Itemwise Missing at Random Modeling for Incomplete Multivariate Data¹

Mauricio Sadinle

Duke University and NISS

Supported by NSF grant SES-11-31897

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

¹Joint work with Jerry Reiter

Incomplete Multivariate Data

Gender	Age	Income	
F	25	60,000	
М	?	?	
?	51	?	
F	?	150,300	

What This Talk is About

- Most common approach to handle missing data: assume the missing data are missing at random (MAR)
- We developed an alternative: assume the missing data are *itemwise* missing at random (IMAR)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

What This Talk is About

- Most common approach to handle missing data: assume the missing data are missing at random (MAR)
- We developed an alternative: assume the missing data are *itemwise* missing at random (IMAR)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Outline

Inference with Missing Data

Itemwise Missing at Random

Take-Home Message



- ► X: Would you lend me \$1,000?
- Want to estimate $\mathbb{P}(X = \text{yes})$

$$\begin{split} \mathbb{P}(X = \mathsf{yes}) &= \mathbb{P}(X = \mathsf{yes} | \mathsf{response}) \mathbb{P}(\mathsf{response}) \\ &+ \mathbb{P}(X = \mathsf{yes} | \mathsf{non-response}) \mathbb{P}(\mathsf{non-response}) \end{split}$$

• Most likely
$$\mathbb{P}(X = \text{yes}) \ll \mathbb{P}(X = \text{yes}|\text{response})$$

Inference impossible without extra, usually untestable, assumptions on how missingness arises

►

- ► X: Would you lend me \$1,000?
- Want to estimate $\mathbb{P}(X = \text{yes})$

$$\begin{split} \mathbb{P}(X = \mathsf{yes}) &= \mathbb{P}(X = \mathsf{yes} | \mathsf{response}) \mathbb{P}(\mathsf{response}) \\ &+ \mathbb{P}(X = \mathsf{yes} | \mathsf{non-response}) \mathbb{P}(\mathsf{non-response}) \end{split}$$

• Most likely
$$\mathbb{P}(X = \text{yes}) \ll \mathbb{P}(X = \text{yes}|\text{response})$$

Inference impossible without extra, usually untestable, assumptions on how missingness arises

►

- ► X: Would you lend me \$1,000?
- Want to estimate $\mathbb{P}(X = \text{yes})$

$$\begin{split} \mathbb{P}(X = \mathsf{yes}) &= \mathbb{P}(X = \mathsf{yes} | \mathsf{response}) \mathbb{P}(\mathsf{response}) \\ &+ \mathbb{P}(X = \mathsf{yes} | \mathsf{non-response}) \mathbb{P}(\mathsf{non-response}) \end{split}$$

• Most likely
$$\mathbb{P}(X = \text{yes}) \ll \mathbb{P}(X = \text{yes}|\text{response})$$

Inference impossible without extra, usually untestable, assumptions on how missingness arises

- ► X: Would you lend me \$1,000?
- Want to estimate $\mathbb{P}(X = \text{yes})$

$$\begin{split} \mathbb{P}(X = \mathsf{yes}) &= \mathbb{P}(X = \mathsf{yes} | \mathsf{response}) \mathbb{P}(\mathsf{response}) \\ &+ \mathbb{P}(X = \mathsf{yes} | \mathsf{non-response}) \mathbb{P}(\mathsf{non-response}) \end{split}$$

• Most likely
$$\mathbb{P}(X = \text{yes}) \ll \mathbb{P}(X = \text{yes}|\text{response})$$

Inference impossible without extra, usually untestable, assumptions on how missingness arises

Gender	Age	Income	M_{Gender}	M_{Age}	M _{Income}	
F	25	60,000	0	0	0	
М	?	?	0	1	1	
?	51	?	1	0	1	
F	?	150,300	0	1	0	

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

In general

- Study variables: $\mathbf{X} = (X_1, \dots, X_p) \in \mathcal{X}$
- Missingness indicators: $\mathbf{M} = (M_1, \dots, M_p) \in \{0, 1\}^p$

• Missingness mechanism: $\mathbb{P}(\mathbf{M}|\mathbf{X})$

Gender	Age	Income	M _{Gender}	M_{Age}	M _{Income}	
F	25	60,000	0	0	0	
М	?	?	0	1	1	
?	51	?	1	0	1	
F	?	150,300	0	1	0	

