

Some methods for handling missing values in outcome variables

Roderick J. Little



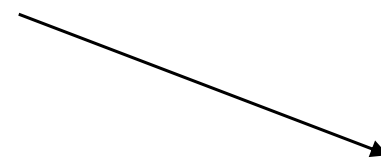
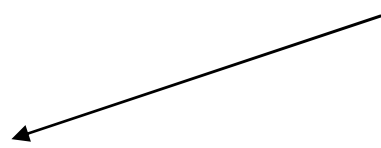
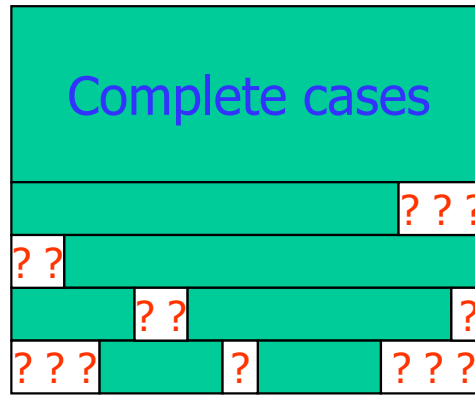
Outline

- Missing data principles
- Likelihood methods
 - ML, Bayes, Multiple Imputation (MI)
- Robust MAR methods
 - Predictive mean matching hot deck,
 - Penalized Spline of Propensity Prediction
- MNAR methods
 - Sensitivity analysis via pattern-mixture models
 - Offsets to chained equation MI
 - Proxy pattern-mixture analysis

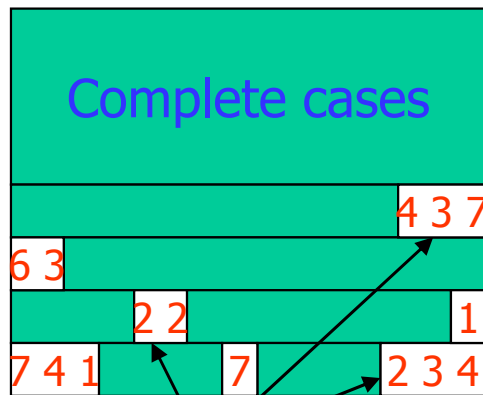
Properties of a good missing-data method

- Makes use of partial information on incomplete cases, for reduced bias, increased efficiency of estimates
 - Goal is better inference from observed data, not best estimates of the missing values
- Valid inferences under plausible model for mechanism and data (e.g. confidence intervals have nominal coverage)
- Propagates missing-data uncertainty
 - Particularly when fraction of missing information is large

General Strategies

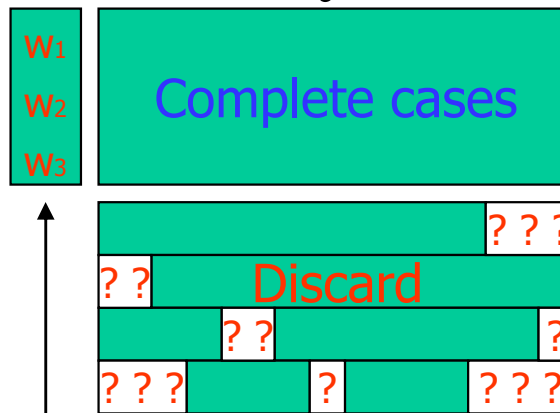


Imputation



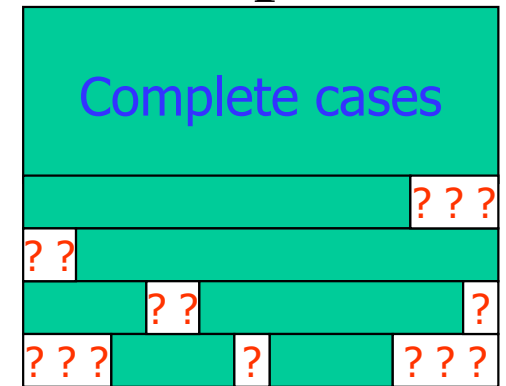
Imputations

Complete-Case Analysis



Weights

Analyze Incomplete



e.g. maximum likelihood

Missing-data mechanisms

Y = data matrix, if no data were missing

M = missing-data indicator matrix

(i,j) th element indicates whether (i,j) th element of Y is missing (1) or observed (0)

– Model mechanism via distribution of M given Y :

– Missing completely at random (MCAR):

$$p(M | Y) = p(M) \text{ for all } Y$$

– Missing at random (MAR):

$$p(M | Y) = p(M | Y_{\text{obs}}) \text{ for all } Y$$

– Missing not at random (MNAR) if missingness depends on missing (as well as perhaps on observed) components of Y (Rubin 1976, Little and Rubin 2002)

MAR for longitudinal dropout

MAR if dropout depends on values recorded prior to drop-out

MNAR if dropout depends on values that are missing (that is, after drop-out)

Censoring by end of study: plausibly MCAR

Designed missing data: generally MCAR or MAR

Unit, item nonresponse: plausibly MAR with good covariate info, otherwise often MNAR

Complete-Case Analysis

- Simple and may be good enough when information in complete cases is limited
 - Depends on context
- Loss of information in incomplete cases has two aspects:
 - Increased variance of estimates
 - Bias when complete cases differ systematically from incomplete cases -- often the case
- Weighting by inverse estimated response rate can reduce bias under MAR, but does not use covariate data efficiently
 - Common for unit nonresponse

Unweighted CC analysis

- CC analysis is inefficient if there is substantial information in incomplete cases
- Information in incomplete cases depends on pattern, estimand and mechanism
- E.g. Figure 1: incomplete cases have:
 - Substantial information for mean of Y , when X 's are predictive
 - No information for regression of Y on X , under MAR;
(but do contain information under NMAR)

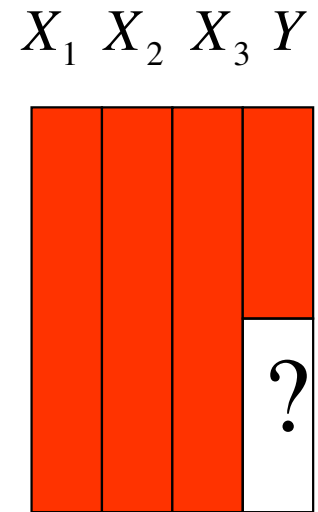


Figure 1

Multiple imputation (MI)

