

# STA 601: Missing data and imputation

Merlise Clyde

Nov 9, 2021



# 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



# 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.



# 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).



# 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).



# 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).



# 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.



# 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.



# 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.



# 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.



# 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.



# 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.



# 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.



# 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.



# 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.



# 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.



# 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  $i = 1, \dots, n$ .



# 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  $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{ip})^T \sim \mathcal{N}_p(\boldsymbol{\theta}, \Sigma)$ .



# 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.



# 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,  $\mathbf{Y} = (\mathbf{Y}_{obs}, \mathbf{Y}_{mis})$ .



# 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.



# 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  $i$ .



# 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  $\mathbf{R}$ .



# 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:

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



# Mathematical Formulation

- MCAR:

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

- MAR:

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



# 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  $\mathbf{R}$ .



# 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

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



# 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

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



# 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

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

- 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.



# 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

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

- 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,  $\psi$ ) using only the observed data.



# Likelihood function: MCAR

- For MCAR, we have:

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



# Likelihood function: MCAR

- For MCAR, we have:

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

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



# Likelihood function: MAR

- For MAR, we have:

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



# Likelihood function: MAR

- For MAR, we have:

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

- For inference on  $(\boldsymbol{\theta}, \Sigma)$ , we can once again focus on  $p(\mathbf{Y}_{obs} | \boldsymbol{\theta}, \Sigma)$  in the likelihood function



# Likelihood function: MAR

- For MAR, we have:

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

- For inference on  $(\boldsymbol{\theta}, \Sigma)$ , we can once again focus on  $p(\mathbf{Y}_{obs} | \boldsymbol{\theta}, \Sigma)$  in the likelihood function
- Also, if we want to infer the missingness mechanism through  $\psi$ , we would need to deal with  $p(\mathbf{R} | \mathbf{Y}_{obs}, \psi)$  anyway.

# Likelihood function: MNAR

- For MNAR, we have:

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



# Likelihood function: MNAR

- For MNAR, we have:

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

- The likelihood under MNAR cannot simplify any further.



# Likelihood function: MNAR

- For MNAR, we have:

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

- The likelihood under MNAR cannot simplify any further.
- In this case, we cannot ignore the missing data when making inferences about  $(\boldsymbol{\theta}, \Sigma)$ .

# Likelihood function: MNAR

- For MNAR, we have:

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

- 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?



# How to tell in practice?

- So how can we tell the type of mechanism we are dealing with?
- In general, we don't know!!!



# 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.



# 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  $R$ .



# 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.



# Bayesian inference with missing data

- As we have seen, for MCAR and MAR, we can focus on  $p(\mathbf{Y}_{obs}|\boldsymbol{\theta}, \Sigma)$  in the likelihood function, when inferring  $(\boldsymbol{\theta}, \Sigma)$ .



# Bayesian inference with missing data

- As we have seen, for MCAR and MAR, we can focus on  $p(\mathbf{Y}_{obs}|\boldsymbol{\theta}, \Sigma)$  in the likelihood function, when inferring  $(\boldsymbol{\theta}, \Sigma)$ .
- While this is great, for posterior sampling under most models (especially multivariate models), we actually do need all the  $\mathbf{y}$ 's to update the parameters.



# Bayesian inference with missing data

- As we have seen, for MCAR and MAR, we can focus on  $p(\mathbf{Y}_{obs}|\boldsymbol{\theta}, \Sigma)$  in the likelihood function, when inferring  $(\boldsymbol{\theta}, \Sigma)$ .
- While this is great, for posterior sampling under most models (especially multivariate models), we actually do need all the  $\mathbf{y}$ 's to update the parameters.
- In addition, we may actually want to learn about the missing values, in addition to inferring  $(\boldsymbol{\theta}, \Sigma)$ .