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

In general

• Study variables: $\mathbf{X} = (X_1, \dots, X_p) \in \mathcal{X}$

• Missingness indicators: $\mathbf{M} = (M_1, \dots, M_p) \in \{0, 1\}^p$

• Missingness mechanism: $\mathbb{P}(\mathbf{M}|\mathbf{X})$

Gender	Age	Income	M _{Gender}	M_{Age}	M _{Income}	
F	25	60,000	0	0	0	
М	?	?	0	1	1	
?	51	?	1	0	1	
F	?	150,300	0	1	0	

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

In general

• Study variables: $\mathbf{X} = (X_1, \dots, X_p) \in \mathcal{X}$

• Missingness indicators: $\mathbf{M} = (M_1, \dots, M_p) \in \{0, 1\}^p$

• Missingness mechanism: $\mathbb{P}(\mathbf{M}|\mathbf{X})$

Gender	Age	Income	M _{Gender}	M_{Age}	M _{Income}	
F	25	60,000	0	0	0	
М	?	?	0	1	1	
?	51	?	1	0	1	
F	?	150,300	0	1	0	

In general

- Study variables: $\mathbf{X} = (X_1, \dots, X_p) \in \mathcal{X}$
- Missingness indicators: $\mathbf{M} = (M_1, \dots, M_p) \in \{0, 1\}^p$
- Missingness mechanism: $\mathbb{P}(\mathbf{M}|\mathbf{X})$

Given $\mathbf{M} = \mathbf{m}$

▶ X_m: missing values (often written as X_{mis})

▶ X_m: observed values (often written as X_{obs})

Example:

$$\blacktriangleright \mathbf{X} = (X_1, X_2, X_3)$$

• If
$$\mathbf{m} = (1, 0, 1)$$
, $\mathbf{X}_{\mathbf{m}} = (X_1, X_3)$, and $\mathbf{X}_{\bar{\mathbf{m}}} = X_2$

Given $\mathbf{M} = \mathbf{m}$

- ▶ X_m: missing values (often written as X_{mis})
- ▶ X_m: observed values (often written as X_{obs})

Example:

$$\blacktriangleright \mathbf{X} = (X_1, X_2, X_3)$$

• If
$$\mathbf{m} = (1, 0, 1)$$
, $\mathbf{X}_{\mathbf{m}} = (X_1, X_3)$, and $\mathbf{X}_{\bar{\mathbf{m}}} = X_2$

Given $\mathbf{M} = \mathbf{m}$

- ▶ X_m: missing values (often written as X_{mis})
- ▶ X_m: observed values (often written as X_{obs})

Example:

•
$$\mathbf{X} = (X_1, X_2, X_3)$$

▶ If
$$\mathbf{m} = (1, 0, 1)$$
, $\mathbf{X}_{\mathbf{m}} = (X_1, X_3)$, and $\mathbf{X}_{\bar{\mathbf{m}}} = X_2$

After Rubin (1976):

Missing at random:

$$\mathbb{P}(\mathsf{M}=\mathsf{m}|\mathsf{X}=\mathsf{x})=\mathbb{P}(\mathsf{M}=\mathsf{m}|\mathsf{X}_{\bar{\mathsf{m}}}=\mathsf{x}_{\bar{\mathsf{m}}})$$

Missing completely at random:

$$\mathbb{P}(\mathsf{M}=\mathsf{m}|\mathsf{X}=\mathsf{x})=\mathbb{P}(\mathsf{M}=\mathsf{m})$$

After Rubin (1976):

► Missing at random:

$$\mathbb{P}(\mathsf{M}=\mathsf{m}|\mathsf{X}=\mathsf{x})=\mathbb{P}(\mathsf{M}=\mathsf{m}|\mathsf{X}_{\bar{\mathsf{m}}}=\mathsf{x}_{\bar{\mathsf{m}}})$$

Missing completely at random:

$$\mathbb{P}(\mathsf{M}=\mathsf{m}|\mathsf{X}=\mathsf{x})=\mathbb{P}(\mathsf{M}=\mathsf{m})$$