- Imputes *draws*, not means, from the predictive distribution of the missing values
- Creates $M > 1$ filled-in data sets with different values imputed
- MI combining rules yield valid inferences under well-specified models – impact of model misspecification increases with the fraction of missing information and deviation from MCAR
- propagate imputation uncertainty, and averaging of estimates over MI data sets avoids the efficiency loss from imputing draws
- Note that MI is valid (under model) even if predictive power of models is weak, since uncertainty is propagated
- MI can also be used for non-MAR models, particularly for *sensitivity analyses* – more later on this

Ex. 1 contd. Tacrine Dataset

IT Analysis, Continuing Dose MI Model: 80mg vs Placebo

<i>MI number</i>	<i>Treat.diff (s.e.)</i>	<i>p-value</i>	<i>95 %C.I.</i>
1	-3.486 (0.951)	0.0003	(-5.35,-1.62)
2	-3.682 (0.876)	0.0000	(-5.40,-1.97)
3	-3.142 (0.944)	0.0009	(-4.99,-1.29)
4	-4.889 (0.908)	0.0000	(-6.67,-3.11)
5	-4.633 (0.910)	0.0000	(-6.42,-2.85)
6	-4.146 (0.920)	0.0000	(-5.95,-2.34)
7	-5.239 (0.925)	0.0000	(-7.05,-3.43)
8	-4.463 (0.933)	0.0000	(-6.29,-2.63)
9	-4.511 (0.953)	0.0000	(-6.38,-2.64)
10	-3.497 (0.899)	0.0001	(-5.26,-1.73)
MI Inference	-4.169 (1.173)	0.0039	(-6.72,-1.62)

Advantages of MI

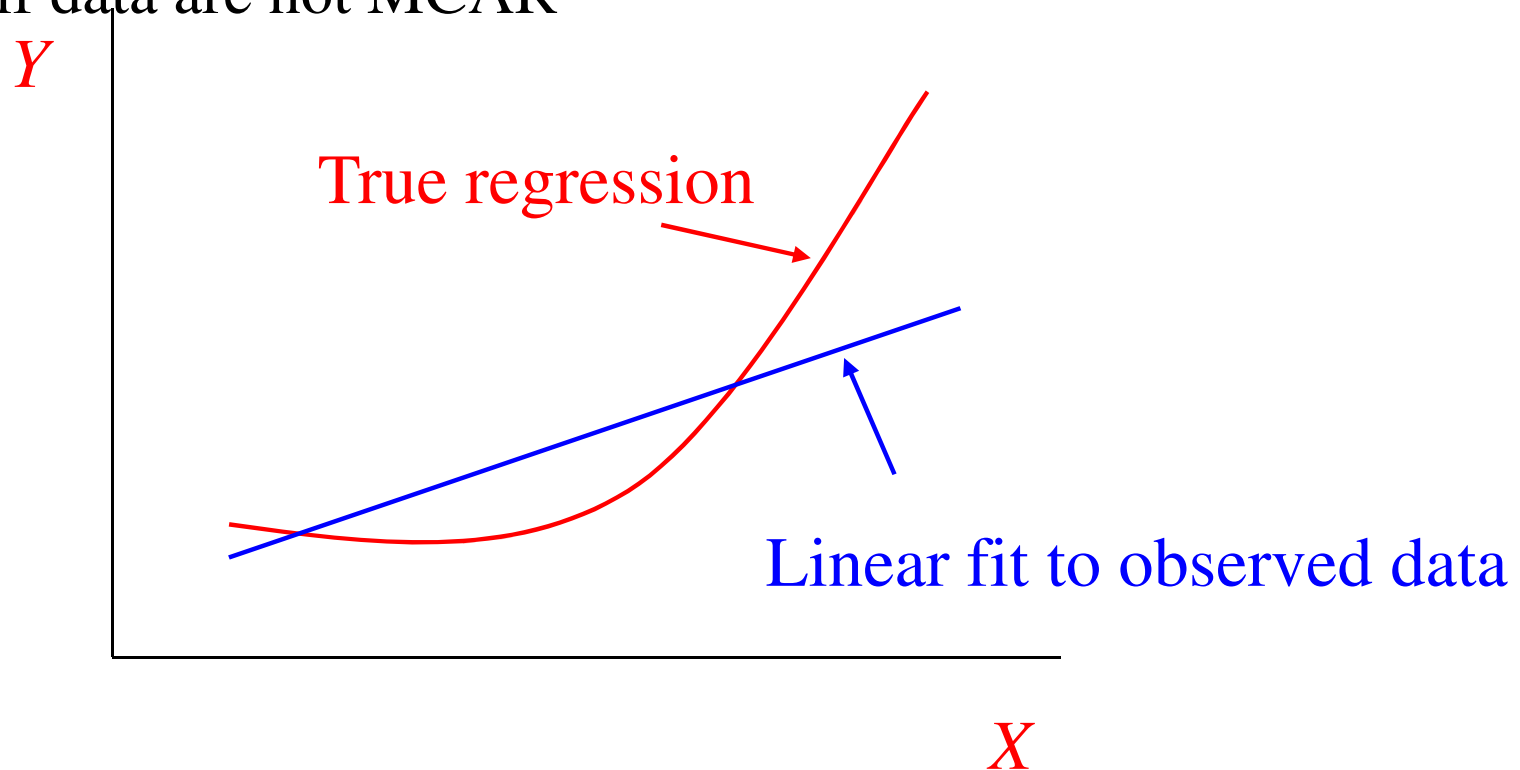
- Imputation model can differ from analysis model
 - By including variables not included in final analysis
 - Promotes consistency of treatment of missing data across multiple analyses
 - Assumptions in imputation model are confined to the imputations – hence with little missing data, simple methods suffice
- Public use data set users can be provided MI's, spared task of building imputation model
 - MI analysis of imputed data is easy, using complete-data methods (e.g. SAS PROC MIANALYZE)

Examples of MI

- Bayes for parametric models, e.g. multivariate normal, general location model (PROC MIXED)
- Sequential regression/chained equations MI (IVEware, MICE)
- Hot deck multiple imputation (more below)

Making MI's under MAR more robust

- Aim to reduce sensitivity of parametric MI's to model misspecification, particularly when data are not MCAR
- Hot deck methods like predictive mean matching
- Weaken regression assumptions of parametric MI's are potentially sensitive to model misspecification, particularly if data are not MCAR



Hot deck MI

For review of hot deck methods see Andridge and Little (2010)

