# How black-box use of imputation can cause bias

Nicole Erler

Erasmus Medical Center, Rotterdam

**Erasmus MC**
University Medical Center Rotterdam

# Handling Missing Values is Easy!

**Functions automatically exclude missing values:**

```
## [...]
## Residual standard error:  2.305 on 69 degrees of freedom
##    (25 observations deleted due to missingness)
## Multiple R-squared:  0.09255, Adjusted R-squared:  0.02679
## F-statistic:  1.407 on 5 and 69 DF, p-value:  0.2325
```

## Handling Missing Values is Easy!

**Functions automatically exclude missing values:**

```
## [...]
## Residual standard error:  2.305 on 69 degrees of freedom
##    (25 observations deleted due to missingness)
## Multiple R-squared:  0.09255, Adjusted R-squared:  0.02679
## F-statistic:  1.407 on 5 and 69 DF, p-value:  0.2325
```

**Imputation is super easy:**

```
library("mice")
imp <- mice(mydata)
```

**However ...**

## Handling Missing Values Correctly is Not So Easy!

(Imputation) methods makes certain **assumptions**, e.g.:

- **missingness** is M(C)AR

# Handling Missing Values Correctly is Not So Easy!

(Imputation) methods makes certain **assumptions**, e.g.:

- **missingness** is M(C)AR
- the incomplete variable has a certain conditional **distribution** (e.g. normal)

## Handling Missing Values Correctly is Not So Easy!

(Imputation) methods makes certain **assumptions**, e.g.:

- **missingness** is M(C)AR
- the incomplete variable has a certain conditional **distribution** (e.g. normal)
- all associations are **linear**
  - no interactions
  - no non-linear effects
  - no transformations

## Handling Missing Values Correctly is Not So Easy!

(Imputation) methods makes certain **assumptions**, e.g.:

- **missingness** is M(C)AR
- the incomplete variable has a certain conditional **distribution**
  (e.g. normal)
- all associations are **linear**
  - no interactions
  - no non-linear effects
  - no transformations
- **compatibility** of the imputation models
- **congeniality** (compatibility between analysis and imputation models)

# Handling Missing Values Correctly is Not So Easy!

(Imputation) methods makes certain **assumptions**, e.g.:

- **missingness** is M(C)AR
- the incomplete variable has a certain conditional **distribution** (e.g. normal)
- all associations are **linear**
  - no interactions
  - no non-linear effects
  - no transformations
- **compatibility** of the imputation models
- **congeniality** (compatibility between analysis and imputation models)

## violation ➡ bias

# Literature: mis-specification in Multiple Imputation

Several authors have

- investigated robustness to mis-specification (of distribution)
    - in MI using FCS / MICE
    - in joint model MI
- and/or proposed to use
    - Tukey's gh distribution
    - Fleishman polynomials
    - GAMs (in FCS)
    - Doubly-robust weighted estimating equations (instead of MI)

# Fully Bayesian Analysis & Imputation

## Joint distribution

$$\underbrace{p(y \mid X, b, \theta)}_{\substack{\text{analysis} \\ \text{model}}} \underbrace{p(X \mid \theta)}_{\substack{\text{imputation} \\ \text{part}}} \underbrace{p(b \mid \theta)}_{\substack{\text{random} \\ \text{effects}}} \underbrace{p(\theta)}_{\text{priors}}$$

## Fully Bayesian Analysis & Imputation

### Joint distribution

$$\underbrace{p(y \mid X, b, \theta)}_{\substack{\text{analysis} \\ \text{model}}} \underbrace{p(X \mid \theta)}_{\substack{\text{imputation} \\ \text{part}}} \underbrace{p(b \mid \theta)}_{\substack{\text{random} \\ \text{effects}}} \underbrace{p(\theta)}_{\text{priors}}$$

### Imputation part

$$
p(\overbrace{x_1, \ldots, x_p, X_{compl.}}^{X} \mid \theta) = p(x_1 \mid X_{compl.}, \theta) \\
p(x_2 \mid X_{compl.}, x_1, \theta) \\
p(x_3 \mid X_{compl.}, x_1, x_2, \theta) \\
\ldots
$$