Under MAR

$$\begin{split} \int \mathbb{P}(\mathsf{M} = \mathsf{m}|\mathsf{X} = \mathsf{x})f(\mathsf{X} = \mathsf{x})d\mathsf{x}_{\mathsf{m}} &= \int \mathbb{P}(\mathsf{M} = \mathsf{m}|\mathsf{X}_{\tilde{\mathsf{m}}} = \mathsf{x}_{\tilde{\mathsf{m}}})f(\mathsf{X} = \mathsf{x})d\mathsf{x}_{\mathsf{m}} \\ &= \mathbb{P}(\mathsf{M} = \mathsf{m}|\mathsf{X}_{\tilde{\mathsf{m}}} = \mathsf{x}_{\tilde{\mathsf{m}}})\int f(\mathsf{X} = \mathsf{x})d\mathsf{x}_{\mathsf{m}} \end{split}$$

 \Rightarrow we can "ignore" the missingness mechanism

$$\mathbb{P}(M_i = 1 | \mathbf{X} = \mathbf{x}) = \mathbb{P}(M_i = 1 | \mathbf{X}_{-i} = \mathbf{x}_{-i})$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Missingness of an item cannot depend on the value of the item

Under MAR

$$\begin{split} \int \mathbb{P}(\mathsf{M} = \mathsf{m} | \mathsf{X} = \mathsf{x}) f(\mathsf{X} = \mathsf{x}) d\mathsf{x}_{\mathsf{m}} &= \int \mathbb{P}(\mathsf{M} = \mathsf{m} | \mathsf{X}_{\bar{\mathsf{m}}} = \mathsf{x}_{\bar{\mathsf{m}}}) f(\mathsf{X} = \mathsf{x}) d\mathsf{x}_{\mathsf{m}} \\ &= \mathbb{P}(\mathsf{M} = \mathsf{m} | \mathsf{X}_{\bar{\mathsf{m}}} = \mathsf{x}_{\bar{\mathsf{m}}}) \int f(\mathsf{X} = \mathsf{x}) d\mathsf{x}_{\mathsf{m}} \end{split}$$

 \Rightarrow we can "ignore" the missingness mechanism

$$\mathbb{P}(M_i = 1 | \mathbf{X} = \mathbf{x}) = \mathbb{P}(M_i = 1 | \mathbf{X}_{-i} = \mathbf{x}_{-i})$$

(ロ)、(型)、(E)、(E)、 E) の(の)

Missingness of an item cannot depend on the value of the item

Under MAR

$$\begin{split} \int \mathbb{P}(\mathsf{M} = \mathsf{m} | \mathsf{X} = \mathsf{x}) f(\mathsf{X} = \mathsf{x}) d\mathsf{x}_{\mathsf{m}} &= \int \mathbb{P}(\mathsf{M} = \mathsf{m} | \mathsf{X}_{\bar{\mathsf{m}}} = \mathsf{x}_{\bar{\mathsf{m}}}) f(\mathsf{X} = \mathsf{x}) d\mathsf{x}_{\mathsf{m}} \\ &= \mathbb{P}(\mathsf{M} = \mathsf{m} | \mathsf{X}_{\bar{\mathsf{m}}} = \mathsf{x}_{\bar{\mathsf{m}}}) \int f(\mathsf{X} = \mathsf{x}) d\mathsf{x}_{\mathsf{m}} \end{split}$$

 \Rightarrow we can "ignore" the missingness mechanism

$$\mathbb{P}(M_i = 1 | \mathbf{X} = \mathbf{x}) = \mathbb{P}(M_i = 1 | \mathbf{X}_{-i} = \mathbf{x}_{-i})$$

(ロ)、(型)、(E)、(E)、 E) の(の)

Missingness of an item cannot depend on the value of the item

Under MAR

• $\mathbb{P}(M_1 = 1, M_2 = 1 | X_1 = x_1, X_2 = x_2) = c$

$$\mathbb{P}(M_1 = 0, M_2 = 1 | X_1 = x_1, X_2 = x_2) = u(x_1)$$

$$\blacktriangleright \mathbb{P}(M_1 = 1, M_2 = 0 | X_1 = x_1, X_2 = x_2) = v(x_2)$$

$$\mathbb{P}(M_1 = 0, M_2 = 0 | X_1 = x_1, X_2 = x_2) = 1 - c - u(x_1) - v(x_2)$$