Hot deck can create multiple imputations, as multiple draws from a donor set close to the recipient on some metric

A preferred metric: predictive mean matching: choose donor with small value of

$$\left(\hat{\mu}_{y \cdot x, \text{donor}} - \hat{\mu}_{y \cdot x, \text{recip}} \right)^T \hat{\Sigma}_{y \cdot x}^{-1} \left(\hat{\mu}_{y \cdot x, \text{donor}} - \hat{\mu}_{y \cdot x, \text{recip}} \right)$$

Extensions:

1. Longitudinal events histories with gaps (Wang et al, 2011)
2. Predictive moment matching (Wang & Little, in progress)

Penalized Spline of Propensity Prediction (PSPP)

- PSPP (Little & An 2004, Zhang & Little 2009, 2011).
- Regression imputation that is
 - Non-parametric (spline) on the propensity to respond
 - Parametric on other covariates
- Exploits the key property of the propensity score that conditional on the propensity score and assuming missing at random, missingness of Y does not depend on the covariates

PSPP method

Estimate: $Y^* = \text{logit}(\Pr(M=0|X_1, \dots, X_p))$

Impute using the regression model:

$$(Y | Y^*, X_1, \dots, X_p; \beta) \sim$$

$$N(s(Y^*) + g(Y^*, X_2, \dots, X_p; \beta), \sigma^2)$$

- Nonparametric part
- Need to be correctly specified
- We choose penalized spline

- Parametric part
- Misspecification does not lead to bias
- Increases precision
- X_1 excluded to prevent multicollinearity

Double Robustness Property

- The PSPP method yields a consistent estimator for the marginal mean of Y , if:

(a) the mean of Y given X is correctly specified,

or

(b1) the propensity is correctly specified, and

(b2) $E(Y | Y^*) = s(Y^*)$

Key idea: the parametric regression $g()$ on the other covariates does not have to be correctly specified

Missing Not at Random Models

- Difficult problem, since information to fit non-MAR is limited and highly dependent on assumptions
- Sensitivity analysis is preferred approach – this form of analysis is not appealing to consumers of statistics, who want clear answers
- Selection vs Pattern-Mixture models
 - Prefer pattern-mixture factorization since it is simpler to explain and implement
 - Offsets, Proxy Pattern-mixture analysis
- Missing covariates in regression
 - Subsample Ignorable Likelihood (talk 3)

A simple missing data pattern

$i = i$ th observation

$x_i =$ baseline covariates (incl constant, treatment)

$y_{0i} =$ baseline value of outcome

$y_{1i} =$ outcome at intermediate time point

$h_i = (x_i, y_{0i}, y_{1i})$ "History" for observation i

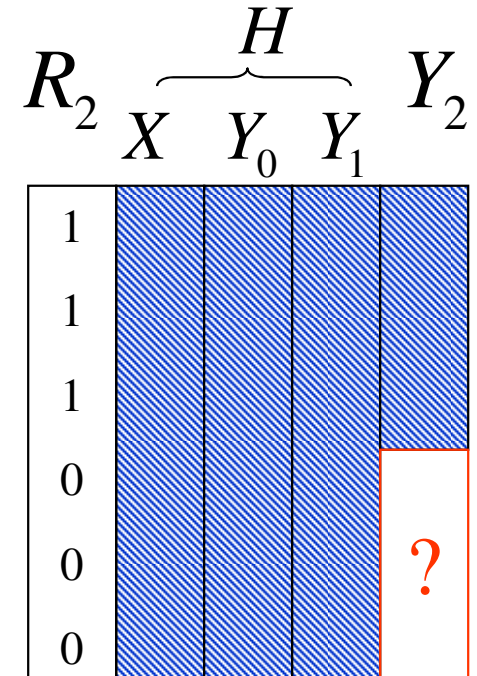
$y_{2i} =$ outcome at final time point

$r_{2i} =$ response indicator for y_{2i}

Target for inference: $E(y_{2i} - y_{0i} \mid z_i)$

$z_i =$ subset of $\{x_i, y_{0i}\}$

Missing data problem: missing values $\{y_{2i}\}$



A simple missing data pattern

$[y_{2i} | h_i, r_{2i} = 1]$: estimated from data

Complete-case analysis: drop $r_{2i} = 0$ cases

Inference then restricted to complete cases

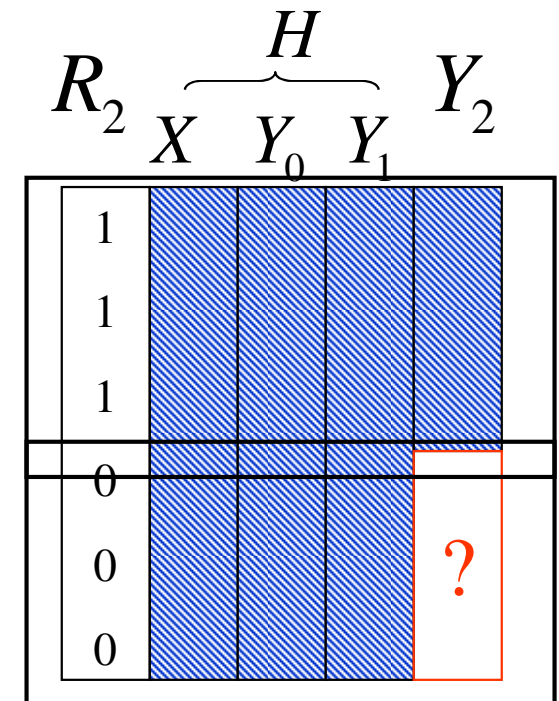
May be OK (e.g. $Y_2 = \text{QOL}$, $R_2 = \text{death}$)

Otherwise need to model or predict

nonrespondent values of y_{2i}

$[y_{2i} | h_i, r_{2i} = 0]$: no information in data

Need to make assumptions, i.e. model!



