrest droptht7 hisp black;
            math10 ON female mothed homeres expect math7
            lunch expel arrest droptht7 hisp black;
OUTPUT:    PATTERNS STANDARDIZED TECH1 TECH8;
    
```

75

## Output Excerpts Path Analysis With A Binary Outcome And A Continuous Mediator With Missing Data Using Monte Carlo Integration

```

MISSING DATA PATTERNS FOR Y
           1      2
MATH10    x
FEMALE    x      x
MOTHED    x      x
HOMERES   x      x
MATH7     x      x
EXPEL     x      x
ARREST    x      x
HISP      x      x
BLACK     x      x
EXPECT    x      x
LUNCH     x      x
DROPTHT7  x      x

MISSING DATA PATTERN FREQUENCIES FOR Y
      Pattern      Frequency      Pattern      Frequency
           1           1639           2           574
    
```

76

**Output Excerpts Path Analysis With A Binary Outcome And A Continuous Mediator With Missing Data Using Monte Carlo Integration (Continued)**

**Tests Of Model Fit**

Loglikelihood

H0 Value -6323.175

Information Criteria

Number of Free Parameters 26  
 Akaike (AIC) 12698.350  
 Bayesian (BIC) 12846.604  
 Sample-Size Adjusted BIC 12763.999  
 (n\* = (n + 2) / 24)

77

**Output Excerpts Path Analysis With A Binary Outcome And A Continuous Mediator With Missing Data Using Monte Carlo Integration (Continued)**

**Model Results**

	Estimates	S.E.	Est./S.E.	Std	StdYX
HSDROP ON					
FEMALE	0.336	0.167	2.012	0.336	0.080
MOTHEd	-0.244	0.101	-2.421	-0.244	-0.117
HOMERES	-0.091	0.054	-1.699	-0.091	-0.072
EXPECT	-0.225	0.063	-3.593	-0.225	-0.147
MATH7	-0.012	0.015	-0.831	-0.012	-0.058
MATH10	-0.031	0.011	-2.816	-0.031	-0.201
LUNCH	0.005	0.004	1.456	0.005	-0.053
EXPEL	1.010	0.216	4.669	1.010	0.129
ARREST	0.033	0.314	0.105	0.033	0.003
DROPTHT7	0.679	0.272	2.499	0.679	0.067
HISP	-0.145	0.265	-0.548	-0.145	-0.019
BLACK	0.038	0.234	0.163	0.038	0.006

78

**Output Excerpts Path Analysis With A Binary Outcome And A Continuous Mediator With Missing Data Using Monte Carlo Integration (Continued)**

	Estimates	S.E.	Est./S.E.	Std	StdYX
MATH10 ON					
FEMALE	-0.973	0.410	-2.372	-0.973	-0.036
MOTHEd	0.343	0.219	1.570	0.343	0.026
HOMERES	0.486	0.140	3.485	0.486	0.059
EXPECT	1.014	0.166	6.111	1.014	0.103
MATH7	0.928	0.023	39.509	0.928	0.687
LUNCH	-0.039	0.011	-3.450	-0.039	-0.059
EXPEL	-1.404	0.851	-1.650	-1.404	-0.028
ARREST	-3.337	1.093	-3.052	-3.337	-0.052
DROPTHT7	-1.077	1.070	-1.007	-1.077	-0.016
HISP	-0.644	0.744	-0.866	-0.644	-0.013
BLACK	-0.809	0.694	-1.165	-0.809	-0.019

79

**Output Excerpts Path Analysis With A Binary Outcome And A Continuous Mediator With Missing Data Using Monte Carlo Integration (Continued)**

	Estimates	S.E.	Est./S.E.	Std	StdYX
Intercepts					
MATH10	10.941	1.269	8.621	10.941	0.809
Thresholds					
HSDROP\$1	-1.207	0.521	-2.319		
Residual Variances					
MATH10	65.128	2.280	28.571	65.128	0.356
Observed					
Variable R-Square					
HSDROP	0.255				
MATH10	0.644				

80



## Path Analysis Of Occupational Destination

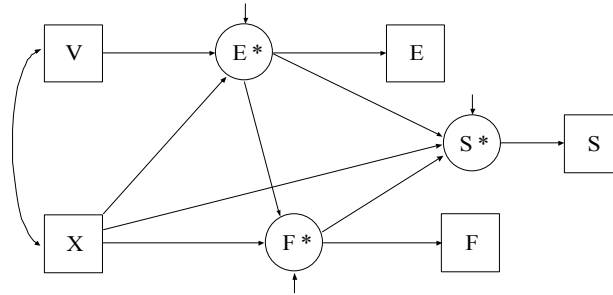


Figure 3: Structural Modeling of the Occupational Destination of Scientist or Engineer, Model 1

Reference: Xie (1989)

Data source: 1962 OCG Survey. The sample size is 14,401.

V: Father's Education. X: Father's Occupation (SEI)

81

## Path Analysis Of Occupational Destination (Continued)

Table 2. Descriptive Statistics of Discrete Dependent Variables

Variable	Code	Meaning	Percent
S: Current Occupation	0	Non-scientific/engineering	96.4
	1	Scientific/engineering	3.6
F: First Job	0	Non-scientific/engineering	98.3
	1	Scientific/engineering	1.7
E: Education	0	0-7 years	13.4
	1	8-11 years	32.6
	2	12 years	29.0
	3	13 and more years	25.0

82

### **Differences Between Weighted Least Squares And Maximum Likelihood Model Estimation For Categorical Outcomes In Mplus**

- Probit versus logistic regression
  - Weighted least squares estimates probit regressions
  - Maximum likelihood estimates logistic or probit regressions
- Modeling with underlying continuous variables versus observed categorical variables for categorical outcomes that are mediating variables
  - Weighted least squares uses underlying continuous variables
  - Maximum likelihood uses observed categorical outcomes

83

### **Differences Between Weighted Least Squares And Maximum Likelihood Model Estimation For Categorical Outcomes In Mplus (Continued)**

- Delta versus Theta parameterization for weighted least squares
  - Equivalent in most cases
  - Theta parameterization needed for models where categorical outcomes are predicted by categorical dependent variables while predicting other dependent variables
- Missing data
  - Weighted least squares allows missingness predicted by covariates
  - Maximum likelihood allows MAR
- Testing of nested models
  - WLSMV uses DIFFTEST
  - Maximum likelihood (ML, MLR) uses regular or special approaches

84

## **Further Readings On Path Analysis With Categorical Outcomes**

- MacKinnon, D.P., Lockwood, C.M., Brown, C.H., Wang, W., & Hoffman, J.M. (2007). The intermediate endpoint effect in logistic and probit regression. Clinical Trials, 4, 499-513.
- Xie, Y. (1989). Structural equation models for ordinal variables. Sociological Methods & Research, 17, 325-352.

85

## **Categorical Observed And Continuous Latent Variables**

86

## Continuous Latent Variable Analysis With Categorical Outcomes

### Model Identification

- EFA, CFA, and SEM the same as for continuous outcomes
- Multiple group and models for longitudinal data require invariance of measurement thresholds and loadings, requiring threshold structure (and scale factor parameters)

### Interpretation

- Estimated coefficients – sign, significance most important
- Estimated coefficients can be converted to probabilities

87

## Continuous Latent Variable Analysis With Categorical Outcomes (Continued)

### Estimation

- Maximum likelihood computational burden increases significantly with number of factors
- Weighted least squares computation burden increases significantly with the number of variables

### Model Fit

- Only chi-square studied
- Simulation studies needed for TLI, CFI, RMSEA, SRMR, and WRMR (see, however, Yu, 2002)

88

**Item Response Theory**

89

**Item Response Theory**

Latent trait modeling  
Factor analysis with categorical outcomes

A graph showing the probability  $P(u_j = 1 | \eta)$  on the y-axis (ranging from 0 to 1) against the latent trait  $\eta$  on the x-axis. Five sigmoidal curves are plotted, each representing a different item. The curves are shifted horizontally, indicating different item difficulties. Below the x-axis, a bell-shaped curve represents the distribution of the latent trait  $\eta$ .

A path diagram illustrating a latent trait model. A latent variable  $\eta$  (represented by a circle) is shown at the bottom, with five arrows pointing upwards to five observed variables  $u_1, u_2, u_3, u_4, u_5$  (represented by rectangles).

90

## Item Response Theory (Continued)

IRT typically does not use the full SEM model

$$u_i^* = \nu + A \eta_i + K x_i + \varepsilon_i, \quad (127)$$

$$\eta_i = \alpha + B \eta_i + \Gamma x_i + \zeta_i, \quad (128)$$

and typically considers a single  $\eta$  (see, however, Bock, Gibbons, & Muraki, 1988). Aims:

- Item parameter estimation (ML): Calibration
- Estimation of  $\eta$  values: Scoring
- Assessment of information function
- Test equating
- DIF analysis

91

## IRT Models And Estimators In Mplus

- ML (full information estimation): Logit and probit links
- WLS (limited information estimation): Probit link

92

### Translating Factor Analysis Parameters In Mplus To IRT Parameters

- IRT calls the continuous latent variable  $\theta$
- 2-parameter logistic IRT model uses

$$P(u = I | \theta) = \frac{I}{I + e^{-Da(\theta-b)}}$$

with  $D = 1.7$  to make  $a, b$  close to those of probit

$a$  discrimination

$b$  difficulty

- 2-parameter normal ogive IRT model uses

$$P(u = I | \theta) = \Phi [a(\theta - b)]$$

- Typically  $\theta \sim N(0,1)$

93

### Translating Factor Analysis Parameters To IRT Parameters (Continued)

- The Mplus factor analysis model uses

$$P(u = I | \eta) = \frac{I}{I + e^{-(\tau + \lambda\eta)}} \quad \text{for logit}$$

$$P(u = I | \eta) = \Phi [(-\tau + \lambda\eta)\theta^{-1/2}] \quad \text{for probit}$$

where  $\theta$  is the residual variance

The logit conversion is:

The probit conversion is:

$$a = \lambda\sqrt{\psi} / D$$

$$a = \lambda\sqrt{\psi} \theta^{-1/2}$$

$$b = (\tau - \lambda\alpha) / \lambda\sqrt{\psi}$$

$$b = (\tau - \lambda\alpha) / \lambda\sqrt{\psi}$$

- Conversion automatically done in Mplus

94

## Testing The Model Against Data

- Model fit to frequency tables. Overall test against data
  - When the model contains only  $\mathbf{u}$ , summing over the cells,

$$\chi_P^2 = \sum_i \frac{(o_i - e_i)^2}{e_i}, \quad (82)$$

$$\chi_{LR}^2 = 2 \sum_i o_i \log o_i / e_i. \quad (83)$$

A cell that has non-zero observed frequency and expected frequency less than .01 is not included in the  $\chi^2$  computation as the default. With missing data on  $\mathbf{u}$ , the EM algorithm described in Little and Rubin (1987; chapter 9.3, pp. 181-185) is used to compute the estimated frequencies in the unrestricted multinomial model. In this case, a test of MCAR for the unrestricted model is also provided (Little & Rubin, 1987, pp. 192-193).

- Model fit to univariate and bivariate frequency tables. Mplus TECH10

95

## Antisocial Behavior (ASB) Data

The Antisocial Behavior (ASB) data were taken from the National Longitudinal Survey of Youth (NLSY) that is sponsored by the Bureau of Labor Statistics. These data are made available to the public by Ohio State University. The data were obtained as a multistage probability sample with oversampling of blacks, Hispanics, and economically disadvantaged non-blacks and non-Hispanics.

Data for the analysis include 15 of the 17 antisocial behavior items that were collected in 1980 when respondents were between the ages of 16 and 23 and the background variables of age, gender and ethnicity. The ASB items assessed the frequency of various behaviors during the past year. A sample of 7,326 respondents has complete data on the antisocial behavior items and the background variables of age, gender, and ethnicity. Following is a list of the 15 items:

96



## Antisocial Behavior (ASB) Data (Continued)

Damaged property	Use other drugs
Fighting	Sold marijuana
Shoplifting	Sold hard drugs
Stole < \$50	“Con” someone
Stole > \$50	Take auto
Seriously threaten	Broken into building
Intent to injure	Held stolen goods
Use marijuana	

These items were dichotomized 0/1 with 0 representing never in the last year. An EFA suggested three factors: property offense, person offense, and drug offense.

97

## Input For IRT Analysis Of Eight ASB Property Offense Items