Under MAR

• $\mathbb{P}(M_1 = 1, M_2 = 1 | X_1 = x_1, X_2 = x_2) = c$

•
$$\mathbb{P}(M_1 = 0, M_2 = 1 | X_1 = x_1, X_2 = x_2) = u(x_1)$$

$$\blacktriangleright \mathbb{P}(M_1 = 1, M_2 = 0 | X_1 = x_1, X_2 = x_2) = v(x_2)$$

$$\blacktriangleright \mathbb{P}(M_1 = 0, M_2 = 0 | X_1 = x_1, X_2 = x_2) = 1 - c - u(x_1) - v(x_2)$$

Under MAR

• $\mathbb{P}(M_1 = 1, M_2 = 1 | X_1 = x_1, X_2 = x_2) = c$

•
$$\mathbb{P}(M_1 = 0, M_2 = 1 | X_1 = x_1, X_2 = x_2) = u(x_1)$$

$$\blacktriangleright \mathbb{P}(M_1 = 1, M_2 = 0 | X_1 = x_1, X_2 = x_2) = v(x_2)$$

 $\triangleright \mathbb{P}(M_1 = 0, M_2 = 0 | X_1 = x_1, X_2 = x_2) = 1 - c - u(x_1) - v(x_2)$

Under MAR

•
$$\mathbb{P}(M_1 = 1, M_2 = 1 | X_1 = x_1, X_2 = x_2) = c$$

►
$$\mathbb{P}(M_1 = 0, M_2 = 1 | X_1 = x_1, X_2 = x_2) = u(x_1)$$

$$\blacktriangleright \mathbb{P}(M_1 = 1, M_2 = 0 | X_1 = x_1, X_2 = x_2) = v(x_2)$$

▶
$$\mathbb{P}(M_1 = 0, M_2 = 0 | X_1 = x_1, X_2 = x_2) = 1 - c - u(x_1) - v(x_2)$$

Non-Ignorable Missingness Mechanisms

Missing not at random:

$$\mathbb{P}(\mathsf{M}=\mathsf{m}|\mathsf{X}=\mathsf{x})\neq\mathbb{P}(\mathsf{M}=\mathsf{m}|\mathsf{X}_{\tilde{\mathsf{m}}}=\mathsf{x}_{\tilde{\mathsf{m}}})$$

 General interest in developing non-ignorable approaches to handle missing data

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Non-Ignorable Missingness Mechanisms

Missing not at random:

$$\mathbb{P}(\mathsf{M}=\mathsf{m}|\mathsf{X}=\mathsf{x})\neq\mathbb{P}(\mathsf{M}=\mathsf{m}|\mathsf{X}_{\bar{\mathsf{m}}}=\mathsf{x}_{\bar{\mathsf{m}}})$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

 General interest in developing non-ignorable approaches to handle missing data

 Generally speaking, inferences should be based on the full data distribution

 $f(\mathbf{X}, \mathbf{M})$

This distribution is not identifiable

Examples of probabilities that cannot be estimated from the data alone, without extra assumptions

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

$$\blacktriangleright \mathbb{P}(X_1 = x_1, M_1 = 1)$$

$$\mathbb{P}(X_1 = x_1, X_2 = x_2, M_1 = 1, M_2 = 0)$$

 Generally speaking, inferences should be based on the full data distribution

 $f(\mathbf{X}, \mathbf{M})$

This distribution is not identifiable

Examples of probabilities that cannot be estimated from the data alone, without extra assumptions

 $\blacktriangleright \mathbb{P}(X_1 = x_1, M_1 = 1)$

$$\mathbb{P}(X_1 = x_1, X_2 = x_2, M_1 = 1, M_2 = 0)$$

 Generally speaking, inferences should be based on the full data distribution

f(**X**,**M**)

- This distribution is not identifiable
- Examples of probabilities that cannot be estimated from the data alone, without extra assumptions

