

Application of Statistical Modeling  
– A very short summary of some key points

why multilevel modeling :

1. with clustering structures, observations are correlated. The iid assumption of OLS is violated.
2. Effective sample size of clustering data is smaller.  $SE = \frac{\sigma^2}{\sqrt{n}}$  over more likely to have significant results.  
under multilevel analyses correct for this.

multilevel model specification

Empty model:

$$y_{ij} = \beta_{0j} + \epsilon_{ij} \quad \text{level 1} \quad \epsilon_{ij} \sim N(0, \sigma^2)$$

$$\beta_{0j} = \beta_0 + u_{0j} \quad \text{level 2} \quad u_{0j} \sim N(0, \sigma_{0j}^2)$$

One covariate:

$$y_{ij} = \beta_{0j} + x_{ij}\beta_{1j} + \epsilon_{ij} \quad \text{level 1} \quad \epsilon_{ij} \sim N(0, \sigma^2)$$

$$\begin{aligned} \beta_{0j} &= \beta_0 + u_{0j} \\ \beta_{1j} &= \beta_1 + u_{1j} \end{aligned} \quad \text{level 2} \quad \begin{bmatrix} u_{0j} \\ u_{1j} \end{bmatrix} \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{0j}^2 & \sigma_{01} \\ \sigma_{10} & \sigma_{1j}^2 \end{bmatrix} = D\right)$$

Mixed model:

$$y_{ij} = X_r \beta_{rj} + X_f \beta_k + \epsilon_{ij} \quad \text{level 1}$$

$$\beta_{rj} = \beta_p + u_j$$

$$u_j = \begin{bmatrix} u_{0j} \\ u_{1j} \\ \vdots \\ u_{pj} \end{bmatrix} \sim N\left(\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{0j}^2 & \cdots & \sigma_{0p} \\ \sigma_{10} & \ddots & \vdots \\ \vdots & & \sigma_{pj}^2 \end{bmatrix}\right)$$

where  $X_r$  are variables that we specify random effects for

$X_f$  are variables that we specify fixed effect for

$\beta_k$  are parameters to be estimated

$$\epsilon_{ij} \sim N(0, \sigma^2)$$

$$y_{ij} | u_j \sim N(x_{ij}\beta + x_{ri}u_j, \sigma^2)$$

$$y_j | u_j \sim N_{nj}(x_j\beta + x_{rj}u_j, R_j)$$

$$y_{ij} = X_r \beta_{pj} + X_f \beta_k + \epsilon_{ij}$$

where  $X_r$  are (intercept + covariates) with random effects  
 $X_f$  are covariates with fixed effects  
 $\epsilon_{ij}$  is the random error

### Conditional Specification

$$y_j | u_j \sim N(X_j \beta + X_{rj} u_j, R_j)$$

①  $X_j$  is the matrix of covariates

$$\begin{array}{c} n_j \\ \text{cluster size} \end{array} \left\{ \begin{bmatrix} X_{11} & X_{21} & \cdots & X_{p1} & \cdots & X_{k1} \\ X_{12} & X_{22} & \cdots & X_{p2} & \cdots & X_{k2} \\ \vdots & & & & & \\ X_{1n_j} & X_{2n_j} & \cdots & X_{pn_j} & \cdots & X_{kn_j} \end{bmatrix} \right. \\ \left. \begin{array}{l} p \text{ covariates} \\ \text{with random effect} \end{array} \quad \begin{array}{l} k \text{ covariates} \\ \text{with fixed effect} \end{array} \right\}$$

②  $\beta$  is a vector of coefficients

$$(p+k) \times 1$$

$$\begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \\ \vdots \\ \beta_k \end{bmatrix} \left\{ \begin{array}{l} p \text{ coefficients as the mean of} \\ \text{the random coefficients} \end{array} \right. \\ \left. \begin{array}{l} k \text{ coefficients as fixed effects} \end{array} \right\}$$

③  $X_{rj}$  is part of  $X_j$  — the covariate matrix of the covariates that have random effect

$$\begin{array}{c} n_j \\ \text{cluster size} \end{array} \left\{ \begin{bmatrix} X_{11} & X_{21} & \cdots & X_{p1} \\ X_{12} & X_{22} & \cdots & X_{p2} \\ \vdots & & & \\ X_{1n_j} & X_{2n_j} & \cdots & X_{pn_j} \end{bmatrix} \right. \\ \left. \begin{array}{l} p \text{ covariates with random effect} \end{array} \right\}$$