## Fully Bayesian Analysis & Imputation

### Joint distribution

$$\underbrace{p(y \mid X, b, \theta)}_{\substack{\text{analysis} \\ \text{model}}} \underbrace{p(X \mid \theta)}_{\substack{\text{imputation} \\ \text{part}}} \underbrace{p(b \mid \theta)}_{\substack{\text{random} \\ \text{effects}}} \underbrace{p(\theta)}_{\text{priors}}$$

### Imputation part

$$p(\overbrace{x_1, \ldots, x_p, X_{compl.}}^{X} \mid \theta) = p(x_1 \mid X_{compl.}, \theta)$$
$$p(x_2 \mid X_{compl.}, x_1, \theta)$$
$$p(x_3 \mid X_{compl.}, x_1, x_2, \theta)$$
$$\ldots$$

### Software

Implemented in the **R** package **JointAI**

### Imputation in MICE

$$p(x_1 \mid y, X_{compl.}, x_2, x_3, x_4, \ldots, \theta)$$
$$p(x_2 \mid y, X_{compl.}, x_1, x_3, x_4, \ldots, \theta)$$
$$p(x_3 \mid y, X_{compl.}, x_1, x_2, x_4, \ldots, \theta)$$
$$\ldots$$

### Imputation in JointAI

$$p(y \mid X_{compl.}, x_1, x_2, x_3, \ldots, \theta)$$
$$p(x_1 \mid X_{compl.}, \theta)$$
$$p(x_2 \mid X_{compl.}, x_1, \theta)$$
$$p(x_3 \mid X_{compl.}, x_1, x_2, \theta)$$
$$\ldots$$

## Imputation in MICE

$$p(x_1 \mid y, X_{compl.}, x_2, x_3, x_4, \ldots, \theta)$$
$$p(x_2 \mid y, X_{compl.}, x_1, x_3, x_4, \ldots, \theta)$$
$$p(x_3 \mid y, X_{compl.}, x_1, x_2, x_4, \ldots, \theta)$$
$$\ldots$$

## Imputation in JointAI

$$p(y \mid X_{compl.}, x_1, x_2, x_3, \ldots, \theta)$$
$$p(x_1 \mid X_{compl.}, \theta)$$
$$p(x_2 \mid X_{compl.}, x_1, \theta)$$
$$p(x_3 \mid X_{compl.}, x_1, x_2, \theta)$$
$$\ldots$$

No issues with

- complex outcomes, e.g.:
  - multi-level
  - survival
- congeniality
- compatibility

## Imputation in MICE

$$p(x_1 \mid y, X_{compl.}, x_2, x_3, x_4, \ldots, \theta)$$
$$p(x_2 \mid y, X_{compl.}, x_1, x_3, x_4, \ldots, \theta)$$
$$p(x_3 \mid y, X_{compl.}, x_1, x_2, x_4, \ldots, \theta)$$
$$\ldots$$

## Imputation in JointAI

$$p(y \mid X_{compl.}, x_1, x_2, x_3, \ldots, \theta)$$
$$p(x_1 \mid X_{compl.}, \theta)$$
$$p(x_2 \mid X_{compl.}, x_1, \theta)$$
$$p(x_3 \mid X_{compl.}, x_1, x_2, \theta)$$
$$\ldots$$

Potential mis-specification of

- association structure
- conditional distribution
- M(C)AR

Analysis model: $y \sim \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots$