 $\blacktriangleright \mathbb{P}(X_1 = x_1, M_1 = 1)$

$$\mathbb{P}(X_1 = x_1, X_2 = x_2, M_1 = 1, M_2 = 0)$$

 Generally speaking, inferences should be based on the full data distribution

f(**X**,**M**)

- This distribution is not identifiable
- Examples of probabilities that cannot be estimated from the data alone, without extra assumptions

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- $\blacktriangleright \mathbb{P}(X_1 = x_1, M_1 = 1)$
- $\mathbb{P}(X_1 = x_1, X_2 = x_2, M_1 = 1, M_2 = 0)$

The observed data distribution is all we can identify from samples

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

•
$$f(\mathbf{X}_{\bar{\mathbf{m}}} = \mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{M} = \mathbf{m}) = \int_{\mathcal{X}_{\mathbf{m}}} f(\mathbf{X} = \mathbf{x}, \mathbf{M} = \mathbf{m}) d\mathbf{x}_{\mathbf{m}}$$

▶
$$\mathbb{P}(M_1 = 1, M_2 = 1)$$

•
$$\mathbb{P}(X_1 = x_1, M_1 = 0, M_2 = 1)$$

- $\mathbb{P}(X_2 = x_2, M_1 = 1, M_2 = 0)$
- $\mathbb{P}(X_1 = x_1, X_2 = x_2, M_1 = 0, M_2 = 0)$

The observed data distribution is all we can identify from samples

•
$$f(\mathbf{X}_{\bar{\mathbf{m}}} = \mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{M} = \mathbf{m}) = \int_{\mathcal{X}_{\mathbf{m}}} f(\mathbf{X} = \mathbf{x}, \mathbf{M} = \mathbf{m}) d\mathbf{x}_{\mathbf{m}}$$

•
$$\mathbb{P}(M_1 = 1, M_2 = 1)$$

•
$$\mathbb{P}(X_1 = x_1, M_1 = 0, M_2 = 1)$$

- $\mathbb{P}(X_2 = x_2, M_1 = 1, M_2 = 0)$
- $\mathbb{P}(X_1 = x_1, X_2 = x_2, M_1 = 0, M_2 = 0)$

The observed data distribution is all we can identify from samples

(ロ)、(型)、(E)、(E)、 E) の(の)

•
$$f(\mathbf{X}_{\bar{\mathbf{m}}} = \mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{M} = \mathbf{m}) = \int_{\mathcal{X}_{\mathbf{m}}} f(\mathbf{X} = \mathbf{x}, \mathbf{M} = \mathbf{m}) d\mathbf{x}_{\mathbf{m}}$$

•
$$\mathbb{P}(M_1 = 1, M_2 = 1)$$

•
$$\mathbb{P}(X_1 = x_1, M_1 = 0, M_2 = 1)$$

▶
$$\mathbb{P}(X_2 = x_2, M_1 = 1, M_2 = 0)$$

$$\blacktriangleright \mathbb{P}(X_1 = x_1, X_2 = x_2, M_1 = 0, M_2 = 0)$$

The observed data distribution is all we can identify from samples

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

•
$$f(\mathbf{X}_{\bar{\mathbf{m}}} = \mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{M} = \mathbf{m}) = \int_{\mathcal{X}_{\mathbf{m}}} f(\mathbf{X} = \mathbf{x}, \mathbf{M} = \mathbf{m}) d\mathbf{x}_{\mathbf{m}}$$

For example with two categorical variables:

•
$$\mathbb{P}(M_1 = 1, M_2 = 1)$$

•
$$\mathbb{P}(X_1 = x_1, M_1 = 0, M_2 = 1)$$

• $\mathbb{P}(X_2 = x_2, M_1 = 1, M_2 = 0)$

$$P(X_1 = x_1, X_2 = x_2, M_1 = 0, M_2 = 0)$$

The observed data distribution is all we can identify from samples

•
$$f(\mathbf{X}_{\bar{\mathbf{m}}} = \mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{M} = \mathbf{m}) = \int_{\mathcal{X}_{\mathbf{m}}} f(\mathbf{X} = \mathbf{x}, \mathbf{M} = \mathbf{m}) d\mathbf{x}_{\mathbf{m}}$$

