

# Statistical Analysis with Missing Data

Module 6

Bayes Inference



# Objectives

- Gibbs' sampler to simulate posterior distribution of parameters
- Bayesian theory of Multiple Imputation under explicit models
- Bayes generates proper multiple imputations – propagates error in estimating parameters

# Gibbs sampling for missing-data problems

$y = (y_{(0)}, y_{(1)})$ ,  $y_{(0)}$  = observed data,  $y_{(1)}$  = missing data; assume MAR

Model for full data:  $f_Y(y | \theta)$ ; prior:  $\pi(\theta)$

Objective: draws  $\theta^{(d)}$  from posterior distribution of  $\theta$ , that is:

$$p(\theta | y_{(0)}) \propto \pi(\theta) f(y_{(0)} | \theta)$$

Often easier to draw  $\theta \sim p(\theta | y_{(0)}, y_{(1)})$ , the complete-data posterior distribution, rather than  $\theta \sim p(\theta | y_{(0)})$

Often easier to draw  $y_{(1)} \sim p(y_{(1)} | y_{(0)}, \theta)$  rather than  $y_{(1)} \sim p(y_{(1)} | y_{(0)})$

So, we apply the Gibbs' sampler to  $(y_{(1)}, \theta)$ :

# Gibbs sampler for missing-data

Initial draw of  $\theta = \theta^{(0)}$ ; then draw  $y_{(1)}^{(0)} \sim f_Y(y_{(1)} | y_{(0)}, \theta^{(0)})$

Let  $(\theta^{(t)}, y_{(1)}^{(t)})$  be draws at iteration  $t$ . Then for iteration  $t + 1$  draw:

P step:  $\theta^{(t+1)} \sim p(\theta | y_{(0)}, y_{(1)}^{(t)})$ , posterior for  $\theta$  with  $y_{(1)}^{(t)}$  imputed for  $y_{(1)}$

I step:  $y_{(1)}^{(t+1)} \sim f_Y(y_{(1)} | y_{(0)}, \theta^{(t+1)})$ , predictive dn of  $y_{(1)}$  given  $\theta = \theta^{(t+1)}$

(P for "posterior," I for "imputation."

the order of the P and I steps is not important).

As  $t \rightarrow \infty$ ,  $(\theta^{(t)}, y_{(1)}^{(t)})$  converges to a draw from  $p(\theta, y_{(1)} | y_{(0)})$

After burn-in  $a$ , draws  $\{\theta^{(a+t)}, t = 1, 2, \dots\}$  simulate posterior dn of  $\theta$

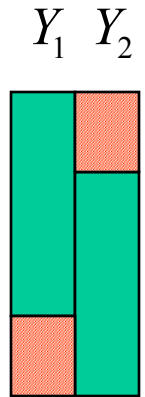
(Recommended: run 2 or more chains to ensure mixing)

# Example: bivariate normal MAR data

- Bivariate normal data with missing data on both variables
- MAR mechanism
- Gibbs' for iteration  $t$  consists of an I step and a P step.
- I-Step is like an E step, except that conditional mean is replaced by a draw:

$$\text{missing } y_{i2} : (y_{i2}^{(t+1)} | y_{i1}, \theta^{(t)}) \sim_{ind} N(\beta_{20.1}^{(t)} + \beta_{21.1}^{(t)} y_{i1}, \sigma_{22.1}^{(t)})$$

$$\text{missing } y_{i1} : (y_{i1}^{(t+1)} | y_{i2}, \theta^{(t)}) \sim_{ind} N(\beta_{10.2}^{(t)} + \beta_{12.2}^{(t)} y_{i2}, \sigma_{11.2}^{(t)})$$



P-Step is like M-Step of EM, with maximization

replaced by draw from complete-data posterior distribution:

$$\Sigma^{(t+1)} \sim \text{Inv-Wishart}(S^{(t+1)}, n-1)$$

$$\mu^{(t+1)} | \Sigma^{(t+1)} \sim N(\bar{x}^{(t+1)}, \Sigma^{(t+1)})$$

# Bayes and multiple imputation

Draws  $y_{(1)}^{(t)}$  from  $p(y_{(1)} | y_{(0)})$  can also be used to create

multiply-imputed data sets  $((y_{(0)}, y_{(1)}^{(d)}), d = 1, \dots, D)$

E.g. impute missing values  $(y_{(1)}^{(a+db)})$  for  $d$ th MI dataset,

$b$  chosen so that imputations are roughly uncorrelated

Or run a separate chain for each MI data set.

