

# STA 360/602L: MODULE 4.5

## MISSING DATA AND IMPUTATION I

DR. OLANREWaju MICHAEL AKANDE

# INTRODUCTION TO MISSING DATA

- Missing data/nonresponse is fairly common in real data. For example,
  - Failure to respond to survey question
  - Subject misses some clinic visits out of all possible
  - Only subset of subjects asked certain questions
- Recall that our posterior computation usually depends on the data through  $p(Y|\theta)$ , which cannot be computed (at least directly) when some of the  $y_i$  values are missing.
- The most common software packages often throw away all subjects with incomplete data (can lead to bias and precision loss).
- Some individuals impute missing values with a mean or some other fixed value (ignores uncertainty).
- As you will see, imputing missing data is actually quite natural in the Bayesian context.

# MISSING DATA MECHANISMS

- Data are said to be **missing completely at random (MCAR)** if the reason for missingness does not depend on the values of the observed data or missing data.
- For example, suppose
  - you handed out a double-sided survey questionnaire of 20 questions to a sample of participants;
  - questions 1-15 were on the first page but questions 16-20 were at the back; and
  - some of the participants did not respond to questions 16-20.
- Then, the values for questions 16-20 for those people who did not respond would be **MCAR** if they simply did not realize the pages were double-sided; they had no reason to ignore those questions.
- **This is rarely plausible in practice!**

# MISSING DATA MECHANISMS

- Data are said to be **missing at random (MAR)** if, conditional on the values of the observed data, the reason for missingness does not depend on the missing data.
- Using our previous example, suppose
  - questions 1-15 include demographic information such as age and education;
  - questions 16-20 include income related questions; and
  - once again, some participants did not respond to questions 16-20.
- Then, the values for questions 16-20 for those people who did not respond would be **MAR** if younger people are more likely not to respond to those income related questions than old people, where age is observed for all participants.
- **This is the most commonly assumed mechanism in practice!**

# MISSING DATA MECHANISMS

- Data are said to be **missing not at random (MNAR or NMAR)** if the reason for missingness depends on the actual values of the missing (unobserved) data.
- Continuing with our previous example, suppose again that
  - questions 1-15 include demographic information such as age and education;
  - questions 16-20 include income related questions; and
  - once again, some of the participants did not respond to questions 16-20.
- Then, the values for questions 16-20 for those people who did not respond would be **MNAR** if people who earn more money are less likely to respond to those income related questions than old people.
- **This is usually the case in real data, but analysis can be complex!**

# MATHEMATICAL FORMULATION

- Consider the multivariate data scenario with  $\mathbf{Y}_i = (\mathbf{Y}_1, \dots, \mathbf{Y}_n)^T$ , where  $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{ip})^T$ , for  $i = 1, \dots, n$ .
- For now, we will assume the multivariate normal model as the sampling model, so that each  $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{ip})^T \sim \mathcal{N}_p(\boldsymbol{\theta}, \Sigma)$ .
- Suppose now that  $\mathbf{Y}$  contains missing values.
- We can separate  $\mathbf{Y}$  into the observed and missing parts, that is,  $\mathbf{Y} = (\mathbf{Y}_{obs}, \mathbf{Y}_{mis})$ .
- Then for each individual,  $\mathbf{Y}_i = (\mathbf{Y}_{i,obs}, \mathbf{Y}_{i,mis})$ .

# MATHEMATICAL FORMULATION

- Let
  - $j$  index variables (where  $i$  already indexes individuals),
  - $r_{ij} = 1$  when  $y_{ij}$  is missing,
  - $r_{ij} = 0$  when  $y_{ij}$  is observed.
- Here,  $r_{ij}$  is known as the missingness indicator of variable  $j$  for person  $i$ .
- Also, let
  - $\mathbf{R}_i = (r_{i1}, \dots, r_{ip})^T$  be the vector of missing indicators for person  $i$ .
  - $\mathbf{R} = (\mathbf{R}_1, \dots, \mathbf{R}_n)$  be the matrix of missing indicators for everyone.
  - $\psi$  be the set of parameters associated with  $\mathbf{R}$ .
- Assume  $\psi$  and  $(\theta, \Sigma)$  are distinct.

# MATHEMATICAL FORMULATION

- MCAR:

$$p(\mathbf{R}|\mathbf{Y}, \boldsymbol{\theta}, \Sigma, \boldsymbol{\psi}) = p(\mathbf{R}|\boldsymbol{\Psi})$$

- MAR:

$$p(\mathbf{R}|\mathbf{Y}, \boldsymbol{\theta}, \Sigma, \boldsymbol{\psi}) = p(\mathbf{R}|\mathbf{Y}_{obs}, \boldsymbol{\Psi})$$

- MNAR:

$$p(\mathbf{R}|\mathbf{Y}, \boldsymbol{\theta}, \Sigma, \boldsymbol{\psi}) = p(\mathbf{R}|\mathbf{Y}_{obs}, \mathbf{Y}_{mis}, \boldsymbol{\Psi})$$



# IMPLICATIONS FOR LIKELIHOOD FUNCTION

- Each type of mechanism has a different implication on the likelihood of the observed data  $\mathbf{Y}_{obs}$ , and the missing data indicator  $\mathbf{R}$ .
- Without missingness in  $\mathbf{Y}$ , the likelihood of the observed data is

$$p(\mathbf{Y}_{obs}|\boldsymbol{\theta}, \Sigma)$$

- With missingness in  $\mathbf{Y}$ , the likelihood of the observed data is instead

$$p(\mathbf{Y}_{obs}, \mathbf{R}|\boldsymbol{\theta}, \Sigma, \boldsymbol{\psi}) = \int p(\mathbf{R}|\mathbf{Y}_{obs}, \mathbf{Y}_{mis}, \boldsymbol{\psi}) \cdot p(\mathbf{Y}_{obs}, \mathbf{Y}_{mis}|\boldsymbol{\theta}, \Sigma) d\mathbf{Y}_{mis}$$

- Since we do not actually observe  $\mathbf{Y}_{mis}$ , we would like to be able to integrate it out so we don't have to deal with it.
- That is, we would like to infer  $(\boldsymbol{\theta}, \Sigma)$  (and sometimes,  $\boldsymbol{\psi}$ ) using only the observed data.

