

Lecture 15

Missing data
imputation

Lecture Notes VII – Missing Data and Imputation

Marina Meila
mmp@stat.washington.edu

Department of Statistics
University of Washington

April 2025

Missing data: the problem

- ▶ Data sampled i.i.d. $\sim F$, but some (parts of each sample) are missing
 - ▶ Here we denote data point $= (X^i, Y^i)$ where X^i is always observed, Y^i may be missing or not.
-
- ▶ A new random variable R , $R^i = 0$ when Y^i not observed, $R^i = 1$ when observed

▶ What is the distribution of R ?

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

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

MNAR Missing Not At Random ← nothing to do
(OR no general recipe)

What is the distribution of R ?

$R=1$ if y obs
 0 if y missing

- 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 $R \perp (X, Y) \Rightarrow R \perp Y | X$

- 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 $R^{n+1:n+m}$ is

$$Y \sim p(Y | X, R = 0) \quad \text{and} \quad p(Y | X, R = 0) = \underbrace{p(Y | X, R = 1)}_{\text{estimate from observed } (X^{1:n}, Y^{1:h})} \quad (1)$$

X = always observed \rightarrow "helper"
 \boxed{Y} = may be missing $\rightarrow \theta(Y)$ of interest $\begin{matrix} \nearrow \text{Bias} \\ \searrow \text{Var} \end{matrix}$

Imputation for MAR

$$p(Y|X, R=0) = p(Y|X, R=1)$$

$q(y|x)$

Idea

1. Estimate $p(Y|X, R=1)$ the distribution of $Y|X$ from the observed data
2. "Guess" (Sample) $\tilde{y}^{n+1:n+m} \sim p(Y|X, R=1)$
3. Use $y^{1:n} \cup \tilde{y}^{n+1:n+m}$ to estimate $\hat{\theta}$

treat them as observed

2! sample $\tilde{y}^{i, 1:B} \sim \hat{p}(Y|X=x^i)$

B values for each $i = n+1: n+m$ $\rightarrow Y \in \{0, 1\}$ or discrete classification

$$\begin{aligned} p(Y|X=0) \\ p(Y|X=1) \end{aligned}$$

(Multiple imputation)

3: use $\{y^{1:n}\} \cup \{\tilde{y}^{n+1:n+m, 1:B}\}$
 $\underbrace{\quad}_{x B \text{ times}}$

obs:

imputed:

1.2	1.5	1.5
1.2	1.5	1.1
1.6	1.1	
1.7	1.2	

Ex: $n=3, m=2, B=2$

$\hat{\theta}$

$y^{1:n} \sim s_{1,1} \quad \tilde{y}^{1:n} \sim s_{1,2}$

$y^{n+1:n+m} \sim s_{2,1} \quad \tilde{y}^{n+1:n+m} \sim s_{2,2}$

Estimate $p(Y | X, R = 1)$ from $\mathcal{D}^1 = \{(x^1, y^1, R^1), \dots (x^n, y^n, R^n)\}$

Example $X \in \{1, \dots, K\}$ discrete variable, $Y \in \mathbb{R}$

- For each $k = 1, \dots, K$, estimate $\underline{p_k(Y)} \equiv p(Y | X = k) \leftarrow$ density estimator

1. Estimation methods

Kernel Density Estimation (KDE)

\leftarrow can estimate any $p(\cdot)$
 \leftarrow is non-parametric

$$p_k(y) = \frac{1}{n_k h} \sum_{i: x^i = k} K_h(y^i - y) \quad (2)$$

n_k = number of $x^i = k$, K = kernel function, $K_h(z) = K(z/h)$, h = kernel width (can depend on k)

- Parametric distributions (e.g. $\text{Normal}(\mu_k, \sigma_k^2)$)

step.2 Sampling

Given $p_1, \dots, p_K, x^i = k, i > n$

1. Sample $i' \sim 1 : n_k$ uniformly
2. Sample $\tilde{y}^{i'} \sim N(x^{i'}, h^2)$

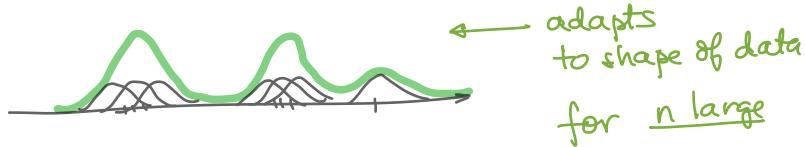
KDE

$$\hat{p}(y) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{y - y_i}{h}\right)$$

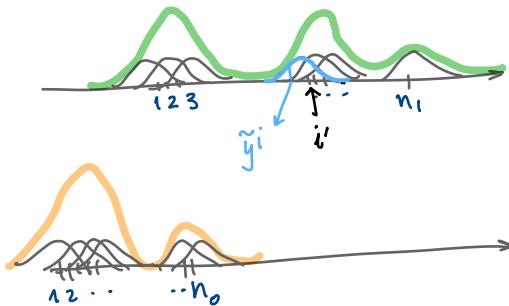
data → place a \sim on each y_i

" $N(0,1)$ " $K(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$ = kernel

h = parameter (kernel width)



$p(y|x=1)$



$p(y|x=0)$

Imputation for y^i

1. $x^i = ?$ $x^i = 1 \rightarrow$ sample from $p(y|x=1)$
2. sample $i' \sim \text{unif}(1 \dots n_1) \Rightarrow i' = 4$
3. sample $\tilde{y}^i \sim N(x^{i'}, h^2) \Rightarrow \underline{\tilde{y}^i}$

Multiple Imputation

Given p_1, \dots, p_k , $x^i = k$, $i = n+1 : n+m$, $B \geq 1$ an integer

for $i = n+1 : n+m$, $b = 1 : B$

1. Sample $i' \sim 1 : n_k$ uniformly
2. Sample $\tilde{y}^{i,b} \sim N(x^{i'}, h^2)$
3. Use data $\underbrace{\{y^{1:n}\}}_{\times B} \cup \{\tilde{y}^{i,b}, \text{ for } i = n+1 : n+m, b = 1 : B\}$