

a) One example of using “replicates” or “reproducibility study” observations to correct for structural measurement error

Suppose we have repeated measures (W_{i1}, W_{i2}, \dots) of an error-prone surrogate for a true exposure (X_i). These repeated measures are taken on all (or at least a reasonable subset) of the study participants

Assume the following simple linear regression “TDM”:

$$Y_i = \beta_0 + \beta_1 X_i + e_i$$

true exposure →

Also, assume the W 's relate to the X 's in the following way:

$$W_{ij} = X_i + \varepsilon_{ij}$$

MEM

where the errors ε_{ij} each have mean 0 and follow a compound-symmetric covariance structure

(PROCEED WITH HANDWRITTEN NOTES FROM HERE)

equivalently, consider a simple mixed linear model

$$W_{ij} = \mu + b_i + \varepsilon_{ij} \sim \text{iid } N(0, \sigma_e^2)$$

b_i 's and ε_{ij} 's mutually independent

$\overset{S}{\text{iid}} N(0, \sigma_b^2)$

$$W_{ij} = \mu + b_i + \epsilon_{ij}$$

SAS proc mixed
R

Let's start w/ balanced data ($j=1, \dots, n$)

Define "X_i" = $\mu_i = \mu + b_i$ "true" latent mean for subject i

$$\text{TDM: } Y_i = \beta_0 + \underbrace{(\beta_1)}_{\text{of interest}} \mu_i + \epsilon_i$$

Typical "surrogate" for μ_i = $\bar{w}_i = \frac{\sum_j w_{ij}}{n}$

Using \bar{w}_i in the "naive" regressor incurs a bias in β_1

How to characterize this bias?

⇒ can show $E(\mu_i | \bar{w}_i) = \gamma \bar{w}_i + (1-\gamma)\mu$
"shrinkage" toward μ dictated by

$$\gamma = \frac{n\sigma_b^2}{n\sigma_b^2 + \sigma_w^2}$$

$$0 < \gamma < 1$$

Let's try regression calibration (RC):

"Naive" OLS of Y vs. \bar{w}_i ⇒ what happens?

$$\begin{aligned} E(Y_i | \bar{w}_i) &= \beta_0 + \beta_1 E(\mu_i | \bar{w}_i) \\ &= \beta_0 + \beta_1 [\gamma \bar{w}_i + (1-\gamma)\mu] \end{aligned}$$

$$= \beta_0^* + \beta_1^* \bar{w}_i \quad \text{"link preserved"}$$

$$\downarrow$$

$$\beta_0 + \beta_1 (1-\gamma) \mu$$

$$\hat{\beta}_1^* \rightarrow \beta_1 \gamma$$

"attenuation"

$$\Rightarrow \hat{\beta}_1 = \frac{\hat{\beta}_1^*}{\hat{\gamma}} \rightarrow \text{plug in ML or REML estimates of } \sigma_b^2 + \sigma_e^2$$

Brunekreet, Noy, Clausen, 1987 (AJE)
design considerations

② What if we had unbalanced data $(j=1, \dots, n_j)$
 $i=1, \dots, k$
 \downarrow
subjects

\Rightarrow Now, $\text{var}(\mu_i | \bar{w}_i)$ is no longer constant across subjects.

\Rightarrow No "clean" corrected closed-form $\hat{\beta}_1$?

options

1) Quasi-likelihood \Rightarrow estimating eqns. weighted by $\text{var}^{-1}(\mu_i | \bar{w}_i)$

Liang + Liu (1991); Lyles + Kupper (1997)
 $(\mu, \sigma_b^2, \sigma_w^2) \rightarrow$ secondary params (from MEM)

2) ML: $\mathcal{L}(\theta, \psi; \mathbf{y}, \mathbf{W})$

\downarrow
primary params $(\beta_0, \beta_1, \sigma_b^2)$ TPM
~~($\mu, \sigma_b^2, \sigma_w^2$)~~

$$\begin{aligned}
 &= \prod_{i=1}^K f(Y_i, \tilde{w}_i; \theta, \psi) = \prod_{i=1}^K \int_{-\infty}^{\infty} f(Y_i, \tilde{w}_i, b_i; \theta, \psi) db_i \\
 &\quad \downarrow \text{latent} \\
 &= \prod_{i=1}^K \int_{-\infty}^{\infty} f(Y_i, \tilde{w}_i | b_i) f(b_i) db_i \\
 &= \prod_{i=1}^K \underbrace{f(Y_i | b_i)}_{\text{TPM}} \underbrace{f(\tilde{w}_i | b_i)}_{\substack{\downarrow \\ \text{w}_{ij}'\text{'s are} \\ \text{independent given } b_i \\ \text{(MEM)}}} \underbrace{f(b_i)}_{\substack{\downarrow \\ \text{MEM}}} \\
 &\quad \downarrow \\
 &\text{* NLMIXED, SAS IML, R}
 \end{aligned}$$

Can we extend this to the "randomized regression" model:

$$\Rightarrow w_{ij} = (\alpha + a_i) + (\beta + b_i) t_{ij} + \epsilon_{ij}$$

\downarrow longitudinal \downarrow random int. deviation \downarrow random slope deviation \downarrow time

$i = 1, \dots, K$ subjects.
 $j = 1, \dots, n_i$ obs. per subject

$$\begin{pmatrix} a_i \\ b_i \\ \epsilon_{ij} \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_a^2 & \sigma_{a\beta} & 0 \\ \sigma_{a\beta} & \sigma_\beta^2 & 0 \\ 0 & 0 & \sigma_e^2 \end{pmatrix} \right)$$

"True" intercept: $\alpha + a_i$
 "True" slope: $\beta + b_i$ (latent)

There could be exposures of interest to relate to health outcomes.

$$TDM: Y_i = \beta_0 + \beta_1 \alpha_i + \beta_2 \beta_i + e_i$$

\downarrow (2+ai) \downarrow (β+bi)

If data are balanced, easy to correct

"naïve" regression

$$Y_i = \theta_0^* + \theta_1^* \hat{\alpha}_{i,OLS} + \theta_2^* \hat{\beta}_{i,OLS} + \epsilon_i^*$$

$\hat{\theta}_1^* \rightarrow \lambda_1 \theta_1$ $\hat{\theta}_2^* \rightarrow \lambda_2 \theta_2$

\downarrow OLS estimates

Lyles and McFarlane (2000)

→ These "clean" results require
not just balanced data on w's, but
↓ equally-timed and equally-spaced w's.

What if not?

$$\mathcal{L}(\theta, \psi; Y, W | t)$$

\downarrow steps ↑ times TDM

$$= \prod_{i=1}^K \int \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(y_i | a_i, b_i, t_i; \theta) \times f(w_i | a_i, b_i; \psi) \times f(a_i, b_i) da_i db_i$$

\downarrow MEM ↑ independent (MEM)

* NLMixed

Wahnenweller
et al., 2005

application
Stats in med.

Daowen Zhang
+ Marie Davidian
↓ (1994?)
theory

b) Misclassification of the outcome (Y) in logistic regression

- Slides to be presented from a 2012 CDC seminar, where the topic has to do with handling (potentially differential) misclassification of the outcome variable in logistic regression via the use of internal validation data. They summarize the contents of the following paper, later extended in various ways by former graduate Dr. Li Tang and others:

Lyles RH, Tang L, Superak HM, King CC, Celentano DD, Lo Y, Sobel JD. Validation data-based adjustments for outcome misclassification in logistic regression: An illustration. *Epidemiology* **22**, 589-598 (2011).