

BIOS 755: Missing Data in Longitudinal Studies II

Alexander McLain

March 27, 2023

Incomplete Data Models

- ▶ Distinguish between dependent variable $\mathbf{y} = \mathbf{y}^O$ for $R = 1$ and $\mathbf{y} = \mathbf{y}^M$ for $R = 0$. and independent variables (all observed): \mathbf{X} time, group,
- ▶ GEE assumes special case of MCAR

$$P(\mathbf{R}|\mathbf{y}, \mathbf{X}) = P(\mathbf{R}|\mathbf{X}) \text{ for all } \mathbf{y}$$

\Rightarrow conditional on covariates, R is independent of both \mathbf{y}^O and \mathbf{y}^M
(“covariate-dependent missingness”)

Incomplete Data Models

- ▶ Likelihood-based methods assume MAR

$$P(\mathbf{R}|\mathbf{y}, \mathbf{X}) = P(\mathbf{R}|\mathbf{X}, \mathbf{y}^O) \text{ for all } \mathbf{y}^M$$

⇒ conditional on covariates and observed values of the dependent variable, \mathbf{R} is independent of \mathbf{y}^M “ignorable non-response” (Laird, 1988)

- ▶ But if baseline covariates are missing, all their \mathbf{y} data are missing from our analyses (whether they were observed or not).
- ▶ In that situation, we'd need to assume that the probability they are missing is independent of the \mathbf{y} data.

Simulation Study

- Data from 5000 subjects were simulated according to:

$$y_{ij} = \beta_0 + \beta_1 T_j + \beta_2 G_i + \beta_3 (G_i \times T_j) + v_{0i} + v_{1i} T_j + \varepsilon_{ij}$$

$T_j = 0, 1, 2, 3, 4$ for five timepoints

$G_i =$ dummy-code (0 or 1) with half in each group

Regression coefficients:

$$\beta_0 = 25, \beta_1 = -1, \beta_2 = 0, \text{ and } \beta_3 = -1$$

\Rightarrow the population means are:

25, 24, 23, 22, and 21 for $Grp = 0$

25, 23, 21, 19, and 17 for $Grp = 1$

Variance parameters:

$$\sigma_{v_0}^2 = 4, \sigma_{v_1}^2 = .25, \sigma_{v_{01}} = -.1 \ (\rho = -.1), \sigma^2 = 4$$

Simulation Study

- The population variance-covariance matrix,

$$V(\mathbf{y}) = \mathbf{Z}\Sigma_v\mathbf{Z}' + \sigma^2\mathbf{I}$$

$$V(\mathbf{y}) = \begin{bmatrix} 8.00 & 3.90 & 3.80 & 3.70 & 3.60 \\ 3.90 & 8.05 & 4.20 & 4.35 & 4.50 \\ 3.80 & 4.20 & 8.60 & 5.00 & 5.40 \\ 3.70 & 4.35 & 5.00 & 9.65 & 6.30 \\ 3.60 & 4.50 & 5.40 & 6.30 & 11.20 \end{bmatrix}$$

Scenarios

- ▶ **Complete data:** no missing data.
- ▶ **50% random missing:** 50% missing data at every timepoint; completely random and unrelated to any variable.
- ▶ **Time related dropout:** dropout rates of 0%, 25%, 50%, 75%, 87.5% for the five timepoints. If a subject was missing at a timepoint, then they were also missing at all later timepoints; these rates indicate the percentage of the original sample that were missing at each of these timepoints.

Scenarios

- ▶ **Group by time related dropout:** dropout rates of
0%, 23%, 46%, 70% and 83% for $G=0$
0%, 27%, 55%, 81% and 91% for $G=1$
⇒ notice that these missing data scenarios are all MCAR (as long as analysis model includes time, G , and G by time)

Results

► MCAR Simulation Results Mixed effects Model estimates (standard errors)

	β_0	β_1	β_2	β_3	$\sigma_{v_0}^2$	$\sigma_{v_{01}}$	$\sigma_{v_1}^2$	σ^2
simulated value:	25	-1	0	-1	4	-.1	.25	4
complete data	24.969 (.050)	-.994 (.016)	-.001 (.071)	-.986 (.023)	3.918 (.129)	-.057 (.032)	.239 (.014)	3.991 (.046)
50% random missing	24.991 (.063)	-1.024 (.023)	-.087 (.089)	-.933 (.032)	3.811 (.193)	-.020 (.056)	.199 (.025)	4.070 (.083)
Time-related dropout	24.989 (.053)	-.968 (.028)	.019 (.075)	-1.021 (.040)	3.853 (.150)	-.062 (.060)	.229 (.032)	4.000 (.074)
Group by Time related dropout	24.991 (.053)	-.977 (.026)	.041 (.075)	-1.014 (.041)	3.872 (.150)	-.048 (.059)	.234 (.031)	3.994 (.073)

MAR and MNAR Scenarios

- ▶ **MAR(a):** if the value of the dependent variable was lower than 23, then the subject dropped out at the next timepoint (i.e., they were missing at the next and all subsequent timepoints).
- ▶ **MAR(b):** the MAR specification was different for the two groups. For $\text{Grp} = 1$, if the dependent variable was lower than 23, then the subject dropped out at the next timepoint, however for $\text{Grp} = 0$, if the dependent variable was greater than 25.5 then the subject dropped out at the next timepoint.
- ▶ **MNAR:** after the first timepoint, if the value of the dependent variable was lower than 21.5, then the subject was missing at that timepoint and all subsequent timepoints.