④  $u_j$  is a vector of the random effect of the cluster

$$p \times 1$$

$$\begin{bmatrix} u_0 \\ u_1 \\ \vdots \\ u_p \end{bmatrix}$$

⑤  $R_j$  is the variance covariance matrix of  $\epsilon_j$

$$\epsilon_j = \begin{bmatrix} \epsilon_{1j} \\ \epsilon_{2j} \\ \vdots \\ \epsilon_{nj} \end{bmatrix} \sim N \left( \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma^2 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \sigma^2 \end{bmatrix} \right)$$

$R_j$

## Marginal specification

$$y_j = \begin{bmatrix} y_{1j} \\ y_{2j} \\ \vdots \\ y_{nj} \end{bmatrix} \quad y_{ij} = X_{ri}\beta_{pj} + X_{fi}\beta_k + \epsilon_{ij} \quad \text{where } \beta_{pj} = \beta_p + u_j$$

At the element level

Variance:

$$\begin{aligned} V(y_{ij}) &= V(X_{ri}\beta_{pj} + X_{fi}\beta_k + \epsilon_{ij}) \\ &= V(X_i\beta + X_{ri}u_j + \epsilon_{ij}) \quad \text{where } X_i = X_{ri} + X_{fi}; \\ &= X_{ri}^2 V(u_j) + V(\epsilon_{ij}) \quad \beta = \beta_p + \beta_k \end{aligned}$$

As for covariance:

$$y_{1j} = X_{1j}\beta + X_{1rj}u_j + \epsilon_{1j} \quad \text{- case 1 in cluster } j$$

$$y_{2j} = X_{2j}\beta + X_{2rj}u_j + \epsilon_{2j} \quad \text{- case 2 in cluster } j$$

$$\text{cov}(y_{1j}, y_{2j}) = \text{cov}(X_{1j}\beta + X_{1rj}u_j + \epsilon_{1j}, X_{2j}\beta + X_{2rj}u_j + \epsilon_{2j})$$

$$= \underbrace{\text{cov}(X_{1rj}u_j, X_{2rj}u_j)}_{\text{covariance between two cases in the same cluster is non zero}} + \text{cov}(\epsilon_{1j}, \epsilon_{2j}) + \text{cov}(\epsilon_{1j}, X_{2rj}u_j) + \text{cov}(\epsilon_{2j}, X_{1rj}u_j)$$

covariance between two cases in the same cluster is non zero  
bc they share the r.v.  $u_j$

At the cluster level

$$y_j = X_{rj}\beta_{pj} + X_{fj}\beta_k + \epsilon_j$$

$$V(y_j) = V(X_{rj}\beta_{pj} + X_{fj}\beta_k + \epsilon_j)$$

$$= V(X_{rj}\beta_p + X_{rj}u_j + X_{fj}\beta_k + \epsilon_j) \quad \text{where } u_j \text{ is a } px1 \text{ matrix} \\ \text{denoting } p \text{ random variables}$$

$$= V(X_{rj}u_j + \epsilon_j)$$

$$= X_{rj}V(u_j)X_{rj}' + V(\epsilon_j)$$

$$\underbrace{X_{rj}D X_{rj}'}_{V_j} + R_j$$

marginal variance-covariance matrix.

Multilevel model implies a marginal model.

$$V(y_j) = X_{rj} D X_{rj}' + R_j \quad \text{where } D = V(u_j)$$

the covariance matrix of random effects

As a covariance matrix, D needs to be positive definite

Consider the empty model:

$$y_{ij} = \beta_0 + \epsilon_{ij} \quad \epsilon_{ij} \sim N(0, \sigma^2)$$

$$\beta_{0j} = \beta_0 + u_{0j} \quad u_{0j} \sim N(0, \sigma_u^2)$$

$$\text{cov}(y_{ij}, y_{sj}) = \text{cov}(\beta_0 + u_{0j} + \epsilon_{ij}, \beta_0 + u_{0j} + \epsilon_{sj}) = \text{cov}(u_{0j}) = \sigma_u^2 \geq 0$$

$$\text{corr}(y_{ij}, y_{sj}) = \frac{\text{cov}(y_{ij}, y_{sj})}{\sqrt{V(y_{ij})} \sqrt{V(y_{sj})}} = \frac{\sigma_u^2}{V(y_{ij})} = \frac{\sigma_u^2}{\sigma_u^2 + \sigma^2} \geq 0$$

↑

ICC : Intra-class correlation.

must larger than 0 bc the correlation is a function of variance

The off-diagonal elements of  $V(y_j)$  matrix are non-negative bc these covariances are a function of the variance.