Quadratic association between covariates: $x_1 \sim \alpha_0 + \alpha_1 x_2 + \alpha_2 x_2^2 + \ldots$
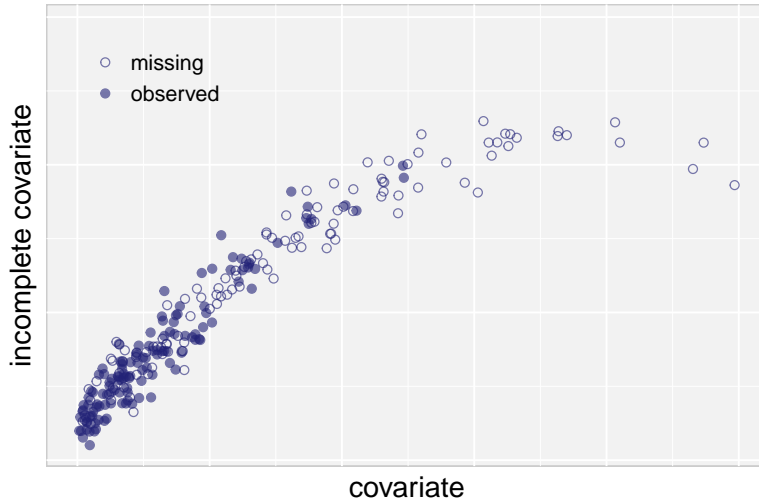
Analysis model: $y \sim \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots$

Quadratic association between covariates: $x_1 \sim \alpha_0 + \alpha_1 x_2 + \alpha_2 x_2^2 + \ldots$

Analysis model: $y \sim \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$

Log-association between covariates: $x_1 \sim \alpha_0 + \alpha_1 \log(x_2) + \dots$

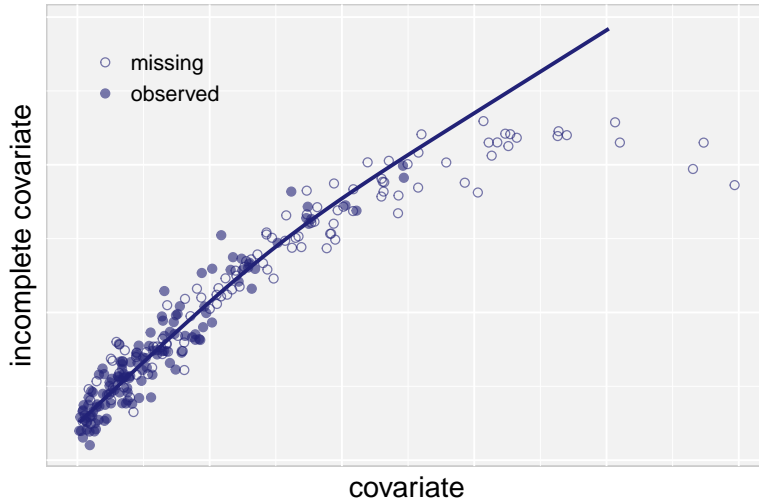Analysis model: $y \sim \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots$

Gamma-distributed covariate: $x_1 \mid x_2, x_3, \ldots \sim Ga()$

## Flexible Bayesian Models

**We need more flexible imputation models!**

**Ideally:** models that fit (almost) any distribution / association structure.

**We need more flexible imputation models!**

**Ideally:** models that fit (almost) any distribution / association structure.

**Ideas:**

- flexible **association** structure: **penalized splines**
- flexible residual **distribution**: **mixture of Polya-Trees**

Instead of $\quad \beta_1 x_2 \quad$ we use $\quad \sum_{\ell=1}^{d} \beta_\ell B_\ell(x_2):$

Instead of $\quad \beta_1 x_2 \quad$ we use $\quad \sum_{\ell=1}^{d} \beta_\ell B_\ell(x_2):$

Instead of $\quad \beta_1 x_2 \quad$ we use $\quad \displaystyle\sum_{\ell=1}^{d} \beta_\ell B_\ell(x_2):$

Analysis model: $y \sim \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots$

Quadratic association between covariates: $x_1 \sim \alpha_0 + \textcolor{red}{\alpha_1 x_2 + \cancel{\alpha_2 x_2^2}} + \ldots$

# Simulation: Bayesian P-Splines

Analysis model: $y \sim \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots$