•
$$\mathbb{P}(M_1 = 1, M_2 = 1)$$

•
$$\mathbb{P}(X_1 = x_1, M_1 = 0, M_2 = 1)$$

- $\mathbb{P}(X_2 = x_2, M_1 = 1, M_2 = 0)$
- $\mathbb{P}(X_1 = x_1, X_2 = x_2, M_1 = 0, M_2 = 0)$

General Strategy

$$f(\mathbf{X}_{\bar{\mathbf{m}}} = \mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{M} = \mathbf{m})$$
Identifying assumption
$$\widetilde{f}(\mathbf{X} = \mathbf{x}, \mathbf{M} = \mathbf{m})$$
Sum over \mathbf{m}

$$\widetilde{f}(\mathbf{X} = \mathbf{x})$$

Non-Parametric Saturated Modeling

lf

$$\int \tilde{f}(\mathbf{X} = \mathbf{x}, \mathbf{M} = \mathbf{m}) d\mathbf{x}_{\mathbf{m}} = f(\mathbf{X}_{\bar{\mathbf{m}}} = \mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{M} = \mathbf{m})$$

 Modeling assumption does not impose restrictions on observed data distribution

- Robins (1997) refers to such modeling approach as being non-parametric saturated
- Important property for sensitivity analysis

Non-Parametric Saturated Modeling

lf

$$\int \tilde{f}(\mathbf{X} = \mathbf{x}, \mathbf{M} = \mathbf{m}) d\mathbf{x}_{\mathbf{m}} = f(\mathbf{X}_{\bar{\mathbf{m}}} = \mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{M} = \mathbf{m})$$

 Modeling assumption does not impose restrictions on observed data distribution

- Robins (1997) refers to such modeling approach as being non-parametric saturated
- Important property for sensitivity analysis

Non-Parametric Saturated Modeling

lf

$$\int \tilde{f}(\mathbf{X} = \mathbf{x}, \mathbf{M} = \mathbf{m}) d\mathbf{x}_{\mathbf{m}} = f(\mathbf{X}_{\bar{\mathbf{m}}} = \mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{M} = \mathbf{m}$$

 Modeling assumption does not impose restrictions on observed data distribution

- Robins (1997) refers to such modeling approach as being non-parametric saturated
- Important property for sensitivity analysis

Outline

Inference with Missing Data

Itemwise Missing at Random

Take-Home Message



Itemwise Missing at Random

DEFINITION. The missing data are itemwise missing at random (IMAR) if

$X_j \perp \!\!\!\perp M_j \mid \mathbf{X}_{-j}, \mathbf{M}_{-j}, \text{ for all } j = 1, \dots, p.$

REMARK. X_j and M_j are conditionally independent but not necessarily marginally independent.

Itemwise Missing at Random

DEFINITION. The missing data are itemwise missing at random (IMAR) if

$$X_j \perp \!\!\!\perp M_j \mid \mathbf{X}_{-j}, \mathbf{M}_{-j}, \text{ for all } j = 1, \dots, p.$$

REMARK. X_j and M_j are conditionally independent but not necessarily marginally independent.

IMAR Distribution

THEOREM 1. For each missingness pattern $\mathbf{m} \in \mathcal{M} \subseteq \{0,1\}^p$, given $f(\mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{m}) > 0$, let the function $\eta_{\mathbf{m}} : \mathcal{X}_{\bar{\mathbf{m}}} \mapsto \mathbb{R}$ be defined recursively as

$$\eta_{\mathbf{m}}(\mathbf{x}_{\bar{\mathbf{m}}}) = \log f(\mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{m}) - \log \int_{\mathcal{X}_{\mathbf{m}}} \exp \left\{ \sum_{\mathbf{m}' \prec \mathbf{m}} \eta_{\mathbf{m}'}(\mathbf{x}_{\bar{\mathbf{m}}'}) I(\mathbf{m}' \in \mathcal{M}) \right\} \mu(d\mathbf{x}_{\mathbf{m}}).$$

Then

$$\widetilde{f}(\mathbf{x},\mathbf{m}) = \exp\left\{\sum_{\mathbf{m}' \preceq \mathbf{m}} \eta_{\mathbf{m}'}(\mathbf{x}_{\widetilde{\mathbf{m}}'}) | (\mathbf{m}' \in \mathcal{M}) \right\}$$

satisfies

$$\int_{\mathcal{X}_{\mathbf{m}}} \tilde{f}(\mathbf{x},\mathbf{m})\mu(d\mathbf{x}_{\mathbf{m}}) = f(\mathbf{x}_{\bar{\mathbf{m}}},\mathbf{m}),$$

for all $(\mathbf{x}, \mathbf{m}) \in \mathcal{X} \times \mathcal{M}$.