- The reason is that the MI combining rules are Bayesian: specifically, as I now discuss, they are simulation approximations of the posterior mean and variance under a Bayesian model

# MI Inference for a Scalar Estimand

$\theta$  = estimand of interest

$\hat{\theta}_d$  = estimate from  $d$ th dataset ( $d = 1, \dots, D$ )

The MI estimate of  $\theta$  is  $\bar{\theta}_D = \frac{1}{D} \sum_{d=1}^D \hat{\theta}_d$

$W_d$  = estimate of variance of  $\hat{\theta}_d$  from  $d$ th dataset

The MI estimate of variance is  $T_D = \bar{W}_D + (1 + 1/D)B_D$

$\bar{W}_D = \frac{1}{D} \sum_{d=1}^D W_d$  = Within-Imputation Variance

$B_D = \frac{1}{D-1} \sum_{d=1}^D (\hat{\theta}_d - \bar{\theta}_D)^2$  = Between-Imputation Variance

# Bayesian Theory of MI

Model:  $f_Y(y | \theta) \Rightarrow$  Likelihood  $L(\theta | y) \propto f_Y(y | \theta)$

Prior distribution:  $\pi(\theta)$ ; md mechanism: MAR

$y = (y_{(0)}, \mathbf{y}_{(1)})$ ,  $y_{(0)}$  = observed data,  $\mathbf{y}_{(1)}$  = missing data

Complete-data posterior distribution,

if there were no missing values:

$$p(\theta | y_{(0)}, \mathbf{y}_{(1)}) \propto \pi(\theta) f_Y(y_{(0)}, \mathbf{y}_{(1)} | \theta)$$

Posterior distribution given observed data:

$$p(\theta | y_{(0)}) \propto \pi(\theta) f(y_{(0)} | \theta)$$

Theory relates these two distributions ...



# Relating the posteriors

- The posterior is related to the complete-data posterior by:

$$p(\theta | y_{(0)}) = \int p(\theta | y_{(0)}, \mathbf{y}_{(1)}) p(\mathbf{y}_{(1)} | y_{(0)}) d\mathbf{y}_{(1)}$$
$$\approx \frac{1}{D} \sum_{d=1}^D p(\theta | y_{(0)}, \mathbf{y}_{(1)}^{(d)}), \text{ where } \mathbf{y}_{(1)}^{(d)} \sim p(\mathbf{y}_{(1)} | y_{(0)})$$

$\mathbf{y}_{(1)}^{(d)}$  is a draw from the predictive distribution of the missing values

The accuracy of the approximation increases with  $D$  and the fraction of observed data

# MI approximation to posterior mean

- Similar approximations for posterior mean and variance yield the MI combining rules given earlier:

$$\begin{aligned} E(\theta | y_{(0)}) \\ &= \int E(\theta | y_{(0)}, \mathbf{y}_{(1)}) p(\mathbf{y}_{(1)} | y_{(0)}) d\mathbf{y}_{(1)} \\ &\approx \frac{1}{D} \sum_{d=1}^D E(\theta | y_{(0)}, \mathbf{y}_{(1)}^{(d)}) = \frac{1}{D} \sum_{d=1}^D \hat{\theta}_d, \end{aligned}$$

where  $\hat{\theta}_d$  = is posterior mean from  $d$ th imputed dataset

# MI approximation to posterior variance

$$\text{Var}(\theta \mid y_{(0)}) = E(\theta^2 \mid y_{(0)}) - \left(E(\theta \mid y_{(0)})\right)^2$$

Apply above approx to  $E(\theta \mid y_{(0)})$  and  $E(\theta^2 \mid y_{(0)})$

Algebra then yields:

$$\text{Var}(\theta \mid y_{(0)}) \approx \bar{V} + B$$

$$\bar{V} = \frac{1}{D} \sum_{d=1}^D V_d = \text{within-imputation variance,}$$

$V_d = \text{Var}(\theta \mid y_{(0)}, y_{(1)}^{(d)})$  is posterior variance from  $d$ th dataset

$$B = \frac{1}{D-1} \sum_{d=1}^D \left( \hat{\theta}_d - \bar{\theta}_D \right)^2 = \text{between-imputation variance}$$

# Refinements of MI combining rules for small $D$

(A):  $\text{Var}(\theta \mid y_{(0)}) \approx \bar{V} + (1 + 1/D) B$

(B) Replace normal reference distribution by t distribution with df

$$\nu = (D-1) \left( 1 + \frac{D}{D+1} \frac{\bar{V}}{B} \right)^2$$

(C) For normal sample with variance based on  $\nu_{\text{com}}$  df, replace  $\nu$  by

$$\nu^* = \left( \nu^{-1} + \hat{\nu}_{\text{obs}}^{-1} \right)^{-1}, \hat{\nu}_{\text{obs}} = (1 - \hat{\gamma}_D) \left( \frac{\nu_{\text{com}} + 1}{\nu_{\text{com}} + 3} \right) \nu_{\text{com}}$$

$$\hat{\gamma}_D = \frac{(1 + D^{-1})B}{\bar{V} + (1 + D^{-1})B} = \text{estimated fraction of missing information}$$