# Bayesian inference with missing data

- As we have seen, for MCAR and MAR, we can focus on  $p(\mathbf{Y}_{obs}|\boldsymbol{\theta}, \Sigma)$  in the likelihood function, when inferring  $(\boldsymbol{\theta}, \Sigma)$ .
- While this is great, for posterior sampling under most models (especially multivariate models), we actually do need all the  $\mathbf{y}$ 's to update the parameters.
- In addition, we may actually want to learn about the missing values, in addition to inferring  $(\boldsymbol{\theta}, \Sigma)$ .
- By thinking of the missing data as **another set of parameters**, we can sample them from the "posterior predictive" distribution of the missing data conditional on the observed data and parameters:

$$p(\mathbf{Y}_{mis}|\mathbf{Y}_{obs}, \boldsymbol{\theta}, \Sigma) \propto \prod_{i=1}^n p(\mathbf{Y}_{i,mis}|\mathbf{Y}_{i,obs}, \boldsymbol{\theta}, \Sigma).$$



# Bayesian inference with missing data

- As we have seen, for MCAR and MAR, we can focus on  $p(\mathbf{Y}_{obs}|\boldsymbol{\theta}, \Sigma)$  in the likelihood function, when inferring  $(\boldsymbol{\theta}, \Sigma)$ .
- While this is great, for posterior sampling under most models (especially multivariate models), we actually do need all the  $\mathbf{y}$ 's to update the parameters.
- In addition, we may actually want to learn about the missing values, in addition to inferring  $(\boldsymbol{\theta}, \Sigma)$ .
- By thinking of the missing data as **another set of parameters**, we can sample them from the "posterior predictive" distribution of the missing data conditional on the observed data and parameters:

$$p(\mathbf{Y}_{mis}|\mathbf{Y}_{obs}, \boldsymbol{\theta}, \Sigma) \propto \prod_{i=1}^n p(\mathbf{Y}_{i,mis}|\mathbf{Y}_{i,obs}, \boldsymbol{\theta}, \Sigma).$$



# Gibbs sampler with missing data

At iteration  $s + 1$ , do the following

1. Sample  $\theta^{(s+1)}$  from its multivariate normal full conditional

$$p(\boldsymbol{\theta}^{(s+1)} | \mathbf{Y}_{obs}, \mathbf{Y}_{mis}^{(s)}, \Sigma^{(s)}).$$



# Gibbs sampler with missing data

At iteration  $s + 1$ , do the following

1. Sample  $\theta^{(s+1)}$  from its multivariate normal full conditional

$$p(\boldsymbol{\theta}^{(s+1)} | \mathbf{Y}_{obs}, \mathbf{Y}_{mis}^{(s)}, \Sigma^{(s)}).$$

2. Sample  $\Sigma^{(s+1)}$  from its inverse-Wishart full conditional

$$p(\Sigma^{(s+1)} | \mathbf{Y}_{obs}, \mathbf{Y}_{mis}^{(s)}, \boldsymbol{\theta}^{(s+1)}).$$



# Gibbs sampler with missing data

At iteration  $s + 1$ , do the following

1. Sample  $\theta^{(s+1)}$  from its multivariate normal full conditional

$$p(\boldsymbol{\theta}^{(s+1)} | \mathbf{Y}_{obs}, \mathbf{Y}_{mis}^{(s)}, \Sigma^{(s)}).$$

2. Sample  $\Sigma^{(s+1)}$  from its inverse-Wishart full conditional

$$p(\Sigma^{(s+1)} | \mathbf{Y}_{obs}, \mathbf{Y}_{mis}^{(s)}, \boldsymbol{\theta}^{(s+1)}).$$

3. For each  $i = 1, \dots, n$ , with at least one "1" value in the missingness indicator vector  $\mathbf{R}_i$ , sample  $\mathbf{Y}_{i,mis}^{(s+1)}$  from the full conditional

$$p(\mathbf{Y}_{i,mis}^{(s+1)} | \mathbf{Y}_{i,obs}, \boldsymbol{\theta}^{(s+1)}, \Sigma^{(s+1)}),$$

which is also multivariate normal, with its form derived from the original sampling model but with the updated parameters, that is,

$$\mathbf{Y}_i^{(s+1)} = (\mathbf{Y}_{i,obs}, \mathbf{Y}_{i,mis}^{(s+1)})^T \sim \mathcal{N}_p(\boldsymbol{\theta}^{(s+1)}, \Sigma^{(s+1)}).$$