Quadratic association between covariates: $x_1 \sim \alpha_0 + \alpha_1 x_2 + \widetilde{\alpha_2 x_2^2} + \ldots$

Potential Issue:

Potential Issue:



covariate

# Mixture of Polya Trees



Beta($a_0$, $a_0$)

# Mixture of Polya Trees

# Mixture of Polya Trees

# Mixture of Polya Trees

# Mixture of Polya Trees

## Practical Issues & Ideas

**Potential / probable issues for practice:**

- flexible fit needs observed data everywhere
- computational time

## Practical Issues & Ideas

**Potential / probable issues for practice:**

- flexible fit needs observed data everywhere
- computational time

**Ideas:**

- check first if simple model fits, e.g.
  posterior predictive checks
    - $\chi^2$ type of tests
    - Kolmogorov-Smirnoff test?
    - discordance tests?

## Practical Issues & Ideas

**Potential / probable issues for practice:**

- flexible fit needs observed data everywhere
- computational time

**Ideas:**

- check first if simple model fits, e.g.
  posterior predictive checks
    - $\chi^2$ type of tests
    - Kolmogorov-Smirnoff test?
    - discordance tests?
- feasibility checks before running the complex model

## Take home message

- **assumptions** of imputation models can easily be **violated** ➡ **bias**
- **more flexible** imputation models are **needed**
- semi- / **non-parametric (Bayesian) methods** can offer a **solution**
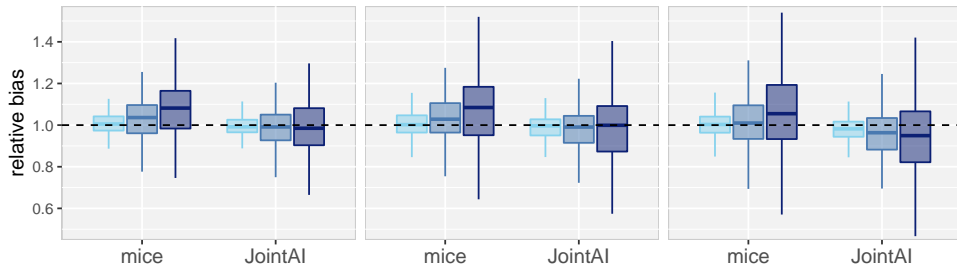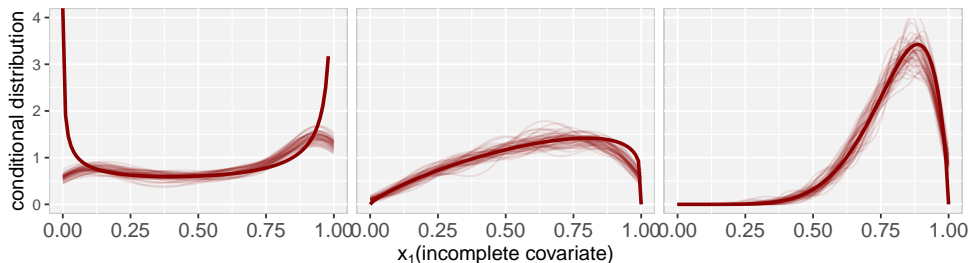- more **flexibility** ➡ more **complexity** ➡ need for **guidance** tools

**Thank you for your attention.**

✉ **n.erler@erasmusmc.nl**
𝕏 **N_Erler**
○ **NErler**
🌐 **www.nerler.com**

Analysis model: $y \sim \beta_0 + \beta_1 x_1 + \beta_2 + \ldots$

Beta-distributed covariate: $x_1 \mid x_2, x_3, \ldots \sim Be()$

Analysis model: $y \sim \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots$

Exclusion of important predictor: $x_1 \sim \alpha_0 + \alpha_1 x_2 + \cancel{\alpha_2 x_3} + \ldots$

Analysis model: $y \sim \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots$

Interaction between covariates: $x_1 \sim \alpha_0 + \alpha_1 x_2 + \alpha_2 x_3 + \cancel{\alpha_3 x_2 x_3} + \ldots$