↓

$V(y_j)$  is unconstrained

↓

multilevel models are less flexible than marginal models

if we directly specify V matrix and do not factor out the  $u_{0j}$  part, we are not subjected to this constraint.

It is possible to specify negative ICC.

But in the context of marginal model, we don't have random effect.

We only have correlated variance structure

Pros of multilevel model. D matrix is of interest.

E.g., variance at the higher level.

Pros of marginal model. Relax the constraint on the positive intra-class correlation.

• Estimation uses two-step maximum likelihood:

1. Estimate  $V$  matrix first

$$2. \text{ with the estimated } V \quad \hat{\beta} = (X'VX)^{-1} X' \hat{V} Y$$

• Variance can be computed using the standard GLS results

$$V(\hat{\beta}) = (X'V^{-1}X)^{-1}$$

$\Rightarrow$  Variance of the coefficients is a direct function of  $V$ .

If  $D$  and  $R$  are specified badly or if  $V$  is specified badly,  
then SE of  $\hat{\beta}$  would be bad.

• The estimations use maximum likelihood.

Nested models can be compared with likelihood ratio test.

$$\Delta -2 \log\text{-likelihood} \sim \chi^2_q \quad \text{where } q \text{ is the different # of parameters}$$

Example:

M1: Random intercept model.

M2: Random intercept + one random slope model

} difference is  $\sigma^2$  and  $\sigma_{01}$

$$df = 2 ?$$

Not exactly,  $\sigma^2$  is not freely  
estimated, has to be  $> 0$

$\therefore$  more like 1.5 df.

$$\therefore 0.5(1 - pchisq(test statistic, df=1)) + 0.5(1 - pchisq(test statistic, df=2))$$

$\Rightarrow$  A mixture p value

- Predicted random effect

Empirically Best Linear Unbiased Predictors (EBLUP) of random effects

The variance and mean of random effects are estimated.

⇒ We can compute their predicted values based on the data.

$$\hat{u}_j = E(u_j | Y_j = y_j) = D Z_j' V_j^{-1} (y_j - X_j \beta)$$

For a demonstration, consider

$$Y_{ij} = \beta_0 + X_{ij} \beta_i + \epsilon_{ij}$$

$$\beta_{0j} = \beta_0 + u_{0j} \quad \begin{bmatrix} u_{0j} \\ u_{ij} \end{bmatrix} \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_0^2 & \sigma_{01} \\ \sigma_{10} & \sigma_1^2 \end{bmatrix} \right)$$

$$\beta_{ij} = \beta_i + u_{ij}$$

$$z_j = X_{rj} = \underbrace{\begin{bmatrix} 1 & x_{1j} \\ 1 & x_{2j} \\ \vdots & \vdots \\ 1 & x_{nj} \end{bmatrix}}_{\geq \text{random variables}} \quad r_j \quad R_j = \begin{bmatrix} \sigma^2 & & & \\ & \ddots & & \\ 0 & & \ddots & \\ & & & \sigma^2 \end{bmatrix}$$

$$\hat{u}_j = \begin{bmatrix} \hat{u}_{0j} \\ \hat{u}_{ij} \end{bmatrix} = \underbrace{\begin{bmatrix} D & Z_j' V_j^{-1} (y_j - X_j \beta) \\ 2 \times 2 & 2 \times n_j \quad n_j \times n_j \quad n_j \times 1 \end{bmatrix}}_{2 \times 1}$$

different clusters' residuals are different  
 $z_j$  are different  
 $v_j$  are different bc  $z_j$  is included  
 ↓  
 get different estimates

- Diagnostics

Residuals of the model  $\hat{\epsilon}_{ij}$  show a normal distribution?

EBLUP computed random effects show a normal distribution?

Conceptually understanding multilevel model as partial pooling:

50 states. intercept only model.