```

TITLE:      2-parameter logistic IRT
            for 8 property offense items
DATA:      FILE = asb.dat;
            FORMAT = 34X 54F2.0;
VARIABLE:  NAMES = property fight shoplift lt50 gt50 force threat
            injure pot drug soldpot solddrug con auto bldg goods
            gambling
            dsml-dsm22 sex black hisp single divorce dropout
            college onset f1 f2 f3
            age94 cohort dep abuse;
            USEVAR = property shoplift lt50 gt50 con auto bldg
            goods;
            CATEGORICAL = property-goods;
ANALYSIS:  ESTIMATOR = MLR;
MODEL:     f BY property-goods*;
            f@1;
OUTPUT:    TECH1 TECH8 TECH10;
PLOT:     TYPE = PLOT3;

```

98

## Output Excerpts IRT Analysis Of Eight ASB Property Offense Items

TESTS OF MODEL FIT

Loglikelihood

H0 Value	-19758.361
H0 Scaling Correction Factor for MLR	0.996

Information Criteria

Number of Free Parameters	16
Akaike (AIC)	39548.722
Bayesian (BIC)	39659.109
Sample-Size Adjusted BIC	39608.265

(n\* = (n + 2) / 24)

Chi-Square Test of Model Fit for the Binary and Ordered  
Categorical (Ordinal) Outcomes

Pearson Chi-Square

Value	324.381
Degrees of Freedom	239
P-Value	0.0002

99

## Output Excerpts IRT Analysis Of Eight ASB Property Offense Items (Continued)

Likelihood Ratio Chi-Square

Value	327.053
Degrees of Freedom	239
P-Value	0.0001

MODEL RESULTS

		Estimate	S.E.	Est./S.E.	Two-Tailed P-Value
F	BY				
	PROPERTY	2.032	0.084	24.060	0.000
	SHOPLIFT	1.712	0.068	25.115	0.000
	LT50	1.850	0.076	24.411	0.000
	GT50	2.472	0.139	17.773	0.000
	CON	1.180	0.051	23.148	0.000
	AUTO	1.383	0.070	19.702	0.000
	BLDG	2.741	0.151	18.119	0.000
	GOODS	2.472	0.116	21.339	0.000

100

**Output Excerpts IRT Analysis Of Eight ASB  
Property Offense Items (Continued)**

	Estimate	S.E.	Est./S.E.	Two-Tailed P-Value
<b>Thresholds</b>				
PROPERTY\$1	2.398	0.073	32.803	0.000
SHOPLIFT\$1	1.529	0.049	31.125	0.000
LT50\$1	2.252	0.065	34.509	0.000
GT50\$1	5.054	0.195	25.912	0.000
CON\$1	1.560	0.041	37.894	0.000
AUTO\$1	3.144	0.079	39.948	0.000
BLDG\$1	5.185	0.208	24.983	0.000
GOODS\$1	3.691	0.126	29.316	0.000
<b>Variances</b>				
F	1.000	0.000	999.000	999.000

101

**Output Excerpts IRT Analysis Of Eight ASB  
Property Offense Items (Continued)**

IRT PARAMETERIZATION IN TWO-PARAMETER LOGISTIC METRIC WHERE THE LOGIT IS  $1.7 * \text{DISCRIMINATION} * (\text{THETA} - \text{DIFFICULTY})$

Item Discriminations	Estimate	S.E.	Est./S.E.	Two-Tailed P-Value
F BY				
PROPERTY	1.195	0.050	24.060	0.000
SHOPLIFT	1.007	0.040	25.115	0.000
LT50	1.088	0.045	24.411	0.000
GT50	1.454	0.082	17.773	0.000
CON	0.694	0.030	23.148	0.000
AUTO	0.813	0.041	19.702	0.000
BLDG	1.612	0.089	18.119	0.000
GOODS	1.454	0.068	21.339	0.000

102

**Output Excerpts IRT Analysis Of Eight ASB  
Property Offense Items (Continued)**

Item Difficulties	Two-Tailed			
	Estimate	S.E.	Est./S.E.	P-Value
PROPERTY\$1	1.180	0.031	38.268	0.000
SHOPLIFT\$1	0.893	0.029	31.309	0.000
LT50\$1	1.217	0.033	36.604	0.000
GT50\$1	2.044	0.053	38.588	0.000
CON\$1	1.322	0.048	27.809	0.000
AUTO\$1	2.274	0.081	28.232	0.000
BLDG\$1	1.891	0.045	42.204	0.000
GOODS\$1	1.493	0.035	43.045	0.000
<b>Variances</b>				
F	1.000	0.000	0.000	1.000

103

**Output Excerpts IRT Analysis Of Eight ASB  
Property Offense Items (Continued)**

TECHNICAL 10 OUTPUT  
MODEL FIT INFORMATION FOR THE LATENT CLASS INDICATOR MODEL PART  
RESPONSE PATTERNS

No.	Pattern	No.	Pattern	No.	Pattern	No.	Pattern
1	00000000	2	10100000	3	00001101	4	00000010
5	01100000	6	00001000	7	10001010	8	00010001
9	10100010	10	11000000	11	10101110	12	11100010
13	11010111	14	10000000	15	11110001	16	10000001

104

## Output Excerpts IRT Analysis Of Eight ASB Property Offense Items (Continued)

RESPONSE PATTERN FREQUENCIES AND CHI-SQURE CONTRIBUTIONS

Response Pattern	Frequency		Standardized Residual (z-score)	Chi-square Pearson	Contribution Loglikelihood
	Observed	Estimated			
1	3581.00	3565.17	0.37	0.07	31.73
2	60.00	57.05	0.39	0.15	6.05
3	2.00	3.12	-0.77	0.59	-2.14
4	18.00	17.65	0.08	0.01	0.71
5	137.00	110.30	2.56	6.46	59.39
6	476.00	495.86	-0.92	0.80	-38.92

105

## Output Excerpts IRT Analysis Of Eight ASB Property Offense Items (Continued)

BIVARIATE MODEL FIT INFORMATION

VARIABLE PROPERTY	VARIABLE SHOPLIFT	Estimated Probabilities		
		H1	H0	Standardized Residual (z-score)
Category 1	Category 1	0.656	0.655	0.157
Category 1	Category 2	0.159	0.160	-0.176
Category 2	Category 1	0.080	0.081	-0.285
Category 2	Category 2	0.105	0.104	0.222
Bivariate Pearson Chi-Square				0.153
Bivariate Log-Likelihood Chi-Square				0.077

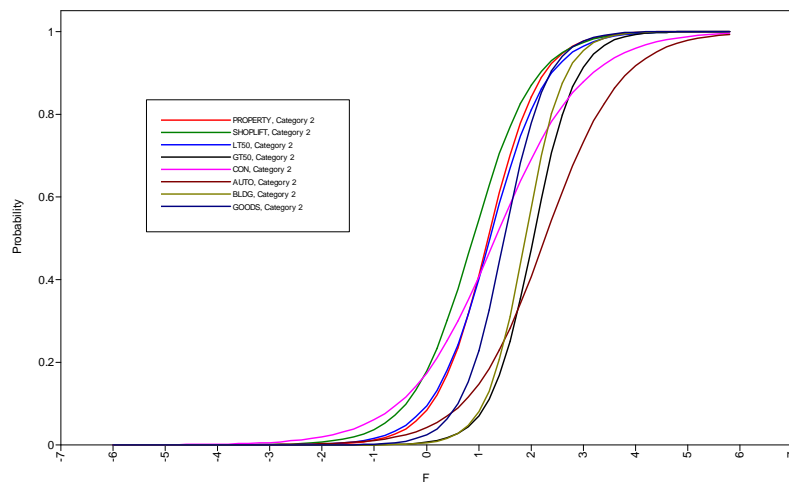
106

## Output Excerpts IRT Analysis Of Eight ASB Property Offense Items (Continued)

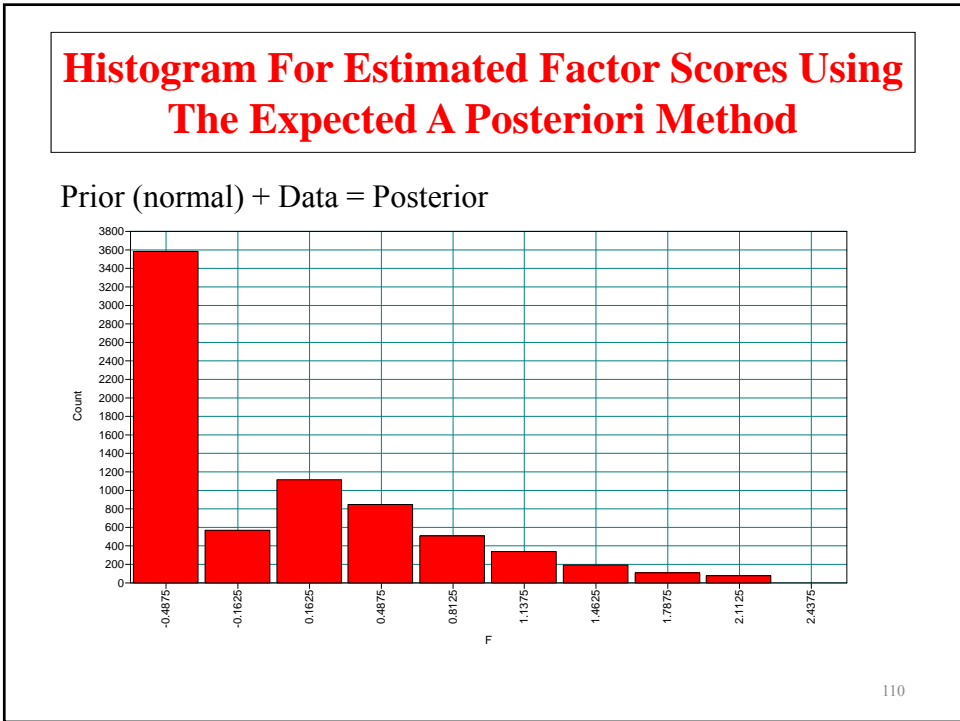
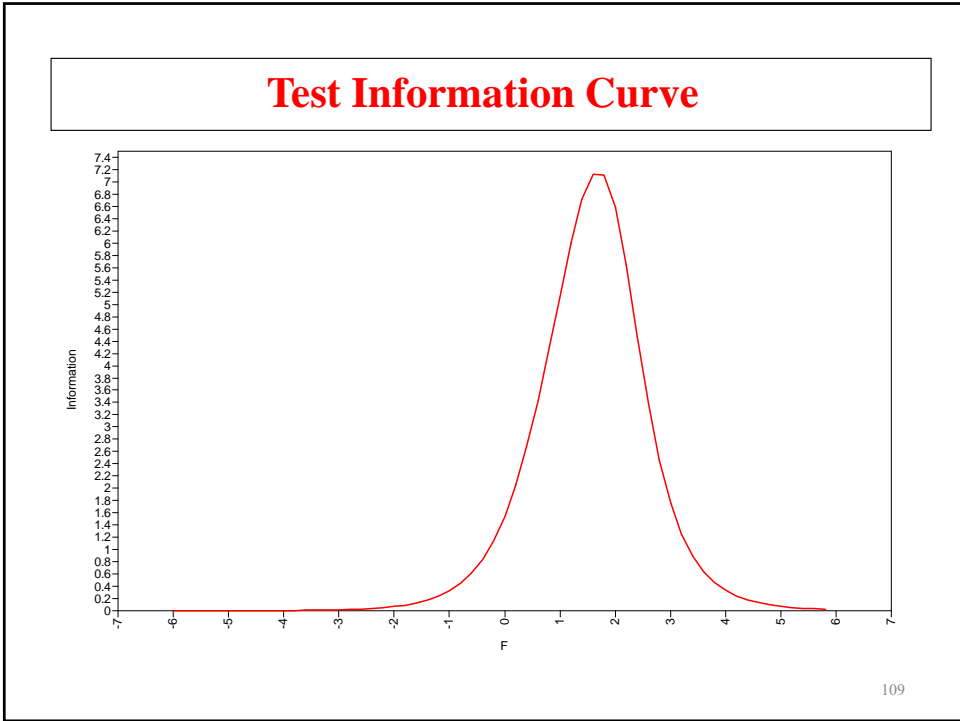
VARIABLE	VARIABLE	Estimated Probabilities		
		H1	H0	Standardized Residual (z-score)
PROPERTY	SHOPLIFT			
LT50	GT50			
Category 1	Category 1	0.799	0.795	0.873
Category 1	Category 2	0.014	0.018	-2.615
Category 2	Category 1	0.152	0.156	-0.945
Category 2	Category 2	0.035	0.032	1.912
Bivariate Pearson Chi-Square				11.167
Bivariate Log-Likelihood Chi-Square				5.806

107

## Item Characteristic Curves



108



## Further Readings On IRT

- Baker, F.B. & Kim, S.H. (2004). Item response theory. Parameter estimation techniques. Second edition. New York: Marcel Dekker.
- Bock, R.D. (1997). A brief history of item response theory. Educational Measurement: Issues and Practice, 16, 21-33.
- du Toit, M. (2003). IRT from SSI. Lincolnwood, IL: Scientific Software International, Inc. (BILOG, MULTILOG, PARSCALE, TESTFACT)
- Embretson, S. E., & Reise, S. P. (2000). Item response theory for psychologists. Mahwah, NJ: Erlbaum.
- Hambleton, R.K. & Swaminathan, H. (1985). Item response theory. Boston: Kluwer-Nijhoff.
- MacIntosh, R. & Hashim, S. (2003). Variance estimation for converting MIMIC model parameters to IRT parameters in DIF analysis. Applied Psychological Measurement, 27, 372-379.
- Muthén, B., Kao, Chih-Fen, & Burstein, L. (1991). Instructional sensitivity in mathematics achievement test items: Applications of a new IRT-based detection technique. Journal of Educational Measurement, 28, 1-22. (#35)

111

## Further Readings On IRT (Continued)

- Muthén, B. & Asparouhov, T. (2002). Latent variable analysis with categorical outcomes: Multiple-group and growth modeling in Mplus. Mplus Web Note #4 ([www.statmodel.com](http://www.statmodel.com)).
- Takane, Y. & DeLeeuw, J. (1987). On the relationship between item response theory and factor analysis of discretized variables. Psychometrika, 52, 393-408.

112



## Exploratory Factor Analysis

113

### Exploratory Factor Analysis For Outcomes That Are Categorical, Censored, Counts

Rotation of the factor loading matrix as with continuous outcomes

- Maximum-likelihood estimation
  - Computationally feasible for only a few factors, but can handle many items
  - Frequency table testing typically not useful
- Limited-information weighted least square estimation
  - Computationally feasible for many factors, but not huge number of items
  - Testing against bivariate tables
  - Modification indices for residual correlations

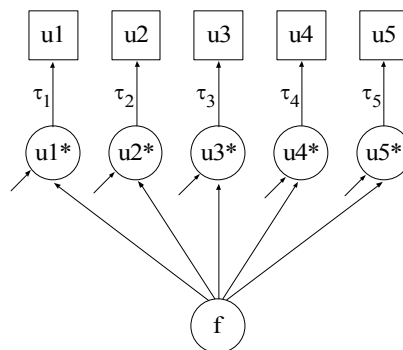
114

## Assumptions Behind ML And WLS

Note that when assuming normal factors and using probit links, ML uses the same model as WLS. This is because normal factors and probit links result in multivariate normal  $u^*$  variables. For model estimation, WLS uses the limited information of first- and second-order moments, thresholds and sample correlations of the multivariate normal  $u^*$  variables (tetrachoric, polychoric, and polyserial correlations), whereas ML uses full information from all moments of the data.

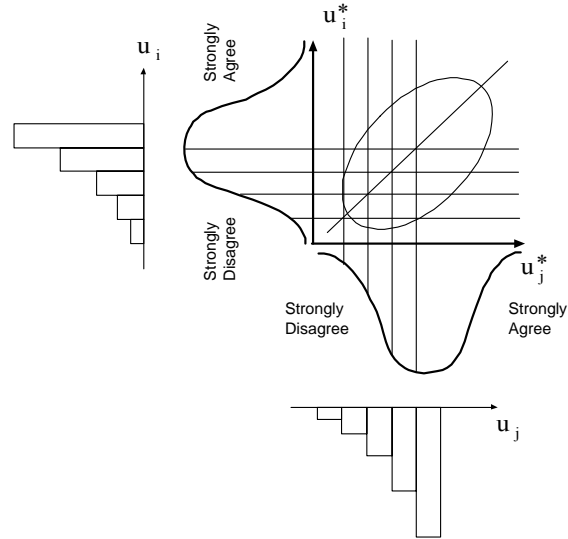
115

## Latent Response Variable Formulation Of A Factor Model



116

## Latent Response Variable Correlations



117

## Sample Statistics With Categorical Outcomes And Weighted Least Squares Estimation

- Types of  $u^*$  correlations (normality assumed)
  - Both dichotomous – tetrachoric
  - Both polytomous – polychoric
  - One dichotomous, one continuous – biserial
  - One polytomous, one continuous – polyserial
- Analysis choices
  - Case A – no  $x$  variables – use  $u^*$  correlations
  - Case B –  $x$  variables present
    - Use  $u^*$  correlations (full normality of  $u^*$  and  $x$  assumed)
    - Use regression-based statistics (conditional normality of  $u^*$  given  $x$  assumed)

118

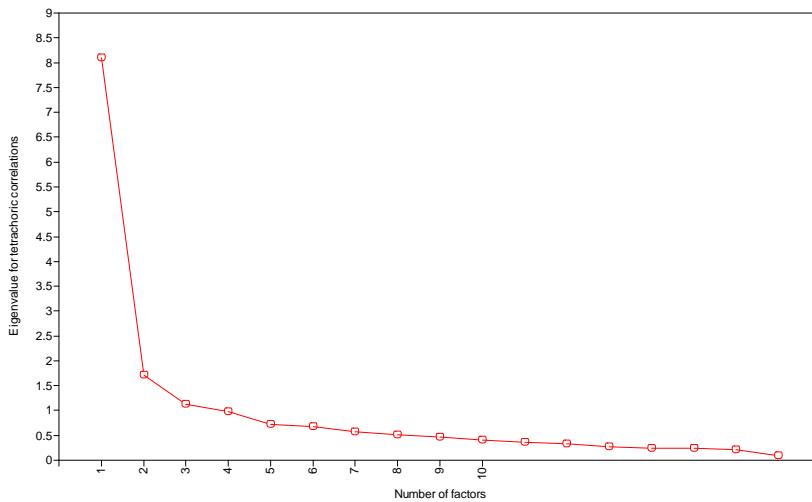
## Exploratory Factor Analysis Of 17 ASB Items Using WLSM