Missing at random assumption

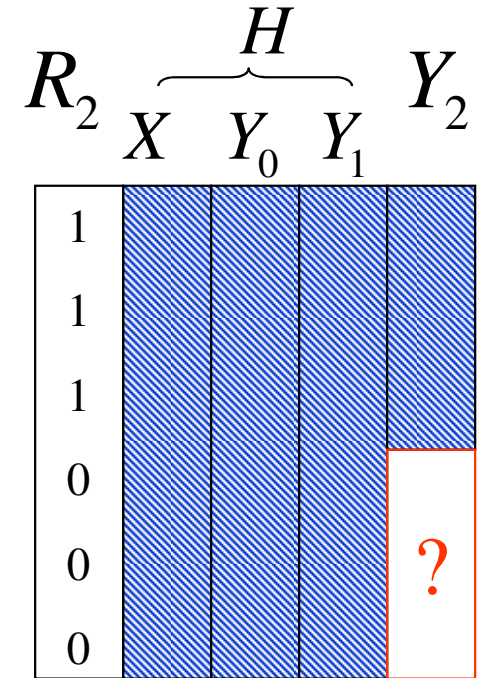
MAR

$$[r_{2i} | h_i, y_{2i}] = [r_{2i} | h_i]$$

or equivalently,

$$[y_{2i} | h_i, r_{2i} = 0] = [y_{2i} | h_i, r_{2i} = 1]$$

Plausibility depends quality of predictors



Missing not at random models

MNAR

$$[y_{2i} | h_i, r_{2i} = 0] \neq [y_{2i} | h_i, r_{2i} = 1]$$

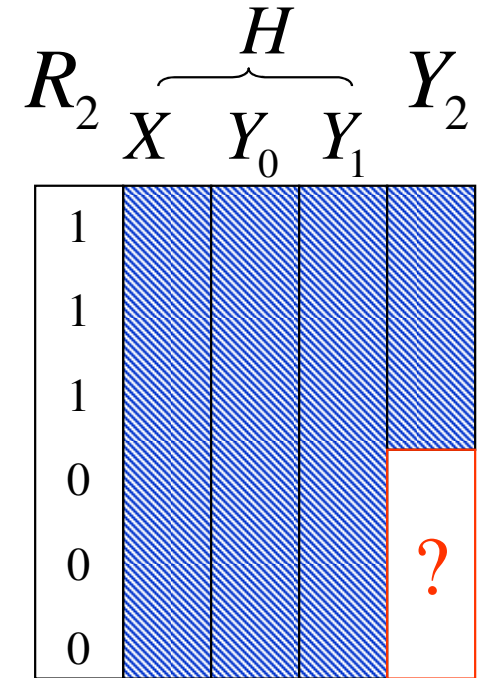
Infinite possibilities!

Two classes (Little & Rubin 2002) are:

Selection Models: $[y_{2i} | h_i] \times [r_{2i} | y_{2i}, h_i]$

Pattern-Mixture Models: $[y_{2i} | r_{2i}, h_i] \times [r_{2i} | h_i]$

I like Pattern-mixture models, since they are more straightforward and transparent



Heckman Selection Model

$$[y_{2i} | h_i] \sim G(\beta^T h_i, \sigma^2)$$

$$r_{2i} = 1 \text{ when } u_{2i} > 0, [u_{2i} | y_{2i}, h_i] \sim G(\alpha^T h_i + \lambda y_{2i}, 1)$$

$$\Rightarrow \Pr(r_{2i} = 1 | y_{2i}, h_i) = \Phi^{-1}(\alpha^T h_i + \lambda y_{2i})$$

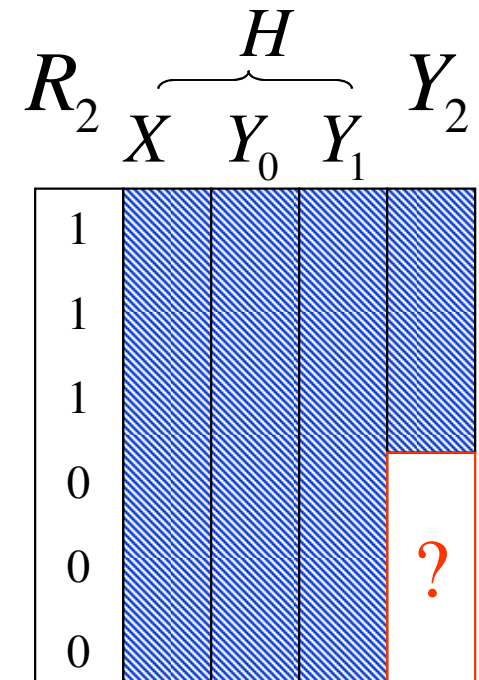
(Heckman 1976)

λ is weakly (practically not) identified without restrictions on β, α

I think attempting to estimate λ is a bad idea

Better to do sensitivity analysis for choices of λ

Pattern-mixture model easier to fit and interpret, since relevant predictive distributions $[y_{2i} | h_i, r_{2i} = 0 \text{ or } 1]$ are modeled directly ...



A simple pattern-mixture model

“*In special cases*, it may be possible to estimate the effect of nonrespondents under accepted models. *More often*, the investigator has to make subjective judgments about the effect of nonrespondents. Given this situation, it seems reasonable to try to formulate these subjective notions *so that they can be easily stated, communicated, and compared*” (Rubin 1977, emphasis added)

$$[y_{2i} \mid h_i, r_{2i} = k] \sim G(\beta^{(k)} h_i, \tau^{2(k)})$$

$$\beta^{(1)} = (\beta_0, \beta), \beta^{(0)} = (\beta_0 + \lambda\tau^{(1)}, \beta)$$

That is, intercept for nonrespondents is perturbed by an offset $\lambda\tau^{(1)}$

Sensitivity analysis, varying λ (Clearly no information about λ here)

Simpler (embarrassingly so?), easier to fit than Heckman model

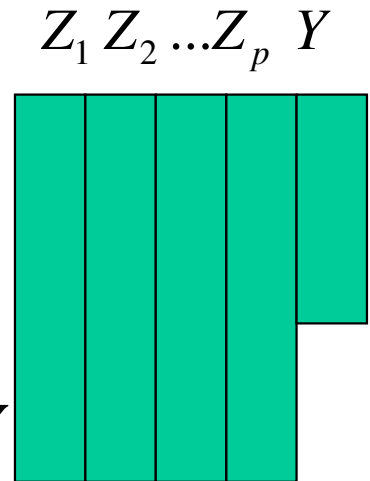
A simple pattern-mixture model

Giusti & Little (2011) extends this idea to a PM model for income nonresponse in a rotating panel survey:

- * Two mechanisms (rotation MCAR, income nonresponse NMAR)
- * Offset includes as a factor the residual sd, so smaller when good predictors are available
- * Complex problem, but PM model is easy to interpret and fit

Proxy pattern-mixture analysis (PPMA)

- Setting: univariate nonresponse
- Y = survey outcome
- Z = auxiliary covariate information
- Goal: nonresponse adjustment of mean of Y
 - (non-MAR as well as MAR)



Create X = single best proxy for Y based on $Z = (Z_1, \dots, Z_p)$

Compute by regression of Y on Z_1, \dots, Z_p using complete cases

$$\rho = \text{Corr}(X, Y) > 0$$

Call X a strong proxy if ρ is high, a weak proxy if ρ is low

Properties of estimators

- Key information about nonresponse bias for Y is:

$d = \bar{x}_R - \bar{x}$, measures deviation from MCAR (is there a problem?)