# Gibbs sampler with missing data

- Rewrite  $\mathbf{Y}_i^{(s+1)} = (\mathbf{Y}_{i,mis}, \mathbf{Y}_{i,obs}^{(s+1)})^T \sim \mathcal{N}_p(\boldsymbol{\theta}^{(s+1)}, \Sigma^{(s+1)})$  as

$$\mathbf{Y}_i = \begin{pmatrix} \mathbf{Y}_{i,mis} \\ \mathbf{Y}_{i,obs} \end{pmatrix} \sim \mathcal{N}_p \left[ \begin{pmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right],$$

so that we can take advantage of the conditional normal results.



# Gibbs sampler with missing data

- Rewrite  $\mathbf{Y}_i^{(s+1)} = (\mathbf{Y}_{i,mis}, \mathbf{Y}_{i,obs}^{(s+1)})^T \sim \mathcal{N}_p(\boldsymbol{\theta}^{(s+1)}, \Sigma^{(s+1)})$  as

$$\mathbf{Y}_i = \begin{pmatrix} \mathbf{Y}_{i,mis} \\ \mathbf{Y}_{i,obs} \end{pmatrix} \sim \mathcal{N}_p \left[ \begin{pmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right],$$

so that we can take advantage of the conditional normal results.

- That is, we have

$$\mathbf{Y}_{i,mis} | \mathbf{Y}_{i,obs} = \mathbf{y}_{i,obs} \sim \mathcal{N} \left( \boldsymbol{\theta}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\mathbf{y}_{i,obs} - \boldsymbol{\theta}_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \right).$$

as the multivariate normal distribution (or univariate normal distribution if  $\mathbf{Y}_i$  only has one missing entry) we need in step 3 of the Gibbs sampler.



# Gibbs sampler with missing data

- Rewrite  $\mathbf{Y}_i^{(s+1)} = (\mathbf{Y}_{i,mis}, \mathbf{Y}_{i,obs}^{(s+1)})^T \sim \mathcal{N}_p(\boldsymbol{\theta}^{(s+1)}, \Sigma^{(s+1)})$  as

$$\mathbf{Y}_i = \begin{pmatrix} \mathbf{Y}_{i,mis} \\ \mathbf{Y}_{i,obs} \end{pmatrix} \sim \mathcal{N}_p \left[ \begin{pmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right],$$

so that we can take advantage of the conditional normal results.

- That is, we have

$$\mathbf{Y}_{i,mis} | \mathbf{Y}_{i,obs} = \mathbf{y}_{i,obs} \sim \mathcal{N} \left( \boldsymbol{\theta}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\mathbf{y}_{i,obs} - \boldsymbol{\theta}_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \right).$$

as the multivariate normal distribution (or univariate normal distribution if  $\mathbf{Y}_i$  only has one missing entry) we need in step 3 of the Gibbs sampler.

- This sampling technique actually encodes MAR since the imputations for  $\mathbf{Y}_{mis}$  depend on the  $\mathbf{Y}_{obs}$ .



# Gibbs sampler with missing data

- Rewrite  $\mathbf{Y}_i^{(s+1)} = (\mathbf{Y}_{i,mis}, \mathbf{Y}_{i,obs}^{(s+1)})^T \sim \mathcal{N}_p(\boldsymbol{\theta}^{(s+1)}, \Sigma^{(s+1)})$  as

$$\mathbf{Y}_i = \begin{pmatrix} \mathbf{Y}_{i,mis} \\ \mathbf{Y}_{i,obs} \end{pmatrix} \sim \mathcal{N}_p \left[ \begin{pmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right],$$

so that we can take advantage of the conditional normal results.