# LIKELIHOOD FUNCTION: MCAR

- For MCAR, we have:

$$\begin{aligned} p(\mathbf{Y}_{obs}, \mathbf{R} | \boldsymbol{\theta}, \Sigma, \boldsymbol{\psi}) &= \int p(\mathbf{R} | \mathbf{Y}_{obs}, \mathbf{Y}_{mis}, \boldsymbol{\psi}) \cdot p(\mathbf{Y}_{obs}, \mathbf{Y}_{mis} | \boldsymbol{\theta}, \Sigma) d\mathbf{Y}_{mis} \\ &= \int p(\mathbf{R} | \boldsymbol{\psi}) \cdot p(\mathbf{Y}_{obs}, \mathbf{Y}_{mis} | \boldsymbol{\theta}, \Sigma) d\mathbf{Y}_{mis} \\ &= p(\mathbf{R} | \boldsymbol{\psi}) \cdot \int p(\mathbf{Y}_{obs}, \mathbf{Y}_{mis} | \boldsymbol{\theta}, \Sigma) d\mathbf{Y}_{mis} \\ &= p(\mathbf{R} | \boldsymbol{\psi}) \cdot p(\mathbf{Y}_{obs} | \boldsymbol{\theta}, \Sigma). \end{aligned}$$

- For inference on  $(\boldsymbol{\theta}, \Sigma)$ , we can simply focus on  $p(\mathbf{Y}_{obs} | \boldsymbol{\theta}, \Sigma)$  in the likelihood function, since  $(\mathbf{R} | \boldsymbol{\psi})$  does not include any  $\mathbf{Y}$ .

# LIKELIHOOD FUNCTION: MAR

- For MAR, we have:

$$\begin{aligned} p(\mathbf{Y}_{obs}, \mathbf{R} | \boldsymbol{\theta}, \Sigma, \boldsymbol{\psi}) &= \int p(\mathbf{R} | \mathbf{Y}_{obs}, \mathbf{Y}_{mis}, \boldsymbol{\psi}) \cdot p(\mathbf{Y}_{obs}, \mathbf{Y}_{mis} | \boldsymbol{\theta}, \Sigma) d\mathbf{Y}_{mis} \\ &= \int p(\mathbf{R} | \mathbf{Y}_{obs}, \boldsymbol{\psi}) \cdot p(\mathbf{Y}_{obs}, \mathbf{Y}_{mis} | \boldsymbol{\theta}, \Sigma) d\mathbf{Y}_{mis} \\ &= p(\mathbf{R} | \mathbf{Y}_{obs}, \boldsymbol{\psi}) \cdot \int p(\mathbf{Y}_{obs}, \mathbf{Y}_{mis} | \boldsymbol{\theta}, \Sigma) d\mathbf{Y}_{mis} \\ &= p(\mathbf{R} | \mathbf{Y}_{obs}, \boldsymbol{\psi}) \cdot p(\mathbf{Y}_{obs} | \boldsymbol{\theta}, \Sigma). \end{aligned}$$

- For inference on  $(\boldsymbol{\theta}, \Sigma)$ , we can once again focus on  $p(\mathbf{Y}_{obs} | \boldsymbol{\theta}, \Sigma)$  in the likelihood function, although there can be some bias if we do not account for  $p(\mathbf{R} | \mathbf{Y}_{obs}, \mathbf{X}, \boldsymbol{\theta})$ , since it contains observed data.
- Also, if we want to infer the missingness mechanism through  $\boldsymbol{\psi}$ , we would need to deal with  $p(\mathbf{R} | \mathbf{Y}_{obs}, \mathbf{X}, \boldsymbol{\theta})$  anyway.

# LIKELIHOOD FUNCTION: MNAR

- For MNAR, we have:

$$p(\mathbf{Y}_{obs}, \mathbf{R} | \boldsymbol{\theta}, \Sigma, \boldsymbol{\psi}) = \int p(\mathbf{R} | \mathbf{Y}_{obs}, \mathbf{Y}_{mis}, \boldsymbol{\psi}) \cdot p(\mathbf{Y}_{obs}, \mathbf{Y}_{mis} | \boldsymbol{\theta}, \Sigma) d\mathbf{Y}_{mis}$$

- The likelihood under MNAR cannot simplify any further.
- In this case, we cannot ignore the missing data when making inferences about  $(\boldsymbol{\theta}, \Sigma)$ .
- We must include the model for  $\mathbf{R}$  and also infer the missing data  $\mathbf{Y}_{mis}$ .

# HOW TO TELL IN PRACTICE?

- So how can we tell the type of mechanism we are dealing with?
- In general, we don't know!!!
- Rare that data are MCAR (unless planned beforehand); more likely that data are MNAR.
- **Compromise:** assume data are MAR if we include enough variables in model for the missing data indicator  $R$ .
- Whenever we talk about missing data in this course, we will do so in the context of MCAR and MAR.

# WHAT'S NEXT?

MOVE ON TO THE READINGS FOR THE NEXT MODULE!