Option 1: Ignore states

$$y_{ij} = \bar{y} + \epsilon_{ij} \quad \epsilon_{ij} \sim N(0, \sigma^2) \quad \begin{matrix} i \text{ is individuals} \\ j \text{ is states.} \end{matrix}$$

Option 2: Separate model for the 50 states

$$y_{i1} = \bar{y}_1 + \epsilon_{i1} \quad \epsilon_{i1} \sim N(0, \sigma_1^2)$$

$$y_{i2} = \bar{y}_2 + \epsilon_{i2} \quad \epsilon_{i2} \sim N(0, \sigma_2^2)$$

:

$$y_{i50} = \bar{y}_{50} + \epsilon_{i50} \quad \epsilon_{i50} \sim N(0, \sigma_{50}^2)$$

Option 3: multilevel

$$y_{ij} = \bar{y}_j + \epsilon_{ij} \quad \epsilon_{ij} \sim N(0, \sigma_e^2)$$

$$\bar{y}_j = \bar{y} + u_j \quad u_j \sim N(0, \sigma_u^2)$$

with EBLUP,  $\hat{\bar{y}}_j$  can be computed for each state.

This  $\hat{\bar{y}}_j$  is a weighted sum of  $\bar{y}$  in option 1 and  $\bar{y}_j$  in option 2.

## Incorporating weights in multilevel models

Model-based approach and including weights as covariates doesn't always make sense.  
 (assuming that weights carry some information relevant to  $\gamma$  and use them as covariates in order to control for it)

Hybrid approach. need weights at both the levels.

- Higher level selection probability  $\rightarrow$  inverse.
- conditioning on the higher level, lower level selection probability  $\rightarrow$  inverse

$$\text{e.g., } P(\text{student} | \text{school}) = \frac{P(\text{student, school})}{P(\text{school})} \begin{matrix} \leftarrow \text{Survey report} \\ \text{individual selection weight} \end{matrix}$$

weights at the lower level need to be scaled.

sum to effective sample size or

sum to sample size.

longitudinal time points nested within individuals

if we follow the multilevel approach

$$y_{ti} = \beta_{0i} + \beta_{1i} x_{ti} + \epsilon_{ti} \quad \epsilon_{ti} \sim N(0, \sigma^2) \Rightarrow R_i = \begin{bmatrix} \sigma^2 & & & \\ 0 & \sigma^2 & & 0 \\ 0 & 0 & \ddots & \\ & & & \sigma^2 \end{bmatrix}$$

$$\beta_{0i} = \beta_0 + u_{0i} \quad [u_{0i}] \sim N([0], [\sigma_{00}^2 \quad \sigma_{01}^2])$$

$$\beta_{1i} = \beta_1 + u_{1i} \quad [u_{1i}] \sim N([0], [\sigma_{10}^2 \quad \sigma_{11}^2])$$

options:

1. First-order autoregressive

- correlations weaken as a power function of time

$$\begin{array}{cccc} T_1 & T_2 & T_3 & T_4 \\ T_1 & \sigma^2 & \sigma_M & \sigma_M \sigma_L \\ T_2 & \sigma_M & \sigma^2 & \\ T_3 & \sigma_M & & \sigma^2 \\ T_4 & \sigma_L & & \sigma^2 \end{array}$$

logitudinal analysis typically allows different error structure than this.

2. simple diagonal (constant variance, zero covariance)
3. compound symmetry (constant variance, constant covariance)
4. heterogeneous variance (variance can also change as a function of time)

## Marginal models

the correlated variance structure is something to be controlled, but not of research interest.

researchers are not interested in the variance of the slope, but only the fixed effect applicable across all subjects.

$$y_i \sim N_{ni}(X_i \beta, V_i) \quad \text{where } \beta \text{ is a vector of fixed effect}$$

$V_i$  is a  $n_i \times n_i$  matrix

Alternative method for fitting this kind of model :

↳ marginal linear model (MLM)  
 ↳ generalized estimating equations (GEE)

## Marginal linear model (MLM)

Same as multilevel model. a two-step approach

First, estimate  $V_i$  matrix

Then, GLS to compute  $\hat{\beta} = (X' V X)^{-1} X' V Y \downarrow$

But MLM cannot deal with non-normal outcome variables

GEE is a more flexible alternative.

solve the score function  $S(\beta) = \sum_{i=1}^n D_i^T V_i^{-1} (y_i - \mu_i) = 0$

where  $D_i$  is a  $n_i \times p$  matrix with the  $(i, j)$ th element being  $\frac{\partial \mu_i}{\partial \beta_j}$   
 $\mu_i$  is the expected value based on the specified model

The score function needs a  $V_j$  which we do not know.