# Logistic regression example revisited

- Imputation Model

$$X_{edi} \sim \text{iid } N(\mu_{ed}, \sigma^2);$$

$e=0,1, d=0,1$ , subject  $i$

- Imputations are draws from the posterior predictive distribution
- Draw  $\sigma^2$ , then  $\mu_{ed}$  and then missing  $X_{edi}$

# Predictive Distributions

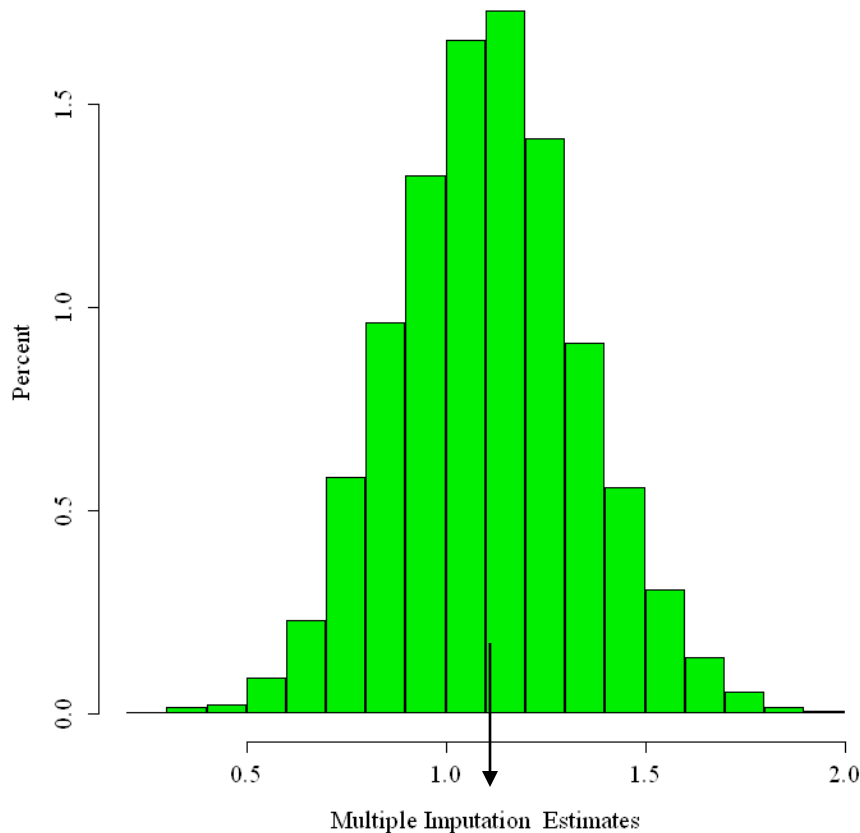
- Draw  $\sigma^2$   
$$\sigma^2 \sim \frac{WSS}{\chi_{r-4}^2}$$
- Draw  $\mu_{ed}$
- $WSS$ =Residual sum of squares
- $r_{ed}$ = Number of respondents in cell  $ed$
- $\bar{x}_{ed}$ = Mean for cell  $ed$

$$\mu_{ed} | X_{obs}, D, E, \sigma^2 \sim N(\bar{x}_{ed}, \sigma^2 / r_{ed})$$

- Draw  $X_{edi} \sim N(\mu_{ed}, \sigma^2)$

# Histogram of Multiple Imputation Estimates

Histogram of 5000 Point Estimates



- 5 Imputations per missing value
- 5 completed Datasets
- Analyze each separately
- Combine using the formulae given earlier

# Coverage and MSE of Various Methods

METHOD	COVERAGE (95% Nominal)	MSE
<i>Before Deletion</i>	94.68	0.0494
Complete-case	37.86	0.4456
Weighted Complete-case	97.42	0.0538
Hot-Deck Single Imputation	90.28	0.0566
Multiple Imputation	94.56	0.0547



# Use of Auxiliary Information in Imputations

- Imputation may involve many more variables though a particular substantive analysis may only use a subset of variables
- Example: Public use data sets or a data set to be used by multiple researchers from different perspectives
- Improve efficiency, reduce bias

# Expanded Simulation Study

- Add auxiliary variable:  $Z \sim N(0,1)$ ,  $\text{Corr}(Z, X)=\rho$

$\rho$	Efficiency of MI Using Z compared to Ignoring Z
0.89	1.42
0.71	1.31
0.55	1.21
0.35	1.12
0	0.97

# Bayes or MI?

- Gibbs sampler can be used to simulate posterior distribution of parameters under a particular model – no need for MI data sets and combining rules
- However, MI data sets are useful for non-Bayesian analyses, or situations where model from MI differs from analysis model, for example by including variables as predictors that are not in the final model.

# Conclusions

- Gibbs sampler useful tool for drawing from the posterior distribution when data are incomplete
- Multiple imputations are a by-product of Gibbs, and can be useful for other analyses
- Other Bayesian simulation methods (SIR, Metropolis-Hastings) can also be useful for handling models where Gibbs is not straightforward