$\hat{\rho}$, measures strength of proxy information (can we fix the problem?)

- PPMA satisfies the following intuitive ranking of scenarios:

1. (Best): ρ high, $|d|$ low "strong evidence of no problem"

2.5 ρ high, $|d|$ high "evidence of a problem, but fixable"

2.5 ρ low, $|d|$ low "weak evidence of no problem"

4. (Worst) ρ low, $|d|$ high "evidence of problem, not fixable"

- PPMA yields least uncertainty for 1, most uncertainty for 4
- Specific choices of $g()$ are based on a pattern-mixture model

Pattern-mixture model

$$((X, Y) | M = m) \sim N_2 \left((\mu_x^{(m)}, \mu_y^{(m)}), \Sigma^{(m)} \right)$$

$$M \sim \text{Bernoulli}(\pi)$$

$$\Sigma^{(m)} = \begin{pmatrix} \sigma_{xx}^{(m)} & \rho^{(m)} \sqrt{\sigma_{xx}^{(m)} \sigma_{yy}^{(m)}} \\ \rho^{(m)} \sqrt{\sigma_{xx}^{(m)} \sigma_{yy}^{(m)}} & \sigma_{yy}^{(m)} \end{pmatrix}$$

$$\Pr(M = 1 | X, Y) = f(Y^*), Y^* = X + \lambda Y, f \text{ unspecified}, \lambda \geq 0$$

$\lambda = 0 \Rightarrow$ missingness depends on X (MAR);

$\lambda = 1 \Rightarrow$ missingness depends on $X + Y$

$\lambda = \infty \Rightarrow$ missingness depends only on Y

Two options: (A) Sensitivity analysis over range of λ

Or (B) specify a prior distribution for λ

Pattern-mixture model

Neat feature: do not need to specify form of f
(Unlike e.g. Heckman selection model)

Model is just identified by parameter restrictions:

$$[X, Y | Y^*, M = 0] = [X, Y | Y^*, M = 1]$$

In particular, ML estimate of mean of Y is

$$\hat{\mu}_y = \bar{y}_R + \frac{s_{xy}^{(0)} + \lambda s_{yy}^{(0)}}{s_{xx}^{(0)} + \lambda s_{xy}^{(0)}} (\bar{x} - \bar{x}_R) \text{ (Little 1994)}$$

Proxy pattern-mixture model

Transform $Z \rightarrow (X, V)$,

$X = Z^T \alpha =$ best predictor of Y , $V =$ other covariates

$$[Y, X, V, M, \alpha] = [Y, X | M, \alpha][M][\alpha][V | Y, X, M, \alpha]$$

$$((X, Y) | M = m) \sim N_2 \left((\mu_x^{(m)}, \mu_y^{(m)}), \Sigma^{(m)} \right)$$

$$M \sim \text{Bernoulli}(\pi)$$

$$\Sigma^{(m)} = \begin{pmatrix} \sigma_{xx}^{(m)} & \rho^{(m)} \sqrt{\sigma_{xx}^{(m)} \sigma_{yy}^{(m)}} \\ \rho^{(m)} \sqrt{\sigma_{xx}^{(m)} \sigma_{yy}^{(m)}} & \sigma_{yy}^{(m)} \end{pmatrix}$$

Unspecified

$$\Pr(M = 1 | X, Y) = f(Y^*), Y^* = X \sqrt{\sigma_{yy}^{(0)} / \sigma_{xx}^{(0)}} + \lambda Y, \lambda \geq 0$$

rescaling X aids interpretation of λ

PPMA ML estimate

ML estimate of mean of Y is

$$\hat{\mu}(\lambda) = \bar{y}_R + g(\hat{\rho}) \sqrt{\frac{s_{yy}^{(0)}}{s_{xx}^{(0)}}} (\bar{y}_{NR}^* - \bar{y}_R^*), \quad g(\hat{\rho}) = \left(\frac{\hat{\rho} + \lambda}{1 + \hat{\rho}\lambda} \right)$$

$\lambda \geq 0$ is a sensitivity parameter,

determined by assumed missing data mechanism

Propose sensitivity analysis with three values of λ :

$\lambda=0$, $g(\hat{\rho}) = \hat{\rho}$ (MAR, usual regression estimator)

$\lambda=1$, $g(\hat{\rho}) = 1$ (NMAR, carries over bias adjustment from proxy)

$\lambda=\infty$, $g(\hat{\rho}) = 1 / \hat{\rho}$ (NMAR, inverse regression estimator)

Note: $g(\hat{\rho})$ varies between $\hat{\rho}$ and $1 / \hat{\rho}$, reduced sensitivity as $\hat{\rho} \uparrow 0$

Estimation methods

- 1. Maximum Likelihood
 - Doesn't incorporate uncertainty in regression parameters used to create the proxy
 - Large-sample variances by Taylor series calculations
- 2. Bayesian, non-informative priors
 - Proxy recreated at each draw of regression parameters, so uncertainty is incorporated
 - Easy to implement, non-iterative
- 3. Multiple Imputation of missing Y 's
 - Allows complex design features to be incorporated in the within-imputation component of variance
 - Easy to implement

