

STAT 403

May 9, 2025

Lecture 15

Imputation

HW5 - TB posted
Project TB posted

Lecture Notes VII – Missing Data and Imputation

Marina Meila
mmp@stat.washington.edu

Department of Statistics
University of Washington

April 2025

The missing data problem

MCAR ✓
MAR ✓
MNAR ?✗

Imputation for MAR

$$P(Y|X, R=1)$$

Imputation for monotone missingness (MAR)



What is the distribution of R ?

- ▶ Let data $\mathcal{D} = \{(x^1, y^1), \dots, (x^n, y^n)\} \cup \{x_{n+1}, \dots, x_{n+m}\}$
 - ▶ complete data $\{(x^1, y^1), \dots, (x^n, y^n)\}$ with $R^{1:n} = 1$
 - ▶ data with missing values $\{x_{n+1}, \dots, x_{n+m}\}$ with $R^{n+1:n+m} = 0$

MCAR Missing Completely At Random $R \perp (X, Y)$

- ▶ no bias if only observed Y values $\{y^1, \dots, y^n\}$ used
- ▶ Can also use imputation as in MAR case

MAR Missing At Random $R \perp Y | X$

- ▶ **bias possible** if only observed Y values $\{y^1, \dots, y^n\}$ are used because $R \not\perp Y$
- ▶ $\{(x^1, y^1, R^1), \dots, (x^n, y^n, R^n)\}$ together with $\{(x^{n+1}, R^{n+1}), \dots, (x^{n+m}, R^{n+m})\}$ has information about $\{y^{n+1}, \dots, y^{n+m}\}$
- ▶ Probability distribution of missing values $Y^{n+1:n+m}$ is

$$Y \sim p(Y | X, R = 0) \quad \text{and} \quad p(Y | X, R = 0) = p(Y | X, R = 1) \quad (1)$$

Imputation for MAR with Monothone Missing data

- ▶ Example Clinical trials with drop-out Each patient i is observed for a $t = 1, 2, \dots, T^i$ with $T^i \leq p$. Hence $y^i = (y_1^i, \dots, y_{T^i}^i)$.
- ▶ T can depend on the previous values $y^{1:T}$ but not on the future values $y^{T+1:p}$.
- ▶ T^i is the missingness variable
- ▶ $x_{1:t}^i$ is the "always observed" data
- ▶ $x_{t+1:p}^i$ is the missing data

Setting ("clinical trial") longitudinal variables data missingness

$[x_1, x_2, \dots, x_t | \dots, x_p]$ $x \in \mathbb{R}^p$

$x_1^i \dots x_t^i | y_{t+1}^i \dots y_p^i$ dropped

$R_1^i \dots R_t^i | R_{t+1}^i \dots R_p^i \Rightarrow t^i$ missing

$\underbrace{R_1^i \dots R_t^i}_{=0}$

MAR $y \perp R | X$

$x_{t+1:p}^i \equiv y_{t+1:p}^i \perp t^i | x_{1:t}^i$

$p(y|X, R=1)$ for obs data levels of missingness

$t = 1, 2, 3, \dots, p-1$

Imputation for MAR with Monothone Missing data

- ▶ Example Clinical trials with drop-out Each patient i is observed for a $t = 1, 2, \dots T^i$ with $T^i \leq p$. Hence $y^i = (y_1^i, \dots y_{T^i}^i)$.
 - ▶ T can depend on the previous values $y^{1:T}$ but not on the future values $y^{T+1:p}$.
 - ▶ T^i is the missingness variable
 - ▶ $x_{1:t}^i$ is the “always observed” data
 - ▶ $x_{t+1:n}^i$ is the missing data

T = stopping time

$$\text{for } t = 1, 2, \dots, P-1$$

want: estimate $p(\underline{x_{t+1}} | x_{1:t}, T=t) = p(\underline{x_{t+1}} | x_{1:t}, T \geq t)$

\uparrow missing want obs

Aq 1

$$\begin{array}{l} i=1 \quad 0 \longrightarrow t^1 = 1 \\ i=2 \quad 0 \quad 1 \quad 0 \qquad t^2 = 3 \end{array}$$