- That is, we have

$$\mathbf{Y}_{i,mis} | \mathbf{Y}_{i,obs} = \mathbf{y}_{i,obs} \sim \mathcal{N} \left( \boldsymbol{\theta}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\mathbf{y}_{i,obs} - \boldsymbol{\theta}_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \right).$$

as the multivariate normal distribution (or univariate normal distribution if  $\mathbf{y}_i$  only has one missing entry) we need in step 3 of the Gibbs sampler.

- This sampling technique actually encodes MAR since the imputations for  $\mathbf{y}_{mis}$  depend on the  $\mathbf{y}_{obs}$ .
- Now let's look at the reading comprehension example in Hoff. We will add missing values to the original data and refit the model.



# Reading example with missing data

```
Y <- as.matrix(dget("http://www2.stat.duke.edu/~pdh10/FCBS/Inline,  
#Add 20% missing data; MCAR  
set.seed(1234)  
Y_WithMiss <- Y #So we can keep the full data  
Miss_frac <- 0.20  
R <- matrix(rbinom(nrow(Y_WithMiss)*ncol(Y_WithMiss),1, Miss_frac),  
nrow(Y_WithMiss),ncol(Y_WithMiss))  
Y_WithMiss[R==1] <- NA  
Y_WithMiss[1:12,]
```

```
##      pretest posttest  
## [1,]      59      77  
## [2,]      43      39  
## [3,]      34      46  
## [4,]      32      NA  
## [5,]      NA      38  
## [6,]      38      NA  
## [7,]      55      NA  
## [8,]      67      86  
## [9,]      64      77  
## [10,]     45      60  
## [11,]     49      50
```



# Compare to inference from full data

With missing data:

```
apply(THETA_WithMiss,2,summary)
```

```
##           theta_1   theta_2
## Min.      30.45459 38.29322
## 1st Qu.   43.65988 51.96991
## Median    45.60829 54.19592
## Mean      45.63192 54.20408
## 3rd Qu.   47.61896 56.48918
## Max.      58.81206 70.49105
```

Based on true data:

```
apply(THETA,2,summary)
```

```
##           theta_1   theta_2
## Min.      34.88365 37.80999
## 1st Qu.   45.29473 51.47834
## Median    47.28229 53.65172
## Mean      47.26301 53.64100
## 3rd Qu.   49.21423 55.81819
```



# Compare to inference from full data

With missing data:

```
apply(SIGMA_WithMiss,2,summary)
```

```
##          sigma_11  sigma_12  sigma_21  sigma_22
## Min.      64.0883 -20.39204 -20.39204  82.55346
## 1st Qu.   149.8338 109.84218 109.84218 190.25962
## Median    182.4496 142.34686 142.34686 233.43312
## Mean      193.9803 152.12898 152.12898 248.67527
## 3rd Qu.   224.0994 182.75082 182.75082 289.47663
## Max.      734.8704 668.77332 668.77332 981.99916
```