Simulations

- Assess confidence coverage and width of ML, Bayes, MI for

$$\rho = 0.2, 0.5, 0.8$$

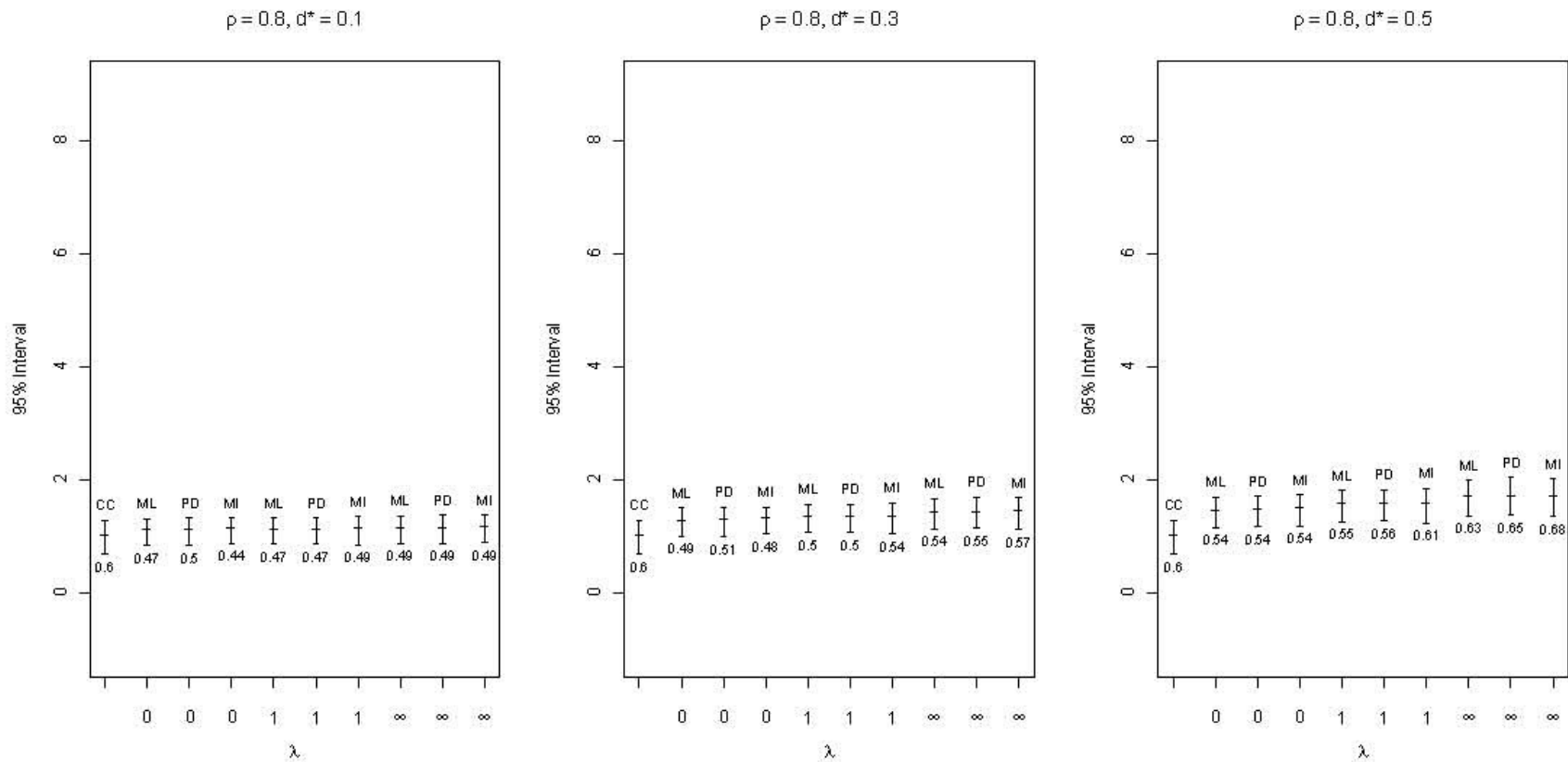
$$d^* = 0.1, 0.3, 0.5$$

$$n = 100, 400$$

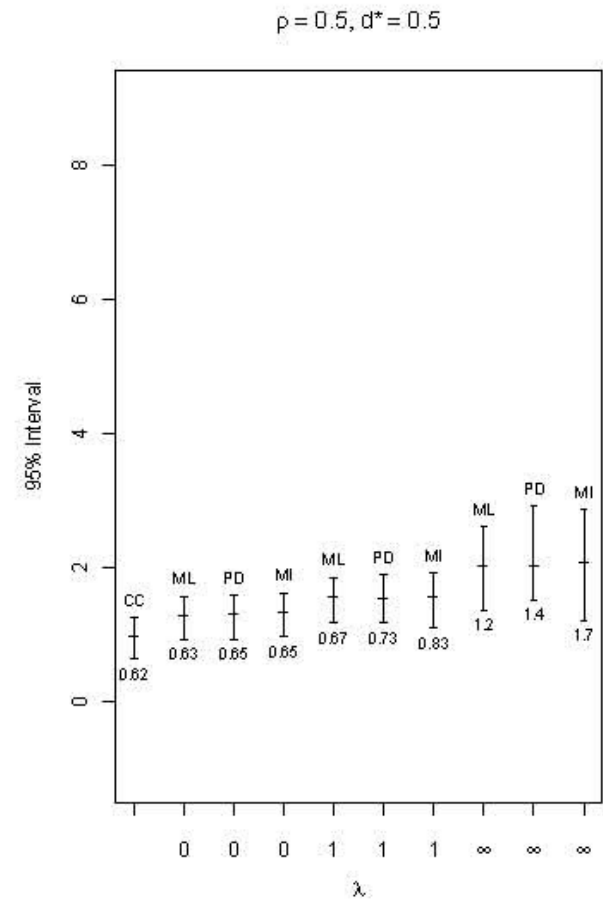
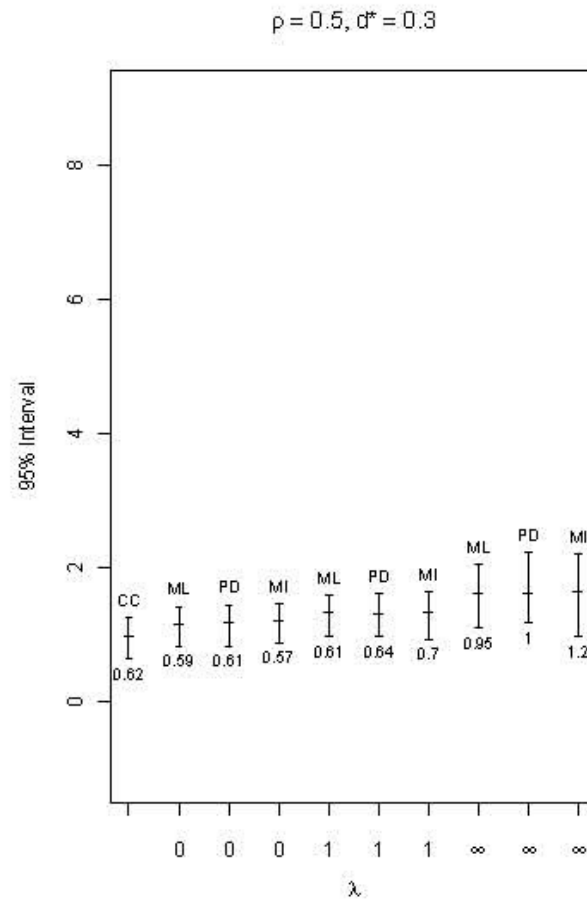
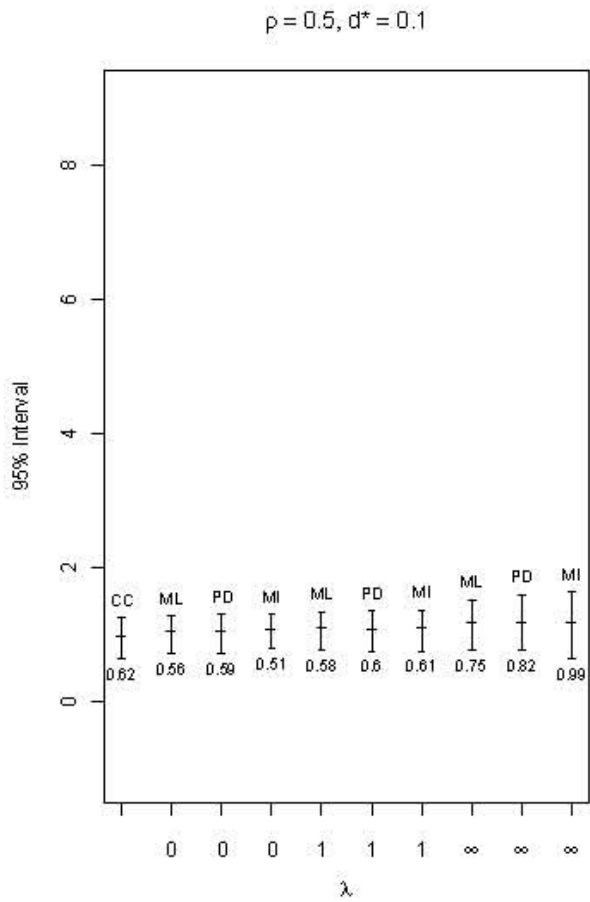
Table 1: Coverage and confidence interval length for eighteen artificial data sets. ML: Maximum likelihood; PD: Posterior distribution; MI: 20 multiply imputed data sets. Results over 500 replicates.