$i = 2$

i=3	1		1
.	0	1	1
.	1	1	0
.	0	0	1
.	0	1	1
.	X ₆	X ₂	5
.	X ₆	X ₂	2
			MISSING

$$P(X_3^6 | X_1^6, X_2^6, T=2) = \left(\frac{1}{2}, \frac{1}{2}\right)$$

$$P(x_2 | x_1, T=1) = \begin{cases} x_1 = 1 & ? \\ x_1 = 0 & \left(\frac{3}{4}, \frac{1}{4}\right) \end{cases}$$

1

$\begin{matrix} 0 \\ \text{---} \\ 0 & 1 & 0 \end{matrix}$ $\xrightarrow{L^2}$
 1

$\begin{matrix} 0 & | & 1 & | & 0 & | & 0 \\ \text{---} & & \text{---} & & \text{---} & & \text{---} \\ 0 & 0 & | & 1 & 1 & | & 0 \end{matrix}$
 $\begin{matrix} 0 & | & 1 \end{matrix}$

use this Eg

$$P(X_4 | X_{1:3}, T=3) = ? \quad \text{no data !!}$$

$\begin{matrix} 0 & 1 & 0 \\ \text{---} \\ 0 & 0 & | & 1 & 0 \end{matrix}$
 no data either

Need additional assumptions

$$\text{Ex. 1 } X_{t+1} | \#\{X_j = 0, j=1:t\}$$

$$\text{Ex. 2. } X_{t+1} | X_t$$

$$\text{Ex. 3. } X_{t+1} | X_t \quad \text{same way for all t}$$

$$P(X_{t+1} | X_t) = P(X_{t+1} | X_t, T) = \begin{pmatrix} 4 & 2 \\ 6 & 6 \\ 1 & 0 \end{pmatrix}$$

Back to original question

$\theta = \text{param of interest}$ Eg

$$\theta = E[X_t]$$

$$\theta = P[\text{longest sequence of 1's}]$$

$$\text{Need } P[X] = P[X_{1:p}]$$

...

$X = X_{1:t} | \underbrace{X_{t+1:p}}_{t \text{ missing}}$

can be imputed

$$P(X) = P(X_{1:t}) P(X_{t+1} | X_{1:t}) P(X_{t+2} | X_{t+1}, X_{1:t}) \cdots P(X_p | X_{1:p-1}) \quad (\text{chain rule})$$

Known ← Alg 1

Need $P[x] = P[x_{1:p}]$

$$x = x_{1:t} \mid \underbrace{x_{t+1:p}}_{\text{is missing}}$$

$$P(x) = P(x_{1:t}) P(x_{t+1} | x_{1:t}) P(x_{t+2} | x_{t+1}, x_{1:t}) \cdots P(x_p | x_{1:p-1})$$

Moreover T also has info about $x_{1:T}$, therefore

$$P(x) = \sum_{t=1}^p p(x_{1:t}) =$$

$$= \sum_{t=1}^p P(T=t) P(x_{1:t} | T=t)$$

$$= \sum_{t=1}^p P(T=t) \underbrace{P(x_{1:t} | T=t)}_{\substack{\text{from} \\ \text{data: sequences of length } t}} \underbrace{P(x_{t+1:p} | x_{1:t})}_{\text{by Alg 1}}$$

from
data

sequences of length t

$$\begin{array}{c} T=1 \mid 0 \\ 0 \mid 0 \\ 1 \\ 0 \ 1 \ 1 \ 1 \ 0 \ 0 \\ 0 \ 0 \ 1 \ 1 \ 0 \\ \parallel 0 \ 1 \end{array}$$

$$P(T) : \frac{1}{6} \ \frac{2}{6} \ \frac{3}{6} \ \frac{4}{6} \ \frac{5}{6} \ \frac{6}{6}$$

$$P(x_1 | T=1) \quad \begin{matrix} 0 & 1 \\ \frac{1}{2} & \frac{1}{2} \end{matrix}$$

$$P(x_1 x_2 | T=2) = \begin{cases} 1 & x_1 x_2 = 0 \\ 0 & \text{otherwise} \end{cases}$$

etc

Finally: ↗ use $P(x)$ to estimate θ
 ↗ use imputation \Rightarrow complete sample \Rightarrow estimate θ