THEOREM 2. The distribution induced by \tilde{f} encodes the IMAR assumption.

IMAR Distribution

THEOREM 1. For each missingness pattern $\mathbf{m} \in \mathcal{M} \subseteq \{0,1\}^p$, given $f(\mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{m}) > 0$, let the function $\eta_{\mathbf{m}} : \mathcal{X}_{\bar{\mathbf{m}}} \mapsto \mathbb{R}$ be defined recursively as

$$\eta_{\mathbf{m}}(\mathbf{x}_{\bar{\mathbf{m}}}) = \log f(\mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{m}) - \log \int_{\mathcal{X}_{\mathbf{m}}} \exp \left\{ \sum_{\mathbf{m}' \prec \mathbf{m}} \eta_{\mathbf{m}'}(\mathbf{x}_{\bar{\mathbf{m}}'}) I(\mathbf{m}' \in \mathcal{M}) \right\} \mu(d\mathbf{x}_{\mathbf{m}}).$$

Then

$$\tilde{f}(\mathbf{x},\mathbf{m}) = \exp\left\{\sum_{\mathbf{m}' \preceq \mathbf{m}} \eta_{\mathbf{m}'}(\mathbf{x}_{\tilde{\mathbf{m}}'}) l(\mathbf{m}' \in \mathcal{M})\right\}$$

satisfies

$$\int_{\mathcal{X}_{\mathsf{m}}} \tilde{f}(\mathsf{x},\mathsf{m})\mu(d\mathsf{x}_{\mathsf{m}}) = f(\mathsf{x}_{\bar{\mathsf{m}}},\mathsf{m}),$$

for all $(\mathbf{x}, \mathbf{m}) \in \mathcal{X} \times \mathcal{M}$.

THEOREM 2. The distribution induced by \tilde{f} encodes the IMAR assumption.

IMAR Distribution for Categorical Variables

• Log-linear model without interactions involving jointly X_j and M_j

In the case of two variables:

$$\log \mathbb{P}(x_1, x_2, m_1, m_1) = \eta_{x_1 x_2}^{X_1 X_2} + \eta_{x_1 m_2}^{X_1 M_2} + \eta_{x_2 m_1}^{X_2 M_1} + \eta_{m_1 m_2}^{M_1 M_2} + \eta_{x_1}^{X_1} + \eta_{x_2}^{X_2} + \eta_{m_1}^{M_1} + \eta_{m_2}^{M_2} + \eta$$

IMAR Distribution for Categorical Variables

- Log-linear model without interactions involving jointly X_j and M_j
- In the case of two variables:

$$\log \mathbb{P}(x_1, x_2, m_1, m_1) = \eta_{x_1 x_2}^{X_1 X_2} + \eta_{x_1 m_2}^{X_1 M_2} + \eta_{x_2 m_1}^{X_2 M_1} + \eta_{m_1 m_2}^{M_1 M_2} + \eta_{x_1}^{X_1} + \eta_{x_2}^{X_2} + \eta_{m_1}^{M_1} + \eta_{m_2}^{M_2} + \eta$$

The Slovenian Plebiscite Data Revisited

- Slovenians voted for independence from Yugoslavia in a plebiscite in 1991
- ► The Slovenian public opinion survey included these questions:
 - 1. Independence: Are you in favor of Slovenian independence?
 - 2. Secession: Are you in favor of Slovenia's secession from Yugoslavia?

- 3. Attendance: Will you attend the plebiscite?
- Rubin, Stern and Vehovar (1995) analyzed these three questions under MAR
- Plebiscite results give us the proportions of non-attendants and attendants in favor of independence

The Slovenian Plebiscite Data Revisited

- Slovenians voted for independence from Yugoslavia in a plebiscite in 1991
- ► The Slovenian public opinion survey included these questions:
 - 1. Independence: Are you in favor of Slovenian independence?
 - 2. Secession: Are you in favor of Slovenia's secession from Yugoslavia?