Population			n=100						n=400					
			Coverage			CI Width			Coverage			CI Width		
ρ	d	λ	ML	PD	MI	ML	PD	MI	ML	PD	MI	ML	PD	MI
0.8	0.1	0	93	94	93	0.46	0.47	0.47	95	94	94	0.23	0.23	0.23
		1	95	95	95	0.47	0.48	0.48	95	95	95	0.24	0.24	0.24
		∞	95	95	96	0.51	0.53	0.53	95	95	94	0.25	0.25	0.26
0.8	0.3	0	94	94	94	0.48	0.49	0.49	96	95	95	0.24	0.24	0.24
		1	96	96	96	0.5	0.51	0.51	96	95	96	0.25	0.25	0.25
		∞	96	95	96	0.55	0.57	0.58	96	95	96	0.27	0.27	0.28
0.8	0.5	0	95	96	95	0.52	0.53	0.53	96	95	95	0.26	0.25	0.26
		1	96	97	96	0.54	0.56	0.56	96	95	97	0.27	0.27	0.28
		∞	97	96	97	0.63	0.65	0.66	97	96	96	0.31	0.31	0.32
0.5	0.1	0	93	93	93	0.52	0.54	0.54	94	93	94	0.26	0.26	0.27
		1	95	96	95	0.57	0.6	0.6	95	95	96	0.29	0.28	0.29
		∞	97	95	97	0.96	1.4	2.9	96	95	95	0.42	0.43	0.44
0.5	0.3	0	93	94	94	0.54	0.56	0.57	94	94	94	0.27	0.27	0.28
		1	96	97	96	0.59	0.66	0.66	95	95	96	0.3	0.31	0.31
		∞	96	96	97	1.3	2.1	6.9	95	95	96	0.52	0.54	0.57
0.5	0.5	0	95	95	95	0.58	0.6	0.61	94	94	95	0.29	0.29	0.3
		1	97	97	98	0.64	0.77	0.76	95	96	97	0.33	0.35	0.36
		∞	96	97	96	1.8	3.2	11	97	96	96	0.68	0.7	0.75
0.2	0.1	0	93	94	94	0.55	0.57	0.58	94	93	93	0.28	0.28	0.28
		1	94	96	96	0.64	0.78	0.78	95	95	95	0.32	0.34	0.35
		∞	94	97	97	111	15	34	94	97	96	2.1	5	14
0.2	0.3	0	94	95	94	0.58	0.59	0.6	95	94	94	0.29	0.29	0.3
		1	87	96	94	0.67	1.1	1.1	95	96	97	0.34	0.45	0.46
		∞	87	96	93	1309	31	75	90	97	94	6.8	13	33
0.2	0.5	0	95	95	95	0.62	0.63	0.65	96	95	94	0.31	0.31	0.32
		1	86	98	97	0.73	1.7	1.6	95	97	98	0.36	0.63	0.63
		∞	85	96	94	2509	50	126	90	97	96	12	22	54