```

TITLE:      EFA using WLSM
DATA:      FILE = asb.dat;
           FORMAT = 34X 54F2.0;
VARIABLE:  NAMES = property fight shoplift lt50 gt50 force threat
           injure pot drug
           soldpot solddrug con auto bldg goods gambling
           dsml-dsm22 sex black hisp single divorce dropout
           college onset f1 f2 f3
           age94 cohort dep abuse;
           USEVAR = property-gambling;
           CATEGORICAL = property-gambling;
ANALYSIS:  TYPE = EFA 1 5;
OUTPUT:    MODINDICES;
PLOT:      TYPE = PLOT3;
    
```

119

## Eigenvalue Plot For Tetrachoric Correlations Among 17 ASB Items



120

**Output Excerpts 3- And 4-Factor WLSM EFA  
Of 17 ASB Items**

EXPLORATORY FACTOR ANALYSIS WITH 3 FACTOR(S):

TESTS OF MODEL FIT

Chi-Square Test of Model Fit

Value	584.356*
Degrees of Freedom	88
P-Value	0.0000

\* The chi-square value for MLM, MLMV, MLR, ULSMV, WLSM and WLSMV cannot be used for chi-square difference tests. MLM, MLR and WLSM chi-square difference testing is described in the Mplus Technical Appendices at [www.statmodel.com](http://www.statmodel.com). See chi-square difference testing in the index of the Mplus User's Guide.

Chi-Square Test of Model Fit for the Baseline Model

Value	53652.583
Degrees of Freedom	136
P-Value	0.0000

121

**Output Excerpts 3- And 4-Factor WLSM EFA  
Of 17 ASB Items (Continued)**

CFI/TLI

CFI	0.991
TLI	0.986

Number of Free Parameters 48

RMSEA (Root Mean Square Error Of Approximation)

Estimate	0.028
----------	-------

SRMR (Standardized Root Mean Square Residual)

Value	0.045
-------	-------

MINIMUM ROTATION FUNCTION VALUE 0.08510

122

**Output Excerpts 3- And 4-Factor WLSM EFA  
Of 17 ASB Items (Continued)**

	QUARTIMIN ROTATED LOADINGS		
	1	2	3
PROPERTY	<b>0.669</b>	0.179	-0.036
FIGHT	0.266	<b>0.548</b>	-0.121
SHOPLIFT	<b>0.600</b>	-0.028	0.185
LT50	<b>0.818</b>	-0.185	0.046
GT50	<b>0.807</b>	0.003	0.016
FORCE	0.379	0.344	0.000
THREAT	-0.008	<b>0.821</b>	0.049
INJURE	-0.022	<b>0.761</b>	0.101
POT	-0.051	0.001	<b>0.903</b>
DRUG	-0.021	-0.020	<b>0.897</b>
SOLDPOT	0.126	0.058	<b>0.759</b>
SOLDDRUG	0.175	0.083	<b>0.606</b>
CON	0.460	0.228	-0.065

123

**Output Excerpts 3- And 4-Factor WLSM EFA  
Of 17 ASB Items (Continued)**

	1	2	3
AUTO	<b>0.460</b>	0.139	0.073
BLDG	<b>0.797</b>	0.033	0.017
GOODS	0.700	0.109	0.066
GAMBLING	0.314	0.327	0.092

QUARTIMIN FACTOR CORRELATIONS			
	1	2	3
1	1.000		
2	0.598	1.000	
3	0.614	0.371	1.000

124

**Output Excerpts 3- And 4-Factor WLSM EFA  
Of 17 ASB Items (Continued)**

EXPLORATORY FACTOR ANALYSIS WITH 4 FACTOR(S):

TESTS OF MODEL FIT

Chi-Square Test of Model Fit

Value	303.340*
Degrees of Freedom	74
P-Value	0.0000

\* The chi-square value for MLM, MLMV, MLR, ULSMV, WLSM and WLSMV cannot be used for chi-square difference tests. MLM, MLR and WLSM chi-square difference testing is described in the Mplus Technical Appendices at [www.statmodel.com](http://www.statmodel.com). See chi-square difference testing in the index of the Mplus User's Guide.

125

**Output Excerpts 3- And 4-Factor WLSM EFA  
Of 17 ASB Items (Continued)**

Chi-Square Test of Model Fit for the Baseline Model

Value	53652.583
Degrees of Freedom	136
P-Value	0.0000

CFI/TLI

CFI	0.996
TLI	0.992

Number of Free Parameters 62

RMSEA (Root Mean Square Error Of Approximation)

Estimate	0.021
----------	-------

SRMR (Standardized Root Mean Square Residual)

Value	0.026
-------	-------

MINIMUM ROTATION FUNCTION VALUE 0.19546

126

**Output Excerpts 3- And 4-Factor WLSM EFA  
Of 17 ASB Items (Continued)**

	QUARTIMIN ROTATED LOADINGS			
	1	2	3	4
PROPERTY	<b>0.670</b>	0.191	-0.006	-0.043
FIGHT	0.290	<b>0.537</b>	-0.060	-0.098
SHOPLIFT	<b>0.679</b>	-0.001	0.225	-0.159
LT50	<b>0.817</b>	-0.152	0.066	-0.049
GT50	<b>0.762</b>	-0.008	-0.036	0.154
FORCE	0.257	0.288	-0.195	0.491
THREAT	0.003	<b>0.858</b>	0.101	-0.078
INJURE	-0.036	<b>0.728</b>	0.056	0.162
POT	0.041	0.074	<b>0.923</b>	-0.069
DRUG	0.051	0.007	<b>0.717</b>	0.227
SOLDPOT	0.149	0.070	<b>0.598</b>	0.281
SOLDDRUG	0.065	-0.037	0.269	<b>0.791</b>
CON	0.420	0.223	-0.072	0.081

127

**Output Excerpts 3- And 4-Factor WLSM EFA  
Of 17 ASB Items (Continued)**

	1	2	3	4
AUTO	<b>0.446</b>	0.138	0.051	0.074
BLDG	<b>0.770</b>	0.042	0.010	0.055
GOODS	<b>0.662</b>	0.109	0.030	0.126
GAMBLING	0.208	0.270	-0.083	0.449

QUARTIMIN FACTOR CORRELATIONS				
	1	2	3	4
1	1.000			
2	0.571	1.000		
3	0.485	0.230	1.000	
4	0.481	0.312	0.376	1.000

128



## **Practical Issues In The Analysis Of Categorical Outcomes**

129

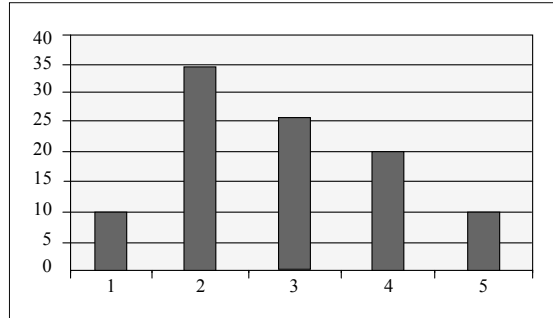
## **Overview Of Practical Issues In The Analysis Of Categorical Outcomes**

- When Is A Variable Best Treated As Categorical?
  - Less dependent on number of categories than the presence of floor and ceiling effects
  - When the aim is to estimate probabilities or odds
- What's Wrong With Treating Categorical Variables As Continuous Variables?
  - Correlations will be attenuated particularly when there are floor and ceiling effects
  - Can lead to factors that reflect item difficulty extremeness
  - Predicted probabilities can be outside the 0/1 range

130

## Approaches To Use With Categorical Data

- Data that lead to incorrect standard errors and chi-square under normality assumption

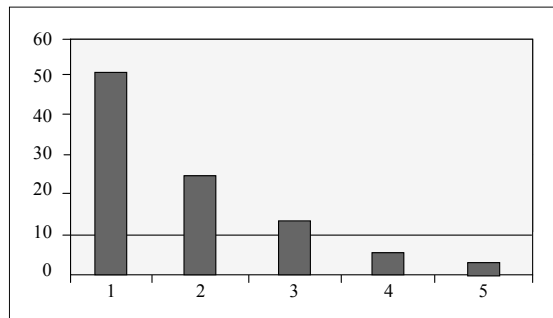


- Transform variable and treat as a continuous variable
- Treat as a continuous variable and use non-normality robust maximum likelihood estimation

131

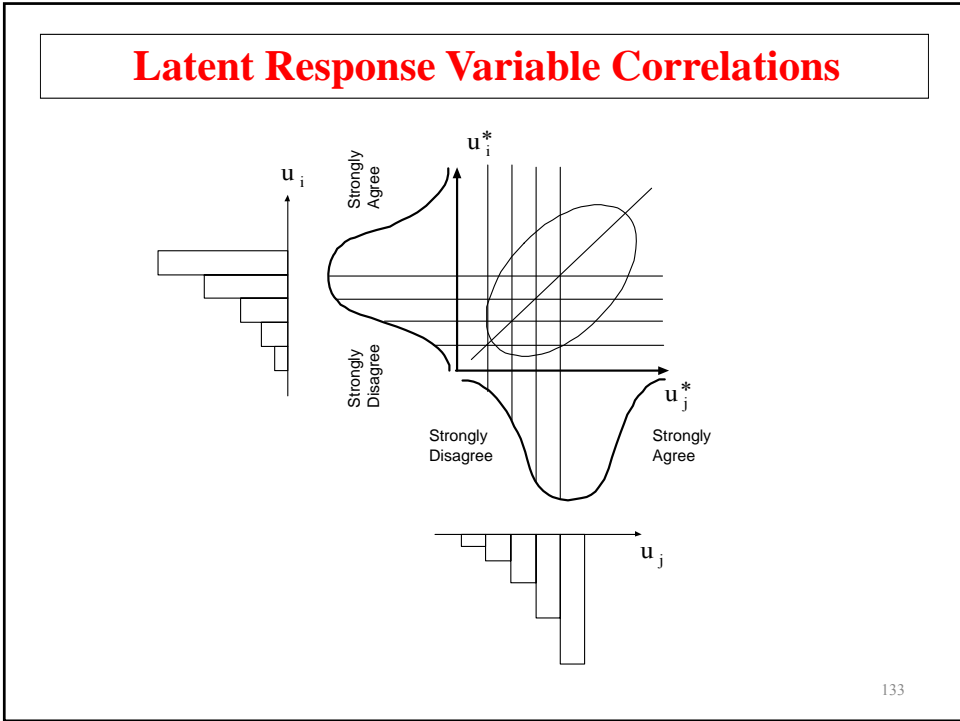
## Approaches To Use With Categorical Data (Continued)

- Data that lead to incorrect standard errors, chi-square, and parameter estimates under normality assumption



- Treat as a categorical variable

132



### Distortions Of Underlying Correlation Structure

Pearson product-moment correlations unsuited to categorical variables due to limitation in range.

Example:  $P(u_1) = 0.5, P(u_2 = 1) = 0.2$   
 Gives max Pearson correlation = 0.5

		Variable 1	
		0	1
Variable 2	0	50	30
	1	0	20
		50	100

134

## Distortions Of Underlying Correlation Structure (Continued)

Phi coefficient (Pearson correlation):

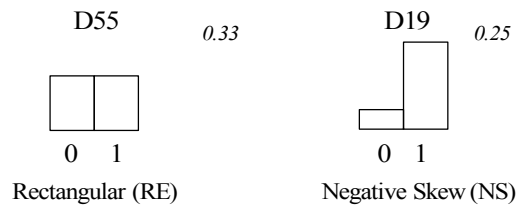
$$R = \frac{\text{Cov}(u_1, u_2)}{\text{SD}(u_1)\text{SD}(u_2)} = \frac{P(u_1 = 1 \text{ and } u_2 = 1) - P(u_1 = 1)P(u_2 = 1)}{\sqrt{P(u_1 = 1)[1 - P(u_1 = 1)]} \sqrt{P(u_2 = 1)[1 - P(u_2 = 1)]}}$$