$\therefore$  we work with a "working correlation matrix"

The estimate of  $\beta$  remains robust even with a bad choice of  $V$   
but standard error of the estimate would be influenced.

Options for the working correlation matrix:

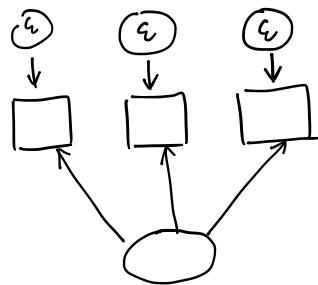
- independence
- exchangeable (constant correlation)
- first-order auto-regressive. decaying correlation over time
- unstructured.

with  $\mu_i$  being the expected value, the score function can be applied to situations where outcome variables are not normal.

structural equation modeling  
 measurement model  
 structural.

variables

observed variables  
 latent variables  
 endogenous variables  
 exogenous variables.



SEM examines plausibility of the hypothesized model.

If a model fits well, it just means that it's plausible.

But there may well be equivalent or better models

latent variables are unobserved.

∴ we need to fix its scale and location. For example,

- 1) fix one loading to be 1 and intercept to be 0
- 2) fix latent variable's variance directly and specify its mean.

This model would imply a set of variance and covariance for the observed variables. We can compare the implied variance covariance matrix and the sample variance-covariance matrix. If close, the model fits well.

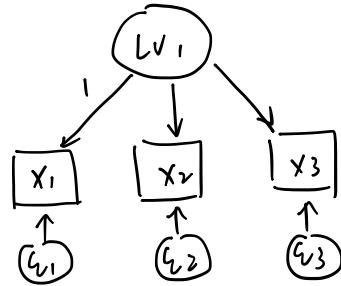
In fact, the maximum likelihood seeks to minimize the difference between the actual and the implied covariance matrix.

(if the model is saturated, there is no difference between the implied and the sample covariance matrix  
 → there is only one saturated model and there is a closed-form solution to the saturated model)

likelihood ratio test can be used to compare nested models  
 $\text{RMSEA} \leq 0.05 \checkmark \Rightarrow$  closeness between the implied and the sample covariance.

CFI: compare the fitted model to a null model that assume no relationship.

## Example 1



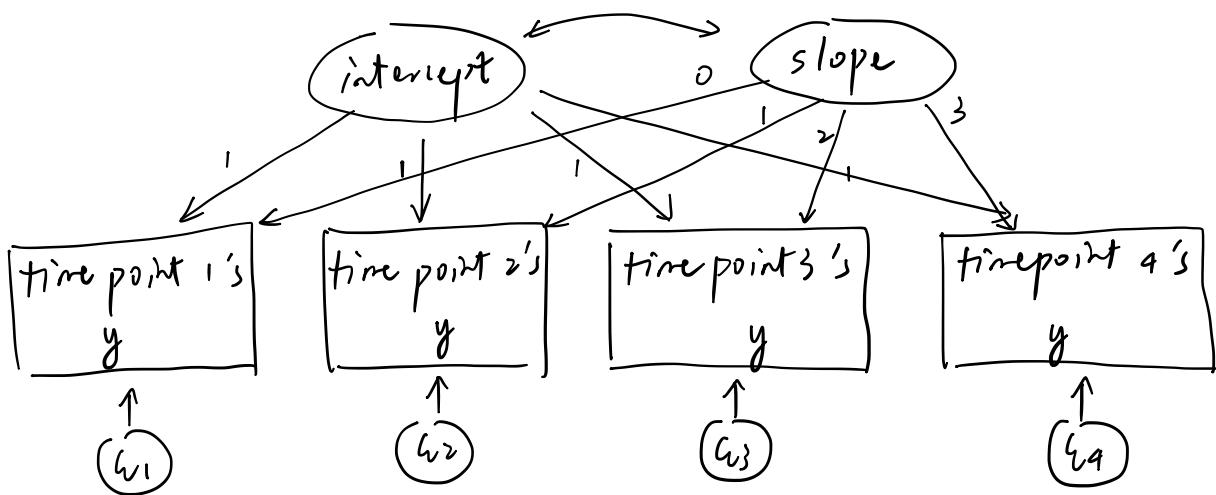
$$\begin{aligned}X_1 &= v_1 + \lambda_1 LV_1 + \epsilon_1 \\X_2 &= v_2 + \lambda_2 LV_1 + \epsilon_2 \\X_3 &= v_3 + \lambda_3 LV_1 + \epsilon_3\end{aligned}$$