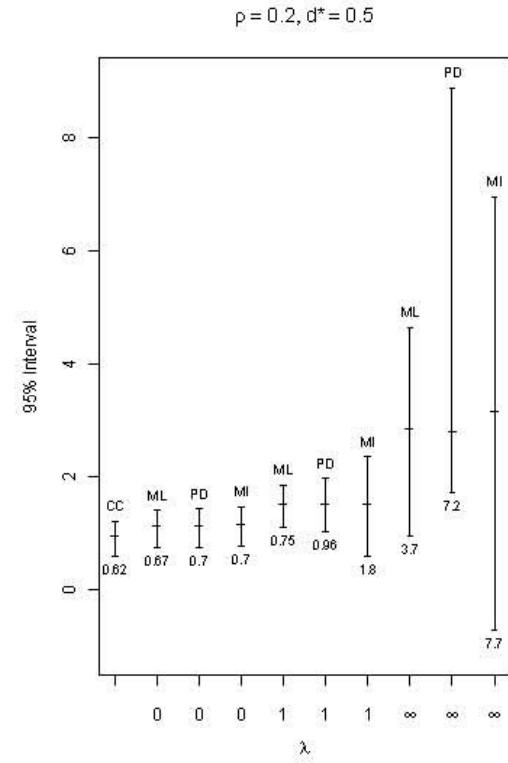
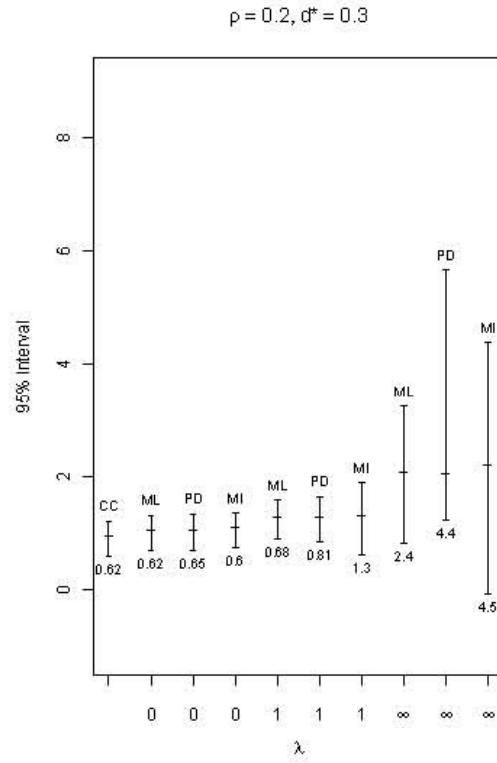
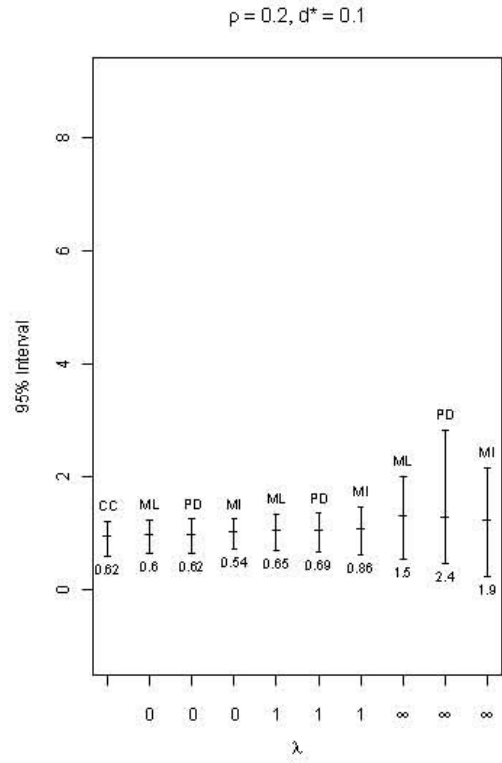
Figure 2: 95% confidence intervals for nine generated data sets ($n = 100$) for $\lambda = (0, 1, \infty)$. Numbers below intervals are the interval length. CC: Complete case; ML: Maximum likelihood; PD: Posterior distribution; MI: 20 multiply imputed data sets. **Rho = 0.8**



Rho = 0.5



Rho = 0.2



Simulation findings

- ML/Bayes are similar with good confidence coverage if n is large, or for strong proxies (assuming λ is correctly chosen)
- For small n , weak proxies, Bayes is more conservative and has better confidence coverage
- Weak proxies lead to greatly increased uncertainty under this framework

Extensions of normal PPMMA

- Non-normal outcomes
 - Transformation may improve normality
 - Extensions to categorical variables via probit models, (Andridge thesis)
- Incomplete covariates
 - Incomplete covariates can be handled by appending procedure to MI of the missing covariates via “chained equations” (IVEware, MICE)
 - Run a chained equation for each choice of lambda

Attractive features of PPMa

- Integrates various components of nonresponse into a single sensitivity analysis reflecting the hierarchy of evidence about bias in the mean
- Easy to implement
- Includes but does not assume MAR; sensitivity analysis is preferred method of assessing NMAR nonresponse
- Gives appropriate credit to the existence of good predictors of the observed outcomes
 - Reinforces that emphasis should be on collecting strong auxiliary data, not solely on obtaining the highest possible response rate

Potential “Disadvantages”

- The interpretation of lambda is complicated by choosing a best proxy for Y
 - This is the price for limiting deviation from MAR to a single parameter
- Analysis needs to be repeated on each of the key outcomes -- no single measure is readily available
 - BUT this is a reflection of reality, not a limitation
- Gives bad news unless covariates are correlated with outcome
 - Including MNAR situations results in more uncertainty

References

- Andridge, R.H. & Little, R. J. (2010). *Int. Statist. Rev.* 78, 1, 40-64.
- Andridge, R.H. & Little, R.J. (2011). To appear in *JOS*.
- Giusti, C. & Little, R.J. (2011). To appear in *JOS*.
- Heckman, J.J. (1976). *Ann. Econ. Soc. Meas.* 5, 475–492.
- IVEware. See <http://www.isr.umich.edu/src/smp/ive/>
- Little, R.J. (1993) *JOS* 9(2), 407-426.
- Little, R.J. (1994). *Biometrika* 81, 3, 471-483.
- Little, R.J. & An, H. (2004). *Statist. Sinica.*, 14, 949-968.
- Little, R.J., & Rubin, D.B. (2002). *Statistical Analysis with Missing Data*, 2nd. ed. Wiley.
- MICE . See <http://web.inter.nl.net/users/S.van.Buuren/mi/html/mice.htm>
- Rubin, D.B. (1976). *Biometrika* **63**, 581-592.
- Rubin, D.B. (1977). *JASA* **72**, 538-543.
- SAS Procs: Mixed, MI, MIANALYZE
- Wang, C., Little, R.J., Nan, B. & Harlow, S. (2011). To appear in *Biometrics*.
- Zhang, G. & Little, R. J. (2009). *Biometrics*, 65, 911-918.
- Zhang, G. & Little, R. J. (2011). To appear in *J. Statist. Comp. Sim.*

and thanks to my recent students...

Hyonggin An, Qi Long, Ying Yuan, Guangyu Zhang, Xiaoxi Zhang, Di An, Yan Zhou, Rebecca Andridge, Qixuan Chen, Ying Guo, Chia-Ning Wang, Nanhua Zhang