$$R_{\text{max.}} = \frac{0.2 - 0.5 \times 0.2}{\sqrt{.5 \times .5} \sqrt{.2 \times .8}} = \frac{0.1}{0.2} = 0.5$$

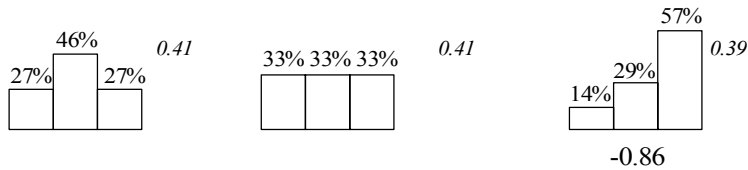
135

## Correlational Attenuation

Correlation between underlying continuous  $u^*$  variables = 0.5



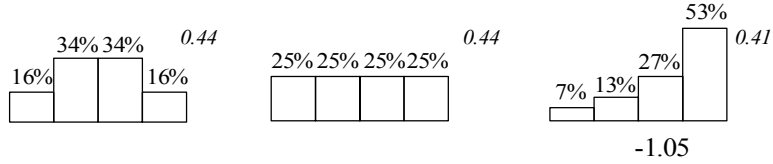
Three Categories



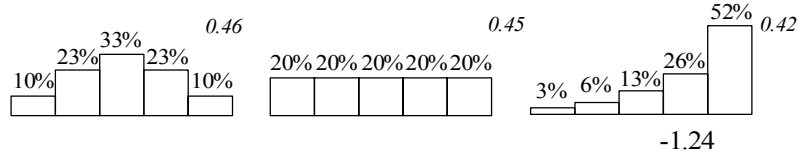
136

## Correlational Attenuation (Continued)

### Four Categories



### Five Categories



137

Table 1 (Part 2)  
Pearson Correlations for True Correlations = 0.50

	D19	D28	D37	D46	D55	D64	D73	D82	D91	3SY	3RE	3NS	3PS
D19	25												
D28	26	30											
D37	26	30	32										
D46	24	30	32	33									
D55	23	28	31	33	33								
D64	20	26	30	23	33	33							
D73	18	23	27	30	31	32	32						
D82	15	20	23	26	28	30	30	30					
D91	10	15	18	20	22	24	26	26	25				
3SY	26	32	35	36	37	36	35	32	26	41			
3RE	25	31	35	36	37	36	35	31	25	41	41		
3NS	29	33	35	36	35	33	30	26	20	39	39	39	
3PS	20	26	30	33	35	36	35	33	29	39	39	34	39
4SY	27	33	36	38	38	38	36	33	27	43	43	40	40
4RE	26	33	36	38	38	38	36	33	26	42	42	40	40
4NS	30	35	36	36	35	33	30	27	20	40	39	40	34
3PS	20	27	31	34	35	36	36	35	30	40	39	34	40
5SY	28	34	37	38	39	38	37	34	28	44	43	41	41
4RE	27	33	37	38	39	38	37	33	27	43	43	41	41
5NS	31	35	36	36	35	33	30	26	20	40	39	40	34
5PS	20	26	30	33	35	36	36	35	31	40	39	34	40
CON	29	35	38	39	40	39	38	35	29	45	45	42	42
	D19	D28	D37	D46	D55	D64	D73	D82	D91	3SY	3RE	3NS	3PS

138

Pearson Correlations for True Correlations = 0.50

	4SY	4RE	4NS	4PS	5SY	5RE	5NS	5PS	CON
4SY	44								
4RE	44	44							
4NS	41	41	41						
4PS	41	41	35	41					
5SY	45	45	42	42	46				
5RE	45	45	41	41	46	45			
5NS	41	40	42	34	42	41	42		
5PS	41	41	34	42	42	41	34	42	
CON	47	46	43	43	48	47	44	44	50
	4SY	4RE	4NS	4PS	5SY	5RE	5NS	5PS	CON

139

**Approaches To Use With Categorical Outcomes**

- Items, Testlets, Sums, Or Factor Scores?
  - A sum of at least 15 unidimensional items is reliable
  - Testlets can be used as continuous indicators
  - Factor scores can be estimated as in IRT
  
- Sample Size
  - Larger than for continuous variables
  - Univariate and bivariate distributions should contain several observations per cell

140

## Further Readings On Factor Analysis Of Categorical Outcomes

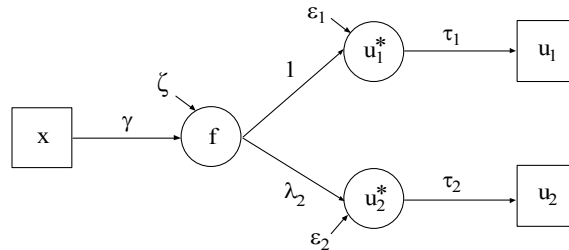
- Bock, R.D., Gibbons, R., & Muraki, E.J. (1998). Full information item factor analysis. Applied Psychological Measurement, 12, 261-280.
- Flora, D.B. & Curran, P.J., (2004). An empirical evaluation of alternative methods of estimation for confirmatory factor analysis with ordinal data. Psychological Methods, 9, 466-491.
- Muthén, B. (1989). Dichotomous factor analysis of symptom data. In Eaton & Bohrnstedt (Eds.), *Latent variable models for dichotomous outcomes: Analysis of data from the epidemiological Catchment Area program* (pp.19-65), a special issue of Sociological Methods & Research, 18, 19-65.
- Muthen, B. & Kaplan, D. (1985). A comparison of some methodologies for the factor analysis of non-normal Likert variables. British Journal of Mathematical and Statistical Psychology, 38, 171-189.
- Muthen, B. & Kaplan, D. (1992). A comparison of some methodologies for the factor analysis of non-normal Likert variables: A note on the size of the model. British Journal of Mathematical and Statistical Psychology, 45, 19-30.

141

## CFA With Covariates (MIMIC)

142

## CFA With Covariates Using WLS



$$u_{ij}^* = \lambda_j f_i + \varepsilon_{ij}, (j = 1, 2)$$

$$f_i = \gamma x_i + \zeta_i$$

Estimate CFA model by fitting to probit / logit regression estimates

143

## CFA With Covariates (MIMIC)

Used to study the effects of covariates or background variables on the factors and outcome variables to understand measurement invariance and heterogeneity

- Measurement non-invariance – direct relationships between the covariates and outcome variables that are not mediated by the factors – if they are significant, this indicates measurement non-invariance due to differential item functioning (DIF)
- Population heterogeneity – relationships between the covariates and the factors – if they are significant, this indicates that the factor means are different for different levels of the covariates.

### Model Assumptions

- Same factor loadings and observed residual variances / covariances for all levels of the covariates
- Same factor variances and covariances for all levels of the covariates

144



## **Steps In CFA With Covariates**

- Establish a CFA or EFA/CFA model
- Add covariates – check that factor structure does not change and study modification indices for possible direct effects
- Add direct effects suggested by modification indices – check that factor structure does not change
- Interpret the model
  - Factors
  - Effects of covariates on factors
  - Effects of covariates on factor indicators

145

## **Antisocial Behavior (ASB) Data**

The Antisocial Behavior (ASB) data were taken from the National Longitudinal Survey of Youth (NLSY) that is sponsored by the Bureau of Labor Statistics. These data are made available to the public by Ohio State University. The data were obtained as a multistage probability sample with oversampling of blacks, Hispanics, and economically disadvantaged non-blacks and non-Hispanics.

Data for the analysis include 15 of the 17 antisocial behavior items that were collected in 1980 when respondents were between the ages of 16 and 23 and the background variables of age, gender and ethnicity. The ASB items assessed the frequency of various behaviors during the past year. A sample of 7,326 respondents has complete data on the antisocial behavior items and the background variables of age, gender, and ethnicity. Following is a list of the 15 items:

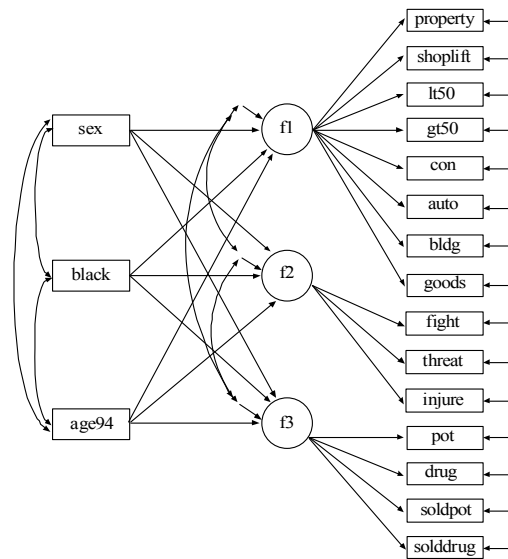
146

## Antisocial Behavior (ASB) Data (Continued)

Damaged property Fighting Shoplifting Stole < \$50 Stole > \$50 Seriously threaten Intent to injure Use marijuana	Use other drugs Sold marijuana Sold hard drugs "Con" someone Take auto Broken into building Held stolen goods
--	---

These items were dichotomized 0/1 with 0 representing never in the last year. An EFA suggested three factors: property offense, person offense, and drug offense.

147



148

## Input For CFA With Covariates With Categorical Outcomes For 15 ASB Items