$$\begin{aligned}LV_1 &\sim N(0, \sigma^2) \\ \begin{cases} \epsilon_1 \sim N(0, \sigma_1^2) \\ \epsilon_2 \sim N(0, \sigma_2^2) \\ \epsilon_3 \sim N(0, \sigma_3^2) \end{cases} \quad \hat{E}(x) = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}\end{aligned}$$

$$\hat{V}(x) = \begin{bmatrix} V(LV_1) + V(\epsilon_1) & & \\ \lambda_2 V(LV_1) & \lambda_2^2 V(LV_1) + V(\epsilon_2) & \\ \lambda_3 V(LV_1) & \lambda_2 \lambda_3 V(LV_1) & \lambda_3^2 V(LV_1) + V(\epsilon_3) \end{bmatrix}$$

## Example 3 latent growth curve model.

wide data. the information on time points are incorporated through the loadings that we specify



For individual  $i$

$$\text{time point 1's } y_i = \text{intercept}_i + 0 \cdot \text{slope}_i + \epsilon_1$$

$$\text{time point 2's } y_i = \text{intercept}_i + 1 \cdot \text{slope}_i + \epsilon_2$$

fix intercepts to 0

allows different people to have different starting points and different growth rates.

multi group analysis.

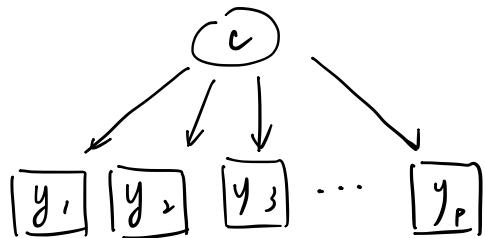
A different way to test / specify interactions.

measurement invariance

starting from free model, constrain the parameters to be the same across groups.

If the constraints do not significantly reduce model fits the invariance can be established.

latent class analysis.



joint probability of p items [ $\leftarrow$  this is the model we specify]

$$\begin{aligned} P(y_1, y_2, \dots, y_p) &= \sum_{k=1}^K P(c=k) P(y_1, y_2, \dots, y_p | c=k) \quad \text{Assuming conditional independence.} \\ &= \sum_{k=1}^K P(c=k) P(y_1 | c=k) P(y_2 | c=k) \dots P(y_p | c=k) \end{aligned}$$

posterior probability

$$P(c=k | y_1, y_2, \dots, y_p) = \frac{P(y_1, y_2, \dots, y_p | c=k) P(c=k)}{P(y_1, y_2, \dots, y_p)}$$

How many classes in C is what we need to figure out.

Likelihood ratio test can compare model with k and k-1 classes.

(doesn't have a  $\chi^2$  distribution.

use bootstrap to get an empirical distribution of LRT statistics)

At least 3 binary indicators are needed for a 2-class model.

with 2 binary indicators : 4 probabilities

$P(X_1=0 \mid X_2=0)$
$P(X_1=0 \mid X_2=1)$
$P(X_1=1 \mid X_2=0)$
$P(X_1=1 \mid X_2=1)$

what we are estimating :

$$P(X_1=1 \mid c=1)$$

/

$$P(X_2=1 \mid c=1)$$

5 parameters. — not identified.

$$P(X_1=1 \mid c=0)$$

$$P(X_2=1 \mid c=0)$$

$$P(c=1)$$

latent class analysis + distal outcomes

different ways to account for the uncertainty of class assignment.

- Mplus/Mixture has a multi-step approach  
that takes this uncertainty into account.

- Can make multiple assignments based on the estimated  
probabilities and use the multiple imputation approach to  
incorporate the uncertainty.