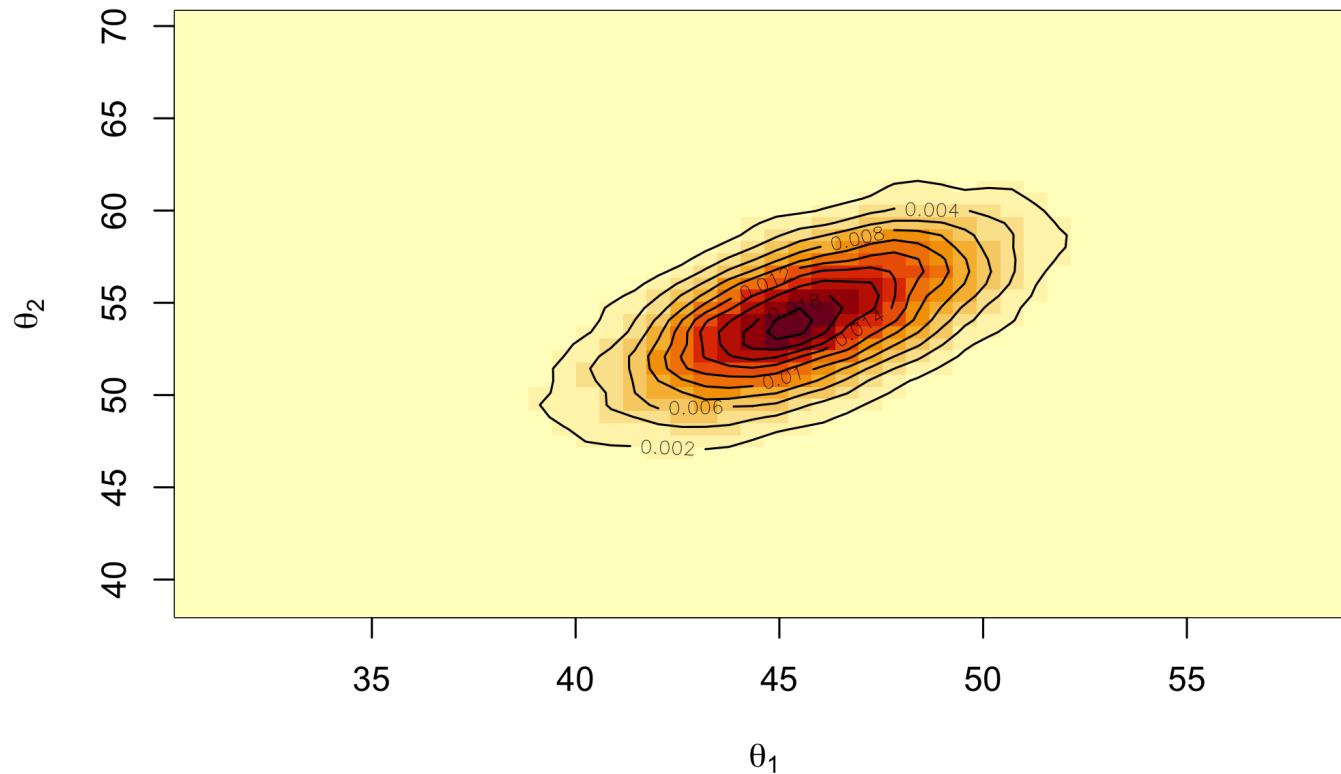
Based on true data:

```
apply(SIGMA,2,summary)
```

```
##          sigma_11  sigma_12  sigma_21  sigma_22
## Min.      76.4661 -38.75561 -38.75561  93.65776
## 1st Qu.   157.5870 113.32529 113.32529 203.69192
## Median    190.6578 145.08962 145.08962 246.08696
## Mean      201.9547 155.20374 155.20374 260.11361
## 3rd Qu.   233.5809 186.36991 186.36991 300.70840
```



# Posterior distribution of the mean



# Missing data vs predictions for new observations

- How about predictions for completely new observations?



# Missing data vs predictions for new observations

- How about predictions for completely new observations?
- That is, suppose your original dataset plus sampling model is

$$\mathbf{y}_i = (y_{i,1}, y_{i,2})^T \sim \mathcal{N}_2(\boldsymbol{\theta}, \Sigma), i = 1, \dots, n.$$



# Missing data vs predictions for new observations

- How about predictions for completely new observations?
- That is, suppose your original dataset plus sampling model is
$$\mathbf{y}_i = (y_{i,1}, y_{i,2})^T \sim \mathcal{N}_2(\boldsymbol{\theta}, \Sigma), i = 1, \dots, n.$$
- Suppose now you have  $n^*$  new observations with  $y_2^*$  values but no  $y_1^*$ .