```

TITLE:      CFA with covariates with categorical outcomes using
            15 antisocial behavior items and 3 covariates

DATA:      FILE IS asb.dat;
            FORMAT IS 34X 54F2.0;

VARIABLE:  NAMES ARE property fight shoplift lt50 gt50 force
            threat injure pot drug soldpot solddrug con auto bldg
            goods gambling dsml-dsm22 sex black hisp single
            divorce dropout college onset fhist1 fhist2 fhist3
            age94 cohort dep abuse;

            USEV ARE property-gt50 threat-goods sex black age94;

            CATEGORICAL ARE property-goods;

```

149

## Input For CFA With Covariates With Categorical Outcomes For 15 ASB Items (Continued)

```

MODEL:      f1 BY property shoplift-gt50 con-
            goods;

            f2 BY fight threat injure;

            f3 BY pot-solddrug;

            f1-f3 ON sex black age94;

            property-goods ON sex-age94@0;

OUTPUT:     STANDARDIZED MODINDICES;

```

150

**Output Excerpts CFA With Covariates With  
Categorical Outcomes For 15 ASB Items**

**Model Results**

		Estimates	S.E.	Est./S.E.	Std	StdYX
F1	BY					
	PROPERTY	1.000	.000	.000	.791	.760
	SHOPLIFT	.974	.023	42.738	.771	.742
	LT50	.915	.023	39.143	.724	.700
	GT50	1.055	.031	33.658	.835	.799
	CON	.752	.024	31.637	.595	.581
	AUTO	.796	.030	26.462	.629	.613
	BLDG	1.084	.030	35.991	.858	.818
	GOODS	1.071	.025	42.697	.847	.809

151

**Output Excerpts CFA With Covariates With  
Categorical Outcomes For 15 ASB Items (Continued)**

F2	BY					
	FIGHT	1.000	.000	.000	.773	.734
	THREAT	1.096	.035	31.382	.847	.797
	INJURE	1.082	.037	28.888	.836	.787
F3	BY					
	POT	1.000	.000	.000	.866	.851
	DRUG	1.031	.023	45.818	.893	.876
	SOLDPOT	1.046	.023	45.844	.905	.888
	SOLDDRUG	.923	.036	25.684	.799	.787

152

**Output Excerpts CFA With Covariates With  
Categorical Outcomes For 15 ASB Items (Continued)**

F1	ON					
	SEX	.516	.024	21.206	.653	.326
	BLACK	-.080	.025	-3.168	-.102	-.047
	AGE94	-.054	.006	-9.856	-.069	-.150
F2	ON					
	SEX	.561	.026	21.715	.726	.363
	BLACK	.174	.025	7.087	.225	.103
	AGE94	-.068	.006	-12.286	-.087	-.191
F3	ON					
	SEX	.229	.026	8.760	.265	.132
	BLACK	-.272	.029	-9.384	-.315	-.144
	AGE94	.039	.006	6.481	.045	.099

153

**Output Excerpts CFA With Covariates With  
Categorical Outcomes For 15 ASB Items (Continued)**

**Tests Of Model Fit**

Chi-Square Test of Model Fit	
Value	1225.266*
Degrees of Freedom	105**
P-Value	0.0000
CFI / TLI	
CFI	0.945
TLI	0.964
RMSEA (Root Mean Square Error Of Approximation)	
Estimate	0.038
WRMR (Weighted Root Mean Square Residual)	
Value	2.498

154

**Output Excerpts CFA With Covariates With  
Categorical Outcomes For 15 ASB Items (Continued)**

**Modification Indices**

PROPERTY ON BLACK	4.479	GT50 ON SEX	12.100
PROPERTY ON AGE94	28.229	GT50 ON BLACK	12.879
FIGHT ON SEX	60.599	GT50 ON AGE94	7.413
FIGHT ON BLACK	26.695	THREAT ON SEX	10.221
<b>FIGHT ON AGE94</b>	<b>64.815</b>	THREAT ON BLACK	26.665
<b>SHOPLIFT ON SEX</b>	<b>131.792</b>	THREAT ON AGE94	3.892
SHOPLIFT ON BLACK	0.039	INJURE ON SEX	22.803
SHOPLIFT ON AGE94	0.038	INJURE ON BLACK	0.089
LT50 ON SEX	0.040	INJURE ON AGE94	42.549
LT50 ON BLACK	22.530	POT ON SEX	10.727
LT50 ON AGE94	24.750	POT ON BLACK	12.177
		POT ON AGE94	17.432

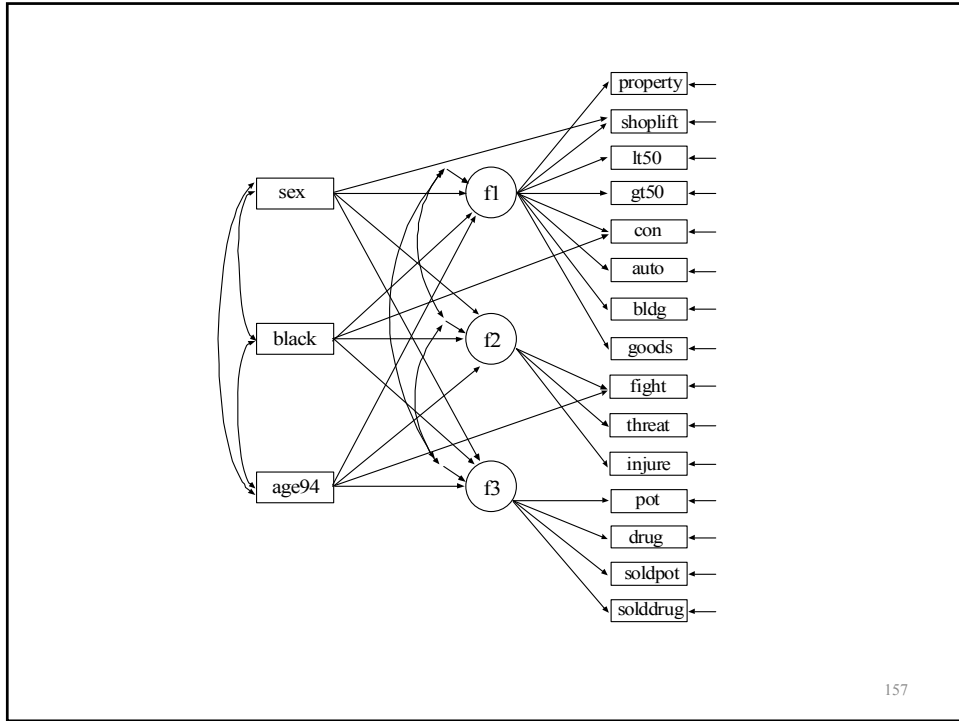
155

**Output Excerpts CFA With Covariates With  
Categorical Outcomes For 15 ASB Items (Continued)**

**Modification Indices**

DRUG ON SEX	15.637	AUTO ON SEX	0.735
DRUG ON BLACK	41.202	AUTO ON BLACK	1.414
DRUG ON AGE94	1.583	AUTO ON AGE94	2.936
SOLDPOT ON SEX	51.496	BLDG ON SEX	37.797
SOLDPOT ON BLACK	1.242	BLDG ON BLACK	7.053
SOLDPOT ON AGE94	29.267	BLDG IB AGE94	0.114
SOLDDRUG ON SEX	3.920	GOODS ON SEX	24.664
SOLDDRUG ON BLACK	7.187	GOODS ON BLACK	0.982
SOLDDRUG ON AGE94	2.956	GOODS ON AGE94	6.061
CON ON SEX	31.521		
<b>CON ON BLACK</b>	<b>80.515</b>		
CON ON AGE94	11.259		

156



## Input Excerpts For ASB CFA With Covariates And Direct Effects

```

MODEL:
    f1 BY property shoplift-gt50 con-goods;
    f2 BY fight threat injure;
    f3 BY pot-solddrug;

    f1-f3 ON sex black age94;

    shoplift ON sex;
    con ON black;
    fight ON age94;
    
