



JOHNS HOPKINS  
BLOOMBERG SCHOOL  
of PUBLIC HEALTH

## Lecture 14

Missing data considerations

R implementation of imputation approaches



# Objectives

- ▶ Throughout the course we have been sub-setting our data such that we are only including rows of data with non-missing outcomes and exposures.
- ▶ Today we will start to explore the possible implications of this practice and think about the underlying assumptions we are making when we do this.
- ▶ Upon completion of this session, you will be able to do the following:
  - ▶ Define mechanisms that generate missing data +
  - ▶ Describe the impact of conducting analyses on complete cases or available data under the different missing data mechanisms ]
  - ▶ Describe imputation procedures to account for missing data ]
  - ▶ Implement several imputation procedures using R mice package ]

multiple  
imputation  
using chained  
equations

# Missing data mechanisms

$$Y \sim f(\mu(x), \sigma^2(x))$$

$$Y^0, R$$

Missing data mechanism	Definition	<u>Ignorable?</u>
<u>Completely at random</u> <div> <div>Covariate dependent missingness</div> <math display="block">R \sim f(x) \neq f(Y^0, Y^m)</math> </div>	<p>Whether or not a value is missing does not depend on observed data OR the missing value we would have observed.</p> <p>- Think a value is missing based a coin toss with some probability of a head that <u>does not depend on anything else</u></p> $Y^0, X, Y^m$	<p>Yes</p> <p>Complete cases represent a random sample of original sample -&gt; analysis of complete case will <u>loss precision</u> but will be unbiased</p>
At random		
Not at random		

# Missing data mechanisms

Missing data mechanism	Definition	Ignorable?
Completely at random	Whether or not a value is missing does not depend on observed data OR the missing value we would have observed.	Yes
<u>At random</u> $f(R   Y^o, X)$	Whether or not a value is missing depends on observed covariates - Think a value is missing is based on a coin toss with probability based on observed characteristics	"Yes" $f(Y   X)$ Complete cases are NOT representative of the original sample. Analysis of complete cases may be biased unless you correctly specify the model, could lose precision
Not at random		

# Missing data mechanisms

Missing data mechanism	Definition	Ignorable?
Completely at random	Whether or not a value is missing does not depend on observed data OR the missing value we would have observed.	Yes
At random	Whether or not a value is missing depends on observed covariates	Yes
Not at random $f(R Y^m)$	Whether or not a value is missing depend on the value of the variable you would observe had it not been missing	<u>No</u> Complete cases are a specific selection of the original sample <u>Bias!</u>

→ sensitivity analysis

# Imputation algorithms

## Single conditional mean imputation

$$Y^0 = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon$$

$$\text{imputed } Y^m = E(\hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p)$$

## Single predicted value imputation

$$Y^0 = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon, \varepsilon \sim N(0, \sigma^2)$$

$$\text{imputed } Y^m = E(\hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p) + \varepsilon_i \quad \left( \begin{array}{l} \text{draw from} \\ N(0, \hat{\sigma}^2) \end{array} \right)$$

Multiple imputation: repeat the single predicted value imputation several times

ID	Y	X	Y <sup>(1)</sup>	Y <sup>(2)</sup>	...	Y <sup>(m)</sup>
1	5	1	5	5		5
2	NA	2	4.9	5.7		6.1 →
3	6	3	6	6		6

## Matching methods

(create subsets of subjects based on  $X$ ):  $Y^0, Y^m$   
 imputed  $Y^m = \text{draw from } Y^0$

replacing missing data under [at random] assumption

# Chained equation approach

mice

- ▶ The idea here is anchored in the desire to estimate the joint distribution of a set of random variables (some values of which are missing). We may be able to derive the exact joint distribution OR we can approximate the joint distribution by deriving the set of full conditional distributions.

E.g.  $Y = (y_1, y_2)$  follows a multivariate normal distribution with mean  $\mu = (\mu_1, \mu_2)$  and variances  $\sigma_1^2$ ,  $\sigma_2^2$  and covariance  $\rho\sigma_1\sigma_2$ .

Then we can write out the two conditional distributions:

- $f(y_1|y_2) \sim N(\mu_1 + \rho\sigma_1\frac{y_2 - \mu_2}{\sigma_2}, \sigma_1^2(1 - \rho^2))$
- $f(y_2|y_1) \sim N(\mu_2 + \rho\sigma_2\frac{y_1 - \mu_1}{\sigma_1}, \sigma_2^2(1 - \rho^2))$

We can use the MCMC algorithm to generate values from each of these two conditional distributions with the end goal of approximating the joint distribution of  $Y$ .



# Chained equation approach

→ most observed → least observed / most missing

Let  $X_1, X_2, \dots, X_p$  be the target imputation variables ordered from most to least observed values.  $Z$  defines a set of prognostic variables that have no missing data. Here I am being generic, the set of target imputation variables may include the outcome variable or not and  $Z$  may include the outcome variable or not, plus any potentially predictive variables for the target imputation variables.

1. Step 1: Setting  $t = 0$ ,  $X_i^{(0)}$  for  $i = 1, \dots, p$  are simulated from

$$f_i(X_i | X_1^{(0)}, X_2^{(0)}, \dots, X_{i-1}^{(0)}; Z, \theta_i)$$

$$X_1^{(0)} = f_1(X_1 | Z, \theta_1) \rightarrow \text{regression parameters}$$

$$X_2^{(0)} = f_2(\overline{X_2} | \underline{X_1^{(0)}}; Z, \theta_2)$$

$$X_3^{(0)} = f_3(X_3 | X_2^{(0)}, X_1^{(0)}; Z, \theta_3)$$

...



## Chained equation approach

2. Step 2: For  $t = \underline{1}$ : obtain simulated values  $X_i^{(1)}$  for  $i = 1, \dots, p$  from

$$X_1^{(1)} = g_1(X_1 | \underbrace{X_2^{(0)}, \dots, X_p^{(0)}}_{\text{previous values}}, Z, \phi_1)$$

$$X_2^{(1)} = g_2(X_2 | \underbrace{X_1^{(1)}, X_3^{(0)}, \dots, X_p^{(0)}}_{\text{previous values}}, Z, \phi_2)$$

through

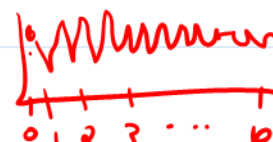
$$X_p^{(1)} = g_p(X_p | X_1^{(1)}, X_2^{(1)}, \dots, X_{p-1}^{(1)}, Z, \phi_p)$$

Then repeat this process for  $t = 2, \dots, b$ .

1 chain

Take as imputed  
values  $X_i^{(b)}$

m chains  $\rightarrow b?$

$X_i$  

Now, let's implement some of these approaches!