# Missing data vs predictions for new observations

- How about predictions for completely new observations?
- That is, suppose your original dataset plus sampling model is
$$\mathbf{y}_i = (y_{i,1}, y_{i,2})^T \sim \mathcal{N}_2(\boldsymbol{\theta}, \Sigma), i = 1, \dots, n.$$
- Suppose now you have  $n^*$  new observations with  $y_2^*$  values but no  $y_1^*$ .
- How can we predict  $y_{i,1}^*$  given  $y_{i,2}^*$ , for  $i = 1, \dots, n^*$ ?



# Missing data vs predictions for new observations

- How about predictions for completely new observations?
- That is, suppose your original dataset plus sampling model is
$$\mathbf{y}_i = (y_{i,1}, y_{i,2})^T \sim \mathcal{N}_2(\boldsymbol{\theta}, \Sigma), i = 1, \dots, n.$$
- Suppose now you have  $n^*$  new observations with  $y_2^*$  values but no  $y_1^*$ .
- How can we predict  $y_{i,1}^*$  given  $y_{i,2}^*$ , for  $i = 1, \dots, n^*$ ?
- Well, we can view this as a "train → test" prediction problem rather than a missing data problem on an original data.



# Missing data vs predictions for new observations

- That is, given the posterior samples of the parameters, and the test values for  $y_{i2}^*$ , draw from the posterior predictive distribution of  $(y_{i,1}^* | y_{i,2}^*, \{(y_{1,1}, y_{1,2}), \dots, (y_{n,1}, y_{n,2})\})$ .



# Missing data vs predictions for new observations

- That is, given the posterior samples of the parameters, and the test values for  $y_{i2}^*$ , draw from the posterior predictive distribution of  $(y_{i,1}^* | y_{i,2}^*, \{(y_{1,1}, y_{1,2}), \dots, (y_{n,1}, y_{n,2})\})$ .
- To sample from this predictive distribution, think of compositional sampling.



# Missing data vs predictions for new observations

- That is, given the posterior samples of the parameters, and the test values for  $y_{i2}^*$ , draw from the posterior predictive distribution of  $(y_{i,1}^* | y_{i,2}^*, \{(y_{1,1}, y_{1,2}), \dots, (y_{n,1}, y_{n,2})\})$ .
- To sample from this predictive distribution, think of compositional sampling.
- That is, for each posterior sample of  $(\theta, \Sigma)$ , sample from  $(y_{i,1} | y_{i,2}, \theta, \Sigma)$ , which is just from the form of the sampling distribution.



# Missing data vs predictions for new observations

- That is, given the posterior samples of the parameters, and the test values for  $y_{i2}^*$ , draw from the posterior predictive distribution of  $(y_{i,1}^* | y_{i,2}^*, \{(y_{1,1}, y_{1,2}), \dots, (y_{n,1}, y_{n,2})\})$ .
- To sample from this predictive distribution, think of compositional sampling.
- That is, for each posterior sample of  $(\theta, \Sigma)$ , sample from  $(y_{i,1} | y_{i,2}, \theta, \Sigma)$ , which is just from the form of the sampling distribution.
- In this case,  $(y_{i,1} | y_{i,2}, \theta, \Sigma)$  is just a normal distribution derived from  $(y_{i,1}, y_{i,2} | \theta, \Sigma)$ , based on the conditional normal formula.



# Missing data vs predictions for new observations

- That is, given the posterior samples of the parameters, and the test values for  $y_{i2}^*$ , draw from the posterior predictive distribution of  $(y_{i,1}^* | y_{i,2}^*, \{(y_{1,1}, y_{1,2}), \dots, (y_{n,1}, y_{n,2})\})$ .
- To sample from this predictive distribution, think of compositional sampling.
- That is, for each posterior sample of  $(\theta, \Sigma)$ , sample from  $(y_{i,1} | y_{i,2}, \theta, \Sigma)$ , which is just from the form of the sampling distribution.
- In this case,  $(y_{i,1} | y_{i,2}, \theta, \Sigma)$  is just a normal distribution derived from  $(y_{i,1}, y_{i,2} | \theta, \Sigma)$ , based on the conditional normal formula.
- No need to incorporate the prediction problem into your original Gibbs sampler!