```

**Input Excerpts For ASB CFA  
With Covariates And Direct Effects (Continued)**

**Tests Of Model Fit**

Chi-Square Test of Model Fit		
Value		946.256 *
Degrees of Freedom		102 **
P-Value		0.0000
CFI/TLI		
CFI		0.959
TLI		0.972
RMSEA (Root Mean Square Error Of Approximation)		
Estimate		0.034
WRMR (Weighted Root Mean Square Residual)		
Value		2.198

159

**Output Excerpts For ASB CFA  
With Covariates And Direct Effects (Continued)**

		Estimates	S.E.	Est./S.E.	Std	StdYX
F1	BY					
	SHOPLIFT	1.002	.024	42.183	.805	.793
F1	ON					
	SEX	.596	.026	22.958	.742	.371
SHOPLIFT	ON					
	SEX	-.385	.033	-11.594	-.385	-.190
CON	ON					
	BLACK	.305	.034	8.929	.305	.136
FIGHT	ON					
	AGE94	-.068	.008	-8.467	-.068	-.138
Thresholds						
	SHOPLIFT\$1	.558	.033	17.015	.558	.558
R-SQUARE						
	Observed Variable	Residual Variance		R-Square		
	SHOPLIFT	.461		.552		

160



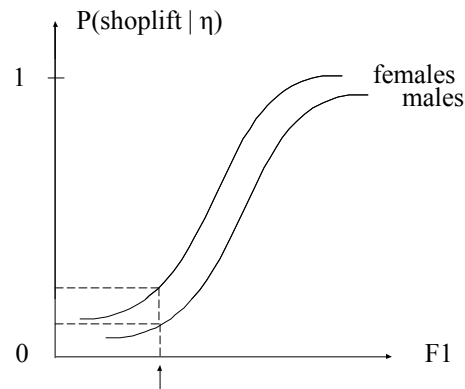
## Interpretation Of Direct Effects

### Shoplift On Gender

- Indirect effect of gender on shoplift
  - F1 has a positive relationship with gender – males have a higher mean than females on the f1 factor
  - Shoplift has a positive loading on the f1 factor
  - Conclusion: males are expected to have a higher probability of shoplifting
- Effect of gender on shoplift
  - Direct effect is negative – for a given factor value, males have a lower probability of shoplifting than females
  - Conclusion – shoplift is not invariant

161

## Calculating Item Probabilities



Graph can be done in Mplus using the PLOT command and the option "Item characteristic curves".

162

## Calculating Item Probabilities (Continued)

The model with a direct effect from  $x$  to item  $u_j$ ,

$$u_{ij}^* = \lambda_j \eta_i + \kappa_j x_i + \varepsilon_{ij}, \quad (45)$$

gives the conditional probability of a  $u = 1$  response given the factor  $\eta_i$  and the covariate  $x_i$

$$P(u_{ij} = 1 \mid \eta_i, x_i) = 1 - F[(\tau_j - \lambda_j \eta_i - \kappa_j x_i) \theta_{jj}^{-1/2}], \quad (46)$$

$$= F[(-\tau_j + \lambda_j \eta_i + \kappa_j x_i) \theta_{jj}^{-1/2}], \quad (47)$$

where  $F$  is the normal distribution function and  $\theta$  is the residual variance.

For example, for the item shoplift,  $\tau_j = 0.558$ ,  $\kappa_j = -0.385$ ,  $\theta_{jj} = 0.461$ . At  $\eta = 0$ , the probability is 0.21 for females ( $x = 0$ ) and 0.08 for males ( $x = 1$ ).

163

## Calculating Item Probabilities (Continued)

Consider

$$P(u_{ij} = 1 \mid \eta_i, x_i) = 1 - F[(\tau_j - \lambda_j \eta_i - \kappa_j x_i) \theta_{jj}^{-1/2}], \quad (47)$$

using  $\tau_j = 0.558$ ,  $\kappa_j = -0.385$ ,  $\theta_{jj} = 0.461$ , and  $\eta = 0$ .

Here,  $\theta_{jj}^{-1/2} = \frac{1}{\sqrt{\theta_{jj}}} = \frac{1}{\sqrt{0.461}} = 1.473$ .

**For females ( $x = 0$ ):**

1.  $(\tau_j - \lambda_j \eta_i - \kappa_j x_i) = 0.558 - 1.002 \times 0 - (-0.385) \times 0 = 0.558$ .

2.  $(\tau_j - \lambda_j \eta_i - \kappa_j x_i) \theta_{jj}^{-1/2} = 0.558 \times 1.473 = 0.822$ .

3.  $F[0.822] = 0.794$  using a z table

4.  $1 - 0.794 = 0.206$ .

**For males ( $x = 1$ ):**

1.  $(\tau_j - \lambda_j \eta_i - \kappa_j x_i) = 0.558 - 1.002 \times 0 - (-0.385) \times 1 = 0.943$ .

2.  $(\tau_j - \lambda_j \eta_i - \kappa_j x_i) \theta_{jj}^{-1/2} = 0.943 \times 1.473 = 1.389$ .

3.  $F[1.389] = 0.918$  using a z table.

4.  $1 - 0.918 = 0.082$ .

164

## **Further Readings On Factor Analysis And MIMIC Analysis With Categorical Outcomes**

- Gallo, J.J., Anthony, J. & Muthen, B. (1994). Age differences in the symptoms of depression: a latent trait analysis. Journals of Gerontology: Psychological Sciences, 49, 251-264. (#52)
- Mislevy, R. (1986). Recent developments in the factor analysis of categorical variables. Journal of Educational Statistics, 11, 3-31.
- Muthén, B. (1978). Contributions to factor analysis of dichotomous variables. Psychometrika, 43, 551-560. (#3)
- Muthén, B. (1989). Dichotomous factor analysis of symptom data. In Eaton & Bohrnstedt (Eds.), Latent variable models for dichotomous outcomes: Analysis of data from the Epidemiological Catchment Area Program (pp. 19-65), a special issue of Sociological Methods & Research, 18, 19-65. (#21)

165

## **Further Readings On Factor Analysis And MIMIC Analysis With Categorical Outcomes (Continued)**

- Muthén, B. (1989). Latent variable modeling in heterogeneous populations. Psychometrika, 54, 557-585. (#24)
- Muthén, B., Tam, T., Muthén, L., Stolzenberg, R. M., & Hollis, M. (1993). Latent variable modeling in the LISCOMP framework: Measurement of attitudes toward career choice. In D. Krebs, & P. Schmidt (Eds.), New directions in attitude measurement, Festschrift for Karl Schuessler (pp. 277-290). Berlin: Walter de Gruyter. (#46)

166

## Multiple Group Analysis With Categorical Outcomes

167

### Steps In Multiple Group Analysis

- Fit the model separately in each group
- Fit the model in all groups allowing all parameters to be free except factor means which are fixed to zero in all groups and scale factors which are fixed to one in all groups
- Fit the model in all groups holding factor loadings and thresholds equal across groups with factor means fixed to zero in the first group and free in the other groups and scale factors fixed to one in the first group and free in the other groups
- Add covariates
- Modify the model

168

## Inputs For Multiple Group Analysis Of 15 ASB Items

### Measurement Non-Invariance

```
MODEL:      f1 BY property shoplift-gt50 con-goods;
            f2 BY fight threat injure;
            f3 BY pot-solddrug;
            [f1-f3@0];
            {property-goods@1};
```

```
MODEL male: f1 BY shoplift-gt50 con-goods;
            f2 BY threat injure;
            f3 BY drug-solddrug;
            [property$1-goods$1];
```

169

## Inputs For Multiple Group Analysis Of 15 ASB Items (Continued)

### Measurement Invariance

```
MODEL:      f1 BY property shoplift-gt50 con-goods;
            f2 BY fight threat injure;
            f3 BY pot-solddrug;
```

### Partial Measurement Invariance

```
MODEL:      f1 BY property shoplift-gt50 con-goods;
            f2 BY fight* threat@1 injure;
            f3 BY pot-solddrug;
```

```
MODEL      f1 BY con lt50;
male:      f2 BY fight;
            f3 BY soldpot pot solddrug;
            [con$1 lt50$1 fight$1 soldpot$1 pot$1 solddrug$1];
            {con@1 lt50@1 fight@1 soldpot@1 pot@1 solddrug@1};
```

170

## **Further Readings On Multiple-Group Analysis Of Categorical Outcomes**

Muthén, B. & Asparouhov, T. (2002). Latent variable analysis with categorical outcomes: Multiple-group and growth modeling in Mplus. Mplus Web Note #4 ([www.statmodel.com](http://www.statmodel.com)).

Muthén, B., & Christofferson, A. (1981). Simultaneous factor analysis of dichotomous variables in several groups. *Psychometrika*, 46, 407-419. (#6)

171

## **Exploratory Structural Equation Modeling**

172

## Overview

- Brief overview of EFA, CFA, and SEM for continuous outcomes
- New approach to structural equation modeling
- Examples

173

## Factor Analysis And Structural Equation Modeling

- Exploratory factor analysis (EFA) is one of the most frequently used multivariate analysis technique in statistics
- 1966 Jennrich solved a significant EFA rotation problem by deriving the direct quartimin rotation
- Jennrich was the first to develop standard errors for rotated solutions although these have still not made their way into most statistical software programs
- 1969 development of confirmatory factor analysis (CFA) by Joreskog
- Joreskog developed CFA further into structural equation modeling (SEM) in LISREL where CFA was used for the measurement part of the model

174

## Structural Equation Model

$$(1) Y_i = \nu + \Lambda \eta_i + K X_i + \varepsilon_i$$

$$(2) \eta_i = \alpha + B \eta_i + \Gamma X_i + \xi_i$$

$\Lambda$  is typically specified as having a "simple structure"

175

## CFA Simple Structure $\Lambda$

$$A = \begin{pmatrix} X & 0 \\ X & 0 \\ X & 0 \\ 0 & X \\ 0 & X \\ 0 & X \end{pmatrix}$$

where X is a factor loading parameter to be estimated

- CFA simple structure is often too restrictive in practice

176



### Quote From Browne (2001)

"Confirmatory factor analysis procedures are often used for exploratory purposes. Frequently a confirmatory factor analysis, with pre-specified loadings, is rejected and a sequence of modifications of the model is carried out in an attempt to improve fit. The procedure then becomes exploratory rather than confirmatory --- In this situation the use of exploratory factor analysis, with rotation of the factor matrix, appears preferable. --- The discovery of misspecified loadings ... is more direct through rotation of the factor matrix than through the examination of model modification indices."

Browne, M.W. (2001). An overview of analytic rotation in exploratory factor analysis. Multivariate Behavioral Research, 36 , 111-150

177

### A New Approach: Exploratory SEM

- Allow EFA measurement model parts (EFA sets)
- Integrated with CFA measurement parts
- Allowing EFA sets access to other SEM parameters, such as
  - Correlated residuals
  - Regressions on covariates
  - Regressions between factors of different EFA sets
  - Regressions between factors of EFA and CFA sets
  - Multiple groups
  - EFA loading matrix equalities across time or group
  - Mean structures
- Available for continuous, categorical, and censored outcomes

178

## Factor Indeterminacy And Rotations

- $\Lambda \Psi \Lambda^T + \Theta$
- $\Lambda$  is  $p \times m$ , so  $m^2$  indeterminacies
- $\Psi = I$  fixes  $m(m+1)/2$  indeterminacies
- $\Lambda \Lambda^T + \Theta = \Lambda^* \Lambda^{*T} + \Theta$   
for  $\Lambda^* = \Lambda H^{-1}$ , where  $H$  is orthogonal
- A starting  $\Lambda^*$  can be rotated using a rotation criterion function that favors simple structure in  $\Lambda$ :

$$f(\Lambda^*) = f(\Lambda H^{-1}) \quad (2a)$$

$$f(\Lambda) = \sum_{i=1}^p \sum_{j=1}^m \sum_{k \neq j}^m \lambda_{ij}^2 \lambda_{ik}^2 \quad (2b)$$

- Common rotation: Quartimin
- Good alternative: Geomin rotation

179

## Rotation Methods

Choice of rotation important when not relying on CFA measurement structure:

- With variable complexity  $> 1$  (“cross-loadings”) Geomin is better than conventional methods such as varimax, promax, quartimin
- Target rotation

180

## Target Rotation

Target rotation:

- Between mechanical rotation and CFA: Rotation guided by judgment
- Choose rotation by specifying target loading values (typically zero)
- Target values not fixed as in CFA – zero targets can come out big if misspecified
- $m - 1$  zeros in each loading column gives EFA ( $m = \#$  factors)
- Mplus language:

