

Identifiability or

"When is your model too big for your data"

Mixed or longitudinal model for i th cluster

$$\tilde{y}_i = X_i \gamma + Z_i \tilde{u}_i + \tilde{\epsilon}_i$$

$$\tilde{u}_i \sim N(\tilde{0}, G) \text{ indep. of } \tilde{\epsilon}_i \sim N_{\tilde{n}_i}(\tilde{0}, R_i)$$

$$i = 1, \dots, K$$

For the fixed part of the model, i.e. X_i 's, the question is similar to 3330.

You need to consider the entire X matrix

$$\tilde{y} = X\gamma + Z\tilde{u} + \tilde{\epsilon}$$

$$\text{where } X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_K \end{bmatrix}, \quad Z = \begin{bmatrix} Z_1 & 0 & & \\ 0 & Z_2 & & \\ & & \ddots & 0 \\ 0 & & & Z_K \end{bmatrix},$$

$$\tilde{u} = \begin{bmatrix} \tilde{u}_1 \\ \vdots \\ \tilde{u}_K \end{bmatrix}, \quad \tilde{\epsilon} = \begin{bmatrix} \tilde{\epsilon}_1 \\ \vdots \\ \tilde{\epsilon}_K \end{bmatrix}$$

- Use: `vif(fit)` in the 'car' package to get collinearity diagnostics.
- DO NOT MECHANICALLY APPLY A SIMPLISTIC RULE TO DECIDE WHAT TO DO!!

- UNDERSTAND THE REASON FOR COLLINEARITY BEFORE ACTING.

e.g. Important confounders may be highly collinear with a causal variable. If you drop the confounder you will have a biased estimate of the causal effect.

- NEVER FOLLOW RECIPES WITHOUT UNDERSTANDING THE CONSEQUENCES OF YOUR CHOICES.

REMINDER:

PRIME DIRECTIVES:

1) KNOW THE PURPOSE OF YOUR ANALYSIS:

- versus
- a) CAUSAL / EXPLANATORY
 - b) FINDING A PREDICTIVE ALGORITHM
 - c) DESCRIPTIVE (often an excuse for ignoring the real purpose)
- } what difference does it make?

2) THE NATURE OF YOUR DATA

- NOT JUST WHAT IT LOOKS LIKE BUT HOW WAS IT OBTAINED?

- WAS THERE:

- RANDOM ASSIGNMENT?

- RANDOM SELECTION?

} what's the difference and why does it matter?

Example: - Longitudinal study with 2 time points
 $T = -1, 1$

Suppose we fit:

$\text{lme}(Y \sim 1 + X, \text{data}, \text{random} = \sim 1 + X | \text{id})$

$$\tilde{y}_i = X_i \gamma + Z_i \tilde{u}_i + \tilde{\varepsilon}_i$$

$$\begin{bmatrix} y_{i1} \\ y_{i2} \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 1 & +1 \end{bmatrix} \begin{bmatrix} \gamma_0 \\ \gamma_1 \end{bmatrix} + \begin{bmatrix} 1 & -1 \\ 1 & +1 \end{bmatrix} \begin{bmatrix} u_{i0} \\ u_{i1} \end{bmatrix} + \begin{bmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \end{bmatrix}$$

$$\begin{pmatrix} u_{i0} \\ u_{i1} \end{pmatrix} \sim N\left(\mathbf{0}, \begin{bmatrix} g_{00} & g_{01} \\ g_{10} & g_{11} \end{bmatrix}\right)$$

$$\begin{pmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \end{pmatrix} \sim N\left(\mathbf{0}, \sigma^2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

The variance model has 4 independent parameters

$$g_{00}, g_{01} = g_{10}, g_{11}, \sigma^2$$

They are not directly observable and can only be estimated through the observable

$$\begin{aligned} \text{Var}(\tilde{y}_i) &= \text{Var}(X_i \gamma + Z_i \tilde{u}_i + \tilde{\varepsilon}_i) \\ &= Z_i G Z_i' + \sigma^2 I \end{aligned}$$

$$\begin{aligned} \begin{pmatrix} v_{11} & v_{12} \\ v_{21} & v_{22} \end{pmatrix} &= \begin{bmatrix} 1 & -1 \\ 1 & +1 \end{bmatrix} \begin{bmatrix} g_{00} & g_{01} \\ g_{10} & g_{11} \end{bmatrix} \begin{bmatrix} 1 & 1 \\ -1 & +1 \end{bmatrix} + \sigma^2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \\ &= \begin{bmatrix} g_{00} - 2g_{01} + g_{11} + \sigma^2 & g_{00} - g_{11} \\ g_{00} - g_{11} & g_{00} + 2g_{01} + g_{11} + \sigma^2 \end{bmatrix} \end{aligned}$$

Problem?

(3 observable functions of variance) = f(4 independent parameters)

So: model is overparametrized and there is no hope of solving for unique optimal values of the G and R parameters.

In practice:

This might not matter if your goal is to estimate $\underline{\gamma}$ since $\text{Var}(\hat{\underline{\gamma}})$ depends on V

$$\underline{y} = X\underline{\gamma} + Z\underline{u} + \underline{\xi}$$

$$y = X\underline{\gamma} + \underline{\delta}$$

$$\text{Var}(\underline{\delta}) = ZGZ^T + R = V$$

GLS est $(X'V^{-1}X)(X'V^{-1}y)$