- 3. Attendance: Will you attend the plebiscite?
- Rubin, Stern and Vehovar (1995) analyzed these three questions under MAR
- Plebiscite results give us the proportions of non-attendants and attendants in favor of independence

The Slovenian Plebiscite Data Revisited



Figure: Samples from joint posterior distributions of pr(Independence = Yes, Attendance = Yes) and pr(Attendance = No). Pattern mixture model (PMM) under the complete-case missing-variable restriction.

・ロト ・ 四ト ・ ヨト ・ ヨト

ъ

IMAR Distribution for Continuous Variables

With continuous variables

$$f(\mathbf{X}_{\mathbf{\tilde{m}}} = \mathbf{x}_{\mathbf{\tilde{m}}}, \mathbf{M} = \mathbf{m}) = f(\mathbf{X}_{\mathbf{\tilde{m}}} = \mathbf{x}_{\mathbf{\tilde{m}}} | \mathbf{M} = \mathbf{m}) \mathbb{P}(\mathbf{M} = \mathbf{m})$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

► f(X_{m̄} = x_{m̄}|M = m) can be estimated parametrically or non-parametrically

▶ IMAR distribution can be obtained following Theorem 1

IMAR Distribution for Continuous Variables

With continuous variables

$$f(\mathbf{X}_{\bar{\mathbf{m}}} = \mathbf{x}_{\bar{\mathbf{m}}}, \mathbf{M} = \mathbf{m}) = f(\mathbf{X}_{\bar{\mathbf{m}}} = \mathbf{x}_{\bar{\mathbf{m}}} | \mathbf{M} = \mathbf{m}) \mathbb{P}(\mathbf{M} = \mathbf{m})$$

(ロ)、(型)、(E)、(E)、 E) の(の)

▶ IMAR distribution can be obtained following Theorem 1

From the National Health and Nutrition Examination Survey — NHANES (1999–2000 and 2001–2002 cycles):

- X_1 : self-reported height
- ► X₂: actual height measured by survey staff

Informally, the IMAR assumptions are:

- the association between self-reported height and the reporting of this value is explained away by the true height and whether this measurement is taken
- the association between the true height and whether this measurement is taken is explained away by the height that would be self-reported and whether this value is reported

From the National Health and Nutrition Examination Survey — NHANES (1999–2000 and 2001–2002 cycles):

- X_1 : self-reported height
- ► X₂: actual height measured by survey staff

Informally, the IMAR assumptions are:

- the association between self-reported height and the reporting of this value is explained away by the true height and whether this measurement is taken
- the association between the true height and whether this measurement is taken is explained away by the height that would be self-reported and whether this value is reported

From the National Health and Nutrition Examination Survey — NHANES (1999–2000 and 2001–2002 cycles):

- X_1 : self-reported height
- ► X₂: actual height measured by survey staff

Informally, the IMAR assumptions are:

- the association between self-reported height and the reporting of this value is explained away by the true height and whether this measurement is taken
- the association between the true height and whether this measurement is taken is explained away by the height that would be self-reported and whether this value is reported



Figure: Estimated IMAR densities. Left: $\hat{\mathbb{P}}(M_j = 1|x_j)$ for actual height (solid line), and for self-reported height (dashed line).

 $f(x_1, x_2|M_1 = 0, M_2 = 0), f(x_2|M_1 = 1, M_2 = 0)$, and $f(x_1|M_1 = 0, M_2 = 1)$ estimated using survey-weighted kernel density estimators

(日)、

э

Further Uses of IMAR Assumption

- Monotone missingness patterns (dropout/attrition)
- Sensitivity analysis to departures from IMAR assumption
- Use marginal information (e.g. from the Census) to parameterize departures from IMAR

Outline

Inference with Missing Data

Itemwise Missing at Random

Take-Home Message



Take-Home Message

- Itemwise missing at random assumption provides an alternative to MAR assumption
- Allows M_i to depend on X_i marginally
- Can be used with arbitrary missingness patterns and types of variables

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶

Questions?

msadinle@stat.duke.edu