```
f1 BY y1-y10 y1~0 (*t);
```

```
f2 BY y1-y10 y5~0 (*t);
```

References: Browne (1972 a, b; Tucker, 1944)

181

## Transformation Of SEM Parameters Based On Rotated $\Lambda$

$$(1) Y_i = v + \Lambda \eta_i + K X_i + \varepsilon_i \quad (2) \quad \eta_i = \alpha + B \eta_i + \Gamma X_i + \xi_i$$

Transformations:

$$(6) v^* = v$$

$$(10) \alpha^* = H \alpha$$

$$(7) \Lambda^* = \Lambda (H^*)^{-1}$$

$$(11) B^* = H^* B (H^*)^{-1}$$

$$(8) K^* = K$$

$$(12) \Gamma^* = H^* \Gamma$$

$$(9) \theta^* = \theta$$

$$(13) \Psi^* = (H^*)^T \Psi H^*$$

182

## Maximum-Likelihood Estimation And Testing

- ML estimation in several steps
  - Compute the unstandardized starting values for  $\Lambda$ ,  $\Psi$ , and  $\Theta$  with identifying restrictions
  - Use the  $\Delta$  method to estimate the asymptotic distribution of the standardized starting value for  $\Lambda$
  - Find the asymptotic distribution of the rotated standardized solution (cf Jennrich, 2003)
- Standard errors for rotated solution of the full SEM
- Pre-specified testing sequence: EFA followed by CFA

183

## Examples

- MIMIC with cross-loadings (see Web Talks)
- Longitudinal EFA (test-retest) (see Web Talks)
- Multiple-group EFA

184

### **Example: Aggressive Behavior Male-Female EFA in Baltimore Cohort 3**

- 261 males and 248 females in third grade
- Teacher-rated aggressive-disruptive behavior
- Outcomes treated as non-normal continuous variables
- Two types of analyses:
  - EFA in each group separately using Geomin rotation
  - Multiple-group EFA analysis of males and females jointly

185

### **EFA-ESEM Variable Scales And Loading Matrix Metrics**

- Sample covariance matrix analyzed, not sample correlation matrix
  - Loadings in original indicator scale
  - Standardized solution gives loadings in regular EFA metric
- Multiple-group EFA allows factor variances and covariances to differ across groups as the default
  - Group 1 has a factor correlation matrix, while other groups have factor covariance matrices
  - Group-invariant loadings still give group-varying standardized loadings due to group-varying indicator variances and group-varying factor variances

186

### Summary Of Separate Male/Female EFAs

Variables	StdYX Loadings for Males			StdYX Loadings for Females		
	Verbal	Person	Property	Verbal	Person	Property
Stubborn	<b><u>0.82</u></b>	-0.05	0.01	<b><u>0.88</u></b>	0.03	-0.22
Breaks Rules	<u>0.47</u>	<u>0.34</u>	0.01	<b><u>0.76</u></b>	0.06	-0.17
Harms Others & Property	-0.01	<b><u>0.63</u></b>	<u>0.31</u>	<u>0.45</u>	0.03	0.36
Breaks Things	-0.02	0.02	<b><u>0.66</u></b>	-0.02	0.19	<b><u>0.43</u></b>
Yells At Others	<b><u>0.66</u></b>	0.23	-0.03	<b><u>0.97</u></b>	-0.23	0.05
Takes Others' Property	<u>0.27</u>	0.08	<b><u>0.52</u></b>	0.02	<b><u>0.79</u></b>	0.10
Fights	<u>0.22</u>	<b><u>0.75</u></b>	-0.00	<b><u>0.81</u></b>	-0.01	0.18
Harms Property	0.03	-0.02	<b><u>0.93</u></b>	0.27	0.20	<b><u>0.57</u></b>
Lies	<b><u>0.58</u></b>	0.01	<u>0.27</u>	<u>0.42</u>	<u>0.50</u>	-0.00
Talks Back to Adults	<b><u>0.61</u></b>	-0.02	<u>0.30</u>	<b><u>0.69</u></b>	0.09	-0.02
Teases Classmates	<u>0.46</u>	<u>0.44</u>	-0.04	<b><u>0.71</u></b>	-0.01	0.10
Fights With Classmates	<u>0.30</u>	<b><u>0.64</u></b>	0.08	<b><u>0.83</u></b>	0.03	<u>0.21</u>
Loses Temper	<b><u>0.64</u></b>	<u>0.16</u>	0.04	<b><u>1.05</u></b>	-0.29	-0.01

187

### Summary Of Separate Male/Female EFAs

Factors	Factor Correlations for Males		Factor Correlations for Females	
	Verbal	Person	Verbal	Person
Person	0.57		0.68	
Property	0.56	0.68	0.32	0.22

188

## Multiple-Group EFA Modeling Results Using MLR

Model	LL0	C	# par. 's	Df	$\chi^2$	CFI	RMSEA
M1	-8122	2.61	84	124	241	0.95	0.061
M2	-8087	2.41	94	114	188	0.97	0.050
M3	-8036	2.38	124	84	146	0.97	0.054

- M1: Loadings and intercepts invariance
- M2: Loadings but not intercepts invariance
- M3: Neither loadings nor intercepts invariance
- LL0: Log likelihood for the H0 (multiple-group EFA) model
- c is a non-normality scaling correction factor

189

## Multiple-Group EFA Modeling Results Using MLR

- Comparing M2 and M1\*:
  - $cd = (84 \cdot 2.61 - 94 \cdot 2.41) / (-10) = 0.704$
  - $TRd = -2(LL0 - LL1) / cd = 98.5$  with 10 df: Not all intercepts are invariant. Choose M2

190

## Multiple-Group EFA Modeling Results Using MLR

- Comparing M3 and M2\*:
  - $cd = (94*2.41-124*2.38)/(-30) = 2.78$
  - $TRd = -2(LL0-LL1)/cd = 36.6$  with 30 df: Loadings are invariant. Choose M2
- $LL1 =$  loglikelihood for unrestricted H1 model (same for all 3) = -7934

\* For loglikelihood difference testing with scaling corrections, see <http://www.statmodel.com/chidiff.shtml>

191

## Male EFA Estimates Compared To Female Estimates From Multiple-Group EFA Using M2

Variables	StdYX Loadings for Males			StdYX Loadings for Females		
	Verbal	Person	Property	Verbal	Person	Property
Stubborn	<b><u>0.82</u></b>	-0.05	0.01	<b><u>0.86</u></b>	-0.00	-0.01
Breaks Rules	<u>0.47</u>	<u>0.34</u>	0.01	<b><u>0.59</u></b>	<u>0.20</u>	0.01
Harms Others & Property	-0.01	<b><u>0.63</u></b>	<u>0.31</u>	0.00	<b><u>0.56</u></b>	<u>0.24</u>
Breaks Things	-0.02	0.02	<b><u>0.66</u></b>	-0.03	-0.03	<b><u>0.63</u></b>
Yells At Others	<b><u>0.66</u></b>	0.23	-0.03	<b><u>0.69</u></b>	0.18	-0.01
Takes Others' Property	<u>0.27</u>	0.08	<u>0.52</u>	<u>0.39</u>	0.03	<u>0.31</u>
Fights	<u>0.22</u>	<b><u>0.75</u></b>	-0.00	<u>0.35</u>	<u>0.61</u>	-0.02
Harms Property	0.03	-0.02	<b><u>0.93</u></b>	0.19	0.04	<b><u>0.68</u></b>
Lies	<b><u>0.58</u></b>	0.01	<u>0.27</u>	<b><u>0.67</u></b>	0.00	<u>0.16</u>
Talks Back to Adults	<b><u>0.61</u></b>	-0.02	<u>0.30</u>	<b><u>0.71</u></b>	-0.02	<u>0.15</u>
Teases Classmates	<u>0.46</u>	<u>0.44</u>	-0.04	<u>0.49</u>	<u>0.30</u>	0.01
Fights With Classmates	<u>0.30</u>	<b><u>0.64</u></b>	0.08	<u>0.41</u>	<u>0.53</u>	0.03
Loses Temper	<u>0.64</u>	<u>0.16</u>	0.04	<b><u>0.74</u></b>	0.14	-0.29

192



**Factor Correlations For  
Males Using EFA And For Females Using Multiple-  
Group Model M2**

Factors	Factor Correlations for Males		Factor Correlations for Females	
	Verbal	Person	Verbal	Person
Person	0.57		0.75	
Property	0.56	0.68	0.42	0.65

193

**Multiple-Group EFA Estimates For M2**

Group	Factor Variances		
	Verbal	Person	Property
Males	1	1	1
Females	1.19 (.18)	2.65 (.56)	5.33 (1.02)

194

## Input Model M1

```

TITLE:      Cohort 3 Case and Class variables
DATA:      FILE = Muthen.dat;
VARIABLE:  NAMES = id race lunch312 gender y301-y313;
           MISSING = ALL (999);
           GROUPING = gender (0=female 1=male);
           USEVARIABLES = y301-y313;
ANALYSIS:  PROCESSORS = 4;
           ESTIMATOR = MLR;
MODEL:     f1-f3 BY y301-y313 (*1);
OUTPUT:    TECH1 SAMPSTAT MODINDICES STANDARDIZED;

```

195

## Input Model M2

```

TITLE:      Cohort 3 Case and Class variables
DATA:      FILE = Muthen.dat;
VARIABLE:  NAMES = id race lunch312 gender y301-y313;
           MISSING = ALL (999);
           GROUPING = gender (0=female 1=male);
           USEVARIABLES = y301-y313;
ANALYSIS:  PROCESSORS = 4;
           ESTIMATOR = MLR;
MODEL:     f1-f3 BY y301-y313 (*1);
           [f1-f3@0];
MODEL MALE: [y301-y313];
OUTPUT:    TECH1 SAMPSTAT MODINDICES STANDARDIZED;

```

196

## Input Model M3