MAR and MNAR Simulation Results - Estimates (standard errors)

	β_0	β_1	β_2	β_3	$\sigma_{v_0}^2$	$\sigma_{v_{01}}$	$\sigma_{v_1}^2$	σ^2
simulated value:	25	-1	0	-1	4	-.1	.25	4
<u>MAR(a)</u>								
MRM	24.996 (.053)	-1.039 (.025)	-.010 (.075)	-.969 (.041)	3.981 (.158)	-.064 (.065)	.233 (.032)	3.873 (.078)
GEE1	25.281 (.058)	-1.164 (.037)	.019 (.097)	-1.001 (.085)				
<u>MAR(b)</u>								
MRM	24.999 (.053)	-1.003 (.022)	-.016 (.075)	-1.004 (.039)	4.050 (.154)	-.082 (.064)	.229 (.027)	3.812 (.073)
GEE1	24.635 (.055)	-.714 (.030)	.634 (.097)	-1.532 (.090)				
<u>MNAR</u>								
MRM	24.956 (.049)	-.233 (.020)	.027 (.070)	-.552 (.035)	3.856 (.131)	-.943 (.051)	.319 (.025)	3.020 (.053)
GEE1	25.051 (.049)	-.386 (.020)	.016 (.071)	-.583 (.034)				

Misspecification of Variance-Covariance

- Re-analysis of the simulated MAR-generated data, however using only a random-intercepts model:

$$y_{ij} = \beta_0 + \beta_1 \text{Time}_j + \beta_2 \text{Grp}_i + \beta_3(\text{Grp}_i \times \text{Time}_j) + v_{0i} + e_{ij}$$

- This is a misspecified model for the variance-covariance structure because the random slope term was omitted from the analysis.

Mis-specified MAR Simulation Results - Estimates (standard errors)

	β_0	β_1	β_2	β_3	$\sigma_{v_0}^2$	$\sigma_{v_{01}}$	$\sigma_{v_1}^2$	σ^2
simulated value:	25	-1	0	-1	4	-.1	.25	4
MRM with MAR(a)	24.938 (.053)	-.891 (.021)	-.009 (.075)	-.949 (.036)	3.722 (.136)			4.329 (.072)
MRM with MAR(b)	24.998 (.053)	-1.048 (.018)	-.069 (.076)	-.805 (.035)	3.880 (.138)			4.288 (.070)

Results

- ▶ These random-intercepts analyses yield biased results, in particular for the time-related parameters β_1 and β_3 .
- ▶ Performing any full-likelihood analysis, even with missing data following an MAR mechanism, does not guarantee that the correct results will be obtained (need to have the mean structure and variance-covariance structure of \mathbf{y} correctly modeled).

Testing MCAR

- ▶ If MCAR, then either GEE or Mixed Effects models are fine (provided that the covariate matrix \mathbf{X}_i includes predictors of missingness).
- ▶ If MAR, then GEE does not perform well, whereas Mixed Effects models analysis is acceptable (as long as the mean and variance-covariance structures are correctly modeled).
- ▶ It is useful to determine whether MCAR is acceptable or not.
- ▶ Distinction between MCAR and MAR is that missingness cannot depend on observed values of the dependent variable, \mathbf{y}_i^O , in the former, but can in the latter.
- ▶ Tests of MCAR are based on analyses involving \mathbf{y}_i^O .

Testing MCAR in 2-timepoint study

- ▶ Suppose all subjects have data at time 1, but some are missing at time 2
- ▶ Define $D_i = 0$ for subjects with data at both timepoints, $D_i = 1$ for subjects with data at first timepoint only
- ▶ Compare y_1 between these two groups ($D_i = 0$ vs. $D_i = 1$); for MCAR y_1 data should not differ between the groups.
- ▶ Logistic regression model

$$\log \left\{ \frac{P(D_i = 1)}{1 - P(D_i = 1)} \right\} = \alpha_0 + \alpha_1 y_{i1} + \alpha_2 \mathbf{x}_i + \alpha_3 (y_{i1} \times \mathbf{x}_i)$$

- ▶ α_1 and α_3 are vectors of regression coefficients for \mathbf{x}_i and their interactions with $y_{i1} \Rightarrow$ MCAR dictates that $\alpha_1 = \alpha_3 = 0$.

Testing MAR vs MNAR

- ▶ The only way to test the MAR assumption (i.e., that \mathbf{R} is independent of \mathbf{y}^M given \mathbf{y}^O and \mathbf{X}_i) is to have some measure on the missing data.
- ▶ This can be accomplished by a follow-up with phone calls to a group of the non-respondents.
- ▶ Mostly we will not have any data.
- ▶ The only quantitative solution in this case is to do the analysis with an MAR model and an MNAR model (i.e., selection, pattern-mixture, or shared parameter).
- ▶ The most common way to “test” the MAR assumption is to **use your scientific knowledge of the data and the field.**