```

TITLE:      Cohort 3 Case and Class variables
DATA:      FILE = Muthen.dat;
VARIABLE:  NAMES = id race lunch312 gender y301-y313;
           MISSING = ALL (999);
           GROUPING = gender (0=female 1=male);
           USEVARIABLES = y301-y313;
ANALYSIS:  PROCESSORS = 4;
           ESTIMATOR = MLR;
MODEL:     f1-f3 BY y301-y313 (*1);
           [f1-f3@0];
MODEL MALE: f1-f3 BY y301-y313 (*1);
           [y301-y313];
OUTPUT:    TECH1 Sampstat Modindices Standardized;

```

197

## Further Readings On ESEM

Asparouhov, T. & Muthén, B. (2008). Exploratory structural equation modeling. Forthcoming in Structural Equation Modeling.

Marsh, H.W., Muthén, B., Asparouhov, A., Lüdtke, O., Robitzsch, A., Morin, A.J.S., & Trautwein, U. (2009). Exploratory Structural Equation Modeling, Integrating CFA and EFA: Application to Students' Evaluations of University Teaching. Forthcoming in Structural Equation Modeling.

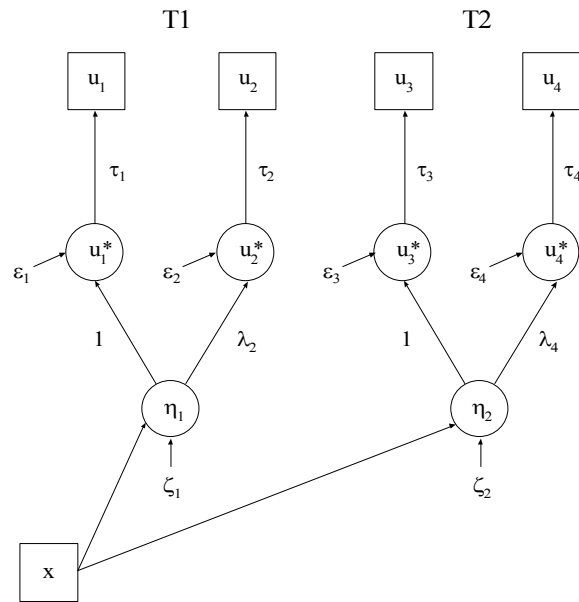
Web talk: Exploratory structural equation modeling. See <http://www.statmodel.com/webtalks.shtml>

Version 5.1 Language Addendum and Examples Addendum covering ESEM. See <http://www.statmodel.com/ugexcerpts.shtml>

198

**Technical Issues For  
Weighted-Least Squares Estimation**

199



200

## Latent Response Variable Modeling

- The analysis considers means (thresholds) and correlations because variances do not contribute further information
  - $E(u) = \pi, V(u) = \pi(1 - \pi)$
- For each  $u$  (see figure)
  - Normality of  $u^*$  given  $x$  (probit)
  - Residual variance fixed at 1 implies  $V(\varepsilon)$  not free,
 
$$V(u^* | x) = \lambda^2 V(\zeta) + V(\varepsilon) = 1, \quad (8)$$

$$i.e. V(\varepsilon) = 1 - \lambda^2 V(\zeta) \quad (9)$$
- For pairs of  $u$ 's
  - Multivariate normal  $u^*$ 's given  $x$
  - Because residual variances are one,  $u^*$  residual correlations are considered, not covariances
  - Normality of  $u^*$ 's given  $x$  is less strong than normal  $u^*$  and normal  $x$ , assumed for polychoric and polyserial correlations

201

## Scale Factors With Measurement Invariance

Problem: Correlations should not be used when comparing relationships for variables with different variances.

Solution: Add scale factors  $\delta$  to the model,  $\delta = 1/\sqrt{V(u^* | x)}$ .

Example (see figure): Aim is to test measurement invariance, e.g.

$$\tau_2 = \tau_4 = \tau, \lambda_2 = \lambda_4 = \lambda.$$

$$V(u_2^* | x) = \lambda^2 V(\zeta_1) + V(\varepsilon_2), \quad (40)$$

$$V(u_4^* | x) = \lambda^2 V(\zeta_2) + V(\varepsilon_4), \quad (41)$$

showing that  $V(u^* | x)$  varies across the two variables if either  $V(\zeta)$  or  $V(\varepsilon)$  varies, even though  $\lambda$  is invariant.

Fixing both  $V(u_2^* | x)$  and  $V(u_4^* | x)$  to 1 is therefore wrong under measurement invariance. Instead, use

$$\delta_2 = 1, \quad (42)$$

$$\delta_4 \text{ free.} \quad (43)$$

By letting  $\delta_4$  be free, the model allows  $V(u_4^* | x) \neq V(u_2^* | x)$ , while still modeling the  $u_2^*, u_4^*$  correlation

$$\text{Cov}(u_2^*, u_4^* | x) \delta_2 \delta_4. \quad (44)$$

202

## Estimation With Categorical Outcomes

Full information maximum-likelihood estimation is heavy for general models.

### Limited-information weighted least squares:

Fitting function:

$$WLS = 1/2 (\mathbf{s} - \boldsymbol{\sigma})' \mathbf{W}^{-1} (\mathbf{s} - \boldsymbol{\sigma})$$

Sample statistics:

- $\mathbf{s}_1$ : probit thresholds
- $\mathbf{s}_2$ : probit regression slopes ( $q > 0$ )
- $\mathbf{s}_3$ : probit residual correlations
- $\mathbf{s}' = (\mathbf{s}_1', \mathbf{s}_2', \mathbf{s}_3')$

Weight matrix:

- Full  $\mathbf{W}$  (GLS/WLS:  $\mathbf{W} = \text{asympt } V(\mathbf{s})$ )
- Diagonal  $\mathbf{W}$  (WLSM, WLSMV)

Robust standard errors and chi-square in line with Satorra

203

## Further Readings On Technical Aspects Of Weighted Least Squares With Categorical Outcomes

- Muthén, B. (1984). A general structural equation model with dichotomous, ordered categorical, and continuous latent variable indicators. *Psychometrika*, 49, 115-132. (#11)
- Muthén, B. (1989). Latent variable modeling in heterogeneous populations. *Psychometrika*, 54, 557-585. (#24)
- Muthén, B. & Satorra, A. (1995). Technical aspects of Muthén's LISCOMP approach to estimation of latent variable relations with a comprehensive measurement model. *Psychometrika*, 60, 489-503.
- Muthén, B. du Toit, S.H.C. & Spisic, D. (1997). Robust inference using weighted least squares and quadratic estimating equations in latent variable modeling with categorical and continuous outcomes. Accepted for publication in *Psychometrika*. (#75)

204

## Levels Of Engagement

- Mplus support for licensed Mplus users
- Mplus Discussion for brief Mplus analysis questions of general interest
- Statistical consulting not available through Mplus
- Research interaction on topics of common interest
- SEMNET

205

## References

### **Analysis With Categorical Outcomes**

#### **General**

- Agresti, A. (2002). Categorical data analysis. Second edition. New York: John Wiley & Sons.
- Agresti, A. (1996). An introduction to categorical data analysis. New York: Wiley.
- Hosmer, D.W. & Lemeshow, S. (2000). Applied logistic regression. Second edition. New York: John Wiley & Sons.
- McKelvey, R.D. & Zavoina, W. (1975). A statistical model for the analysis of ordinal level dependent variables. Journal of Mathematical Sociology, 4, 103-120.

#### **Censored and Poisson Regression**

- Hilbe, J. M. (2007). Negative binomial regression. Cambridge, UK: Cambridge University Press.
- Lambert, D. (1992). Zero-inflated Poisson regression, with an application to defects in manufacturing. Technometrics, 34, 1-13.

206

## References (Continued)

- Long, S. (1997). Regression models for categorical and limited dependent variables. Thousand Oaks: Sage.
- Maddala, G.S. (1983). Limited-dependent and qualitative variables in econometrics. Cambridge: Cambridge University Press.
- Tobin, J (1958). Estimation of relationships for limited dependent variables. Econometrica, 26, 24-36.

### IRT

- Baker, F.B. & Kim, S.H. (2004). Item response theory. Parameter estimation techniques. Second edition. New York: Marcel Dekker.
- Bock, R.D. (1997). A brief history of item response theory. Educational Measurement: Issues and Practice, 16, 21-33.
- du Toit, M. (2003). IRT from SSI. Lincolnwood, IL: Scientific Software International, Inc. (BILOG, MULTILOG, PARSCALE, TESTFACT)
- Embretson, S. E., & Reise, S. P. (2000). Item response theory for psychologists. Mahwah, NJ: Erlbaum.
- Hambleton, R.K. & Swaminathan, H. (1985). Item response theory. Boston: Kluwer-Nijhoff.
- MacIntosh, R. & Hashim, S. (2003). Variance estimation for converting MIMIC model parameters to IRT parameters in DIF analysis. Applied Psychological Measurement, 27, 372-379.

207

## References (Continued)

- Muthén, B. (1985). A method for studying the homogeneity of test items with respect to other relevant variables. Journal of Educational Statistics, 10, 121-132. (#13)
- Muthén, B. (1988). Some uses of structural equation modeling in validity studies: Extending IRT to external variables. In H. Wainer & H. Braun (Eds.), Test Validity (pp. 213-238). Hillsdale, NJ: Erlbaum Associates. (#18)
- Muthén, B. (1989). Using item-specific instructional information in achievement modeling. Psychometrika, 54, 385-396. (#30)
- Muthén, B. (1994). Instructionally sensitive psychometrics: Applications to the Second International Mathematics Study. In I. Westbury, C. Ethington, L. Sosniak & D. Baker (Eds.), In search of more effective mathematics education: Examining data from the IEA second international mathematics study (pp. 293-324). Norwood, NJ: Ablex. (#54)
- Muthén, B. & Asparouhov, T. (2002). Latent variable analysis with categorical outcomes: Multiple-group and growth modeling in Mplus. Mplus Web Note #4 ([www.statmodel.com](http://www.statmodel.com)).
- Muthén, B., Kao, Chih-Fen & Burstein, L. (1991). Instructional sensitivity in mathematics achievement test items: Applications of a new IRT-based detection technique. Journal of Educational Measurement, 28, 1-22. (#35)

208



## References (Continued)

- Muthén, B. & Lehman, J. (1985). Multiple-group IRT modeling: Applications to item bias analysis. Journal of Educational Statistics, 10, 133-142. (#15)
- Takane, Y. & DeLeeuw, J. (1987). On the relationship between item response theory and factor analysis of discretized variables. Psychometrika, 52, 393-408.

### Factor Analysis

- Bartholomew, D.J. (1987). Latent variable models and factor analysis. New York: Oxford University Press.
- Bock, R.D., Gibbons, R., & Muraki, E.J. (1988). Full information item factor analysis. Applied Psychological Measurement, 12, 261-280.
- Blafield, E. (1980). Clustering of observations from finite mixtures with structural information. Unpublished doctoral dissertation, Jyvaskyla studies in computer science, economics, and statistics, Jyvaskyla, Finland.
- Browne, M.W. (2001). An overview of analytic rotation in exploratory factor analysis. Multivariate Behavioral Research, 36, 111-150
- Flora, D.B. & Curran P.J. (2004). An empirical evaluation of alternative methods of estimation for confirmatory factor analysis with ordinal data. Psychological Methods, 9, 466-491.

209

## References (Continued)

- Lord, F.M. & Novick, M.R. (1968). Statistical theories of mental test scores. Reading, Mass.: Addison-Wesley Publishing Co.
- Millsap, R.E. & Yun-Tien, J. (2004). Assessing factorial invariance in ordered-categorical measures. Multivariate Behavioral Research, 39, 479-515.
- Mislevy, R. (1986). Recent developments in the factor analysis of categorical variables. Journal of Educational Statistics, 11, 3-31.
- Muthén, B. (1978). Contributions to factor analysis of dichotomous variables. Psychometrika, 43, 551-560. (#3)
- Muthén, B. (1989). Dichotomous factor analysis of symptom data. In Eaton & Bohrnstedt (Eds.), Latent variable models for dichotomous outcomes: Analysis of data from the Epidemiological Catchment Area program (pp. 19-65), a special issue of Sociological Methods & Research, 18, 19-65. (#21)
- Muthén, B. (1989). Latent variable modeling in heterogeneous populations. Psychometrika, 54, 557-585. (#24)
- Muthén, B. (1996). Psychometric evaluation of diagnostic criteria: Application to a two-dimensional model of alcohol abuse and dependence. Drug and Alcohol Dependence, 41, 101-112. (#66)

210

## References (Continued)

- Muthén, B. & Asparouhov, T. (2002). Latent variable analysis with categorical outcomes: Multiple-group and growth modeling in Mplus. Mplus Web Note #4 ([www.statmodel.com](http://www.statmodel.com)).
- Muthén, B. & Christoffersson, A. (1981). Simultaneous factor analysis of dichotomous variables in several groups. *Psychometrika*, 46, 407-419. (#6)
- Muthén, B. & Kaplan, D. (1985). A comparison of some methodologies for the factor analysis of non-normal Likert variables. *British Journal of Mathematical and Statistical Psychology*, 38, 171-189.
- Muthén, B. & Kaplan, D. (1992). A comparison of some methodologies for the factor analysis of non-normal Likert variables: A note on the size of the model. *British Journal of Mathematical and Statistical Psychology*, 45, 19-30.
- Muthén, B. & Satorra, A. (1995). Technical aspects of Muthén's LISCOMP approach to estimation of latent variable relations with a comprehensive measurement model. *Psychometrika*, 60, 489-503.
- Takane, Y. & DeLeeuw, J. (1987). On the relationship between item response theory and factor analysis of discretized variables. *Psychometrika*, 52, 393-408.
- Yung, Y.F. (1997). Finite mixtures in confirmatory factor-analysis models. *Psychometrika*, 62, 297-330.

211

## References (Continued)

### MIMIC

- Gallo, J.J., Anthony, J. & Muthén, B. (1994). Age differences in the symptoms of depression: a latent trait analysis. *Journals of Gerontology: Psychological Sciences*, 49, 251-264. (#52)
- Muthén, B. (1989). Latent variable modeling in heterogeneous populations. *Psychometrika*, 54, 557-585. (#24)
- Muthén, B., Tam, T., Muthén, L., Stolzenberg, R.M. & Hollis, M. (1993). Latent variable modeling in the LISCOMP framework: Measurement of attitudes toward career choice. In D. Krebs & P. Schmidt (Eds.), *New directions in attitude measurement*, Festschrift for Karl Schuessler (pp. 277-290). Berlin: Walter de Gruyter. (#46)

### SEM

- Browne, M.W. & Arminger, G. (1995). Specification and estimation of mean- and covariance-structure models. In G. Arminger, C.C. Clogg & M.E. Sobel (Eds.), *Handbook of statistical modeling for the social and behavioral sciences* (pp. 311-359). New York: Plenum Press.

212

## References (Continued)

- MacKinnon, D.P., Lockwood, C.M., Brown, C.H., Wang, W., & Hoffman, J.M. (2007). The intermediate endpoint effect in logistic and probit regression. Clinical Trials, 4, 499-513.
- Muthén, B. (1979). A structural probit model with latent variables. Journal of the American Statistical Association, 74, 807-811. (#4)
- Muthén, B. (1983). Latent variable structural equation modeling with categorical data. Journal of Econometrics, 22, 48-65. (#9)
- Muthén, B. (1984). A general structural equation model with dichotomous, ordered categorical, and continuous latent variable indicators. Psychometrika, 49, 115-132. (#11)
- Muthén, B. (1989). Latent variable modeling in heterogeneous populations. Psychometrika, 54, 557-585. (#24)
- Muthén, B. (1993). Goodness of fit with categorical and other non-normal variables. In K.A. Bollen, & J.S. Long (Eds.), Testing structural equation models (pp. 205-243). Newbury Park, CA: Sage. (#45).
- Muthén, B. & Speckart, G. (1983). Categorizing skewed, limited dependent variables: Using multivariate probit regression to evaluate the California Civil Addict Program. Evaluation Review, 7, 257-269. (#3)

213

## References (Continued)

- Muthén, B. du Toit, S.H.C. & Spisic, D. (1997). Robust inference using weighted least squares and quadratic estimating equations in latent variable modeling with categorical and continuous outcomes. Accepted for publication in Psychometrika. (#75)
- Prescott, C.A. (2004). Using the Mplus computer program to estimate models for continuous and categorical data from twins. Behavior Genetics, 34, 17-40.
- Xie, Y. (1989). Structural equation models for ordinal variables. Sociological Methods & Research, 17, 325-352.
- Yu, C.Y. (2002). Evaluating cutoff criteria of model fit indices for latent variable models with binary and continuous outcomes. Doctoral dissertation, University of California, Los Angeles. [www.statmodel.com](http://www.statmodel.com).

214