

EP16: Missing Values in Clinical Research: Multiple Imputation

14. Strategies for using MICE

Nicole Erler

Department of Biostatistics, Erasmus Medical Center

✉ n.erler@erasmusmc.nl

Number of Imputations

In **early publications** on multiple imputation:

- ▶ 3 – 5 imputations are sufficient
- ▶ still is a common assumption in practice (Rubin 2004)

Reasoning behind using a small number of imputations:

- ▶ **storage of imputed data was “expensive”** (no longer the case)
- ▶ larger number of imputations would only have little advantage (Schafer 1997)

Number of Imputations

More **recent work** from various authors considers

- ▶ the **efficiency** of the pooled estimates
- ▶ **reproducibility** of the results,
- ▶ **statistical** power of tests or
- ▶ the **width of the resulting confidence intervals** compared to the width of the true confidence intervals

(White, Royston, and Wood 2011; Van Buuren 2012; Graham, Olchowski, and Gilreath 2007)

Number of Imputations

A **rule of thumb** (White, Royston, and Wood 2011):

The number of imputed datasets should be similar to the **percentage of incomplete cases**.

Alternative:

The percentage of complete cases depends on the size of the dataset
➔ **average percentage of missing values** per variable (Van Buuren 2012)

Generally:

- ▶ using **more imputed datasets should be preferred**
- ▶ especially in settings with acceptable computational burden

Even though results are unlikely to change with a larger number of imputations, it can increase the efficiency and reproducibility of the results.

What to do with large datasets?

General advice: **Include as much information as possible** in the imputation.

Using a large number of predictor variables

- ▶ makes the **MAR assumption more plausible** (⇒ reduces bias due to MNAR missingness)
- ▶ can **reduce uncertainty** about the missing values

What to do with large datasets?

General advice: **Include as much information as possible** in the imputation.

Using a large number of predictor variables

- ▶ makes the **MAR assumption more plausible** (⇒ reduces bias due to MNAR missingness)
- ▶ can **reduce uncertainty** about the missing values

This **works in small or medium sized datasets** (20 – 30 separate variables, i.e. without interactions, variables derived from others, ...).

In large datasets (contain hundreds or thousands of variables) this is **not feasible**. (Van Buuren 2012)

What to do with large datasets?

For large datasets a possible strategy is to

- ▶ Include all **variables used in the analysis model(s)** (including the outcome!)
- ▶ Include auxiliary variables
 - ▶ if they are **strong predictors of missingness**, or
 - ▶ if they have **strong associations with the incomplete variables**
 - ▶ **only if they do not have too many missing values** themselves
 - ▶ only in those imputation models for which they are **relevant**
- ▶ Use **summary scores** when multiple items referring to the same concept.

How much missing is too much?

There is **no clear cut-off** for the proportion of missing values that can be handled adequately using MICE (or any other imputation method).

The amount of missingness that can be handled **depends on the information that is available** to impute it.

- ▶ Are there **strong predictor variables** available & observed?
- ▶ Are there **sufficient observed cases** to get reliable estimates for the predictive distribution?

How much missing is too much?

There is **no clear cut-off** for the proportion of missing values that can be handled adequately using MICE (or any other imputation method).

The amount of missingness that can be handled **depends on the information that is available** to impute it.

- ▶ Are there **strong predictor variables** available & observed?
- ▶ Are there **sufficient observed cases** to get reliable estimates for the predictive distribution?

Example:

- ▶ In a set of $N = 50$ cases, 50% missing values leaves 25 cases to estimate the parameters of the predictive distribution.
- ▶ In a large set of $N = 5000$ subjects, 50% missing cases leaves 2500 observed cases to estimate parameters.

Imputation of Outcomes

Usually, **missing outcome values are not imputed.**

Why?

When there are no auxiliary variables, imputation and analysis model are equal.

- ▶ Parameters of the imputation model are estimated on observed cases of the outcome.
- ▶ Imputed values will fit the assumed model perfectly.
- ▶ Including imputed cases in the analysis does not add any information.

Imputation of Outcomes

Exception:

- ▶ When very strong auxiliary variables are available.
- ▶ Outcomes may be imputed when one imputation is performed for several analysis models, because not imputing the outcome(s) would mean
 - ▶ excluding cases with missing outcome(s) from the imputation, or
 - ▶ excluding the outcome variable(s) as predictor(s).

Notes of Caution & Things to Keep in Mind

Multiple imputation is **not a quick and easy solution for missing data**.

It requires **care and knowledge** about

- ▶ the **data** to be imputed (and the context of the data),
- ▶ the statistical **method** used for imputation, and
- ▶ the **software** implementation used.

Moreover

- ▶ **Never accept default settings of software blindly.**
- ▶ **Question the plausibility of the MAR assumption.**
If it is doubtful, use sensitivity analysis.

Notes of Caution & Things to Keep in Mind

- ▶ **Use as much information as possible**
 - ▶ include all covariates **and the outcome**
 - ▶ use auxiliary information
 - ▶ use the most detailed version of variables if possible
- ▶ **Avoid feedback** from derived variables to their originals.
- ▶ Think carefully how to handle variables that are derived from other variables.
- ▶ Consider the impact the **visit sequence** may have.
- ▶ **Imputation models must fit the data**
(correct assumption of error distribution and functional forms and possible interactions of predictor variables).

Notes of Caution & Things to Keep in Mind

- ▶ Choose an appropriate **number of imputations**.
- ▶ Make sure the imputation algorithm has **converged**.
- ▶ Use **common sense** when evaluating if the imputed values are plausible.
- ▶ Be aware of the assumptions of your analysis model
 - ▶ non-linear effects
 - ▶ interactions
 - ▶ complex outcomes

Tips & Tricks

In complex settings, variables may need to be **re-calculated** or **re-coded** after imputation:


- ▶ Use `complete()` to convert the imputed data from a `mids` object to a `data.frame`.
- ▶ Perform the necessary calculations.
- ▶ Convert the changed `data.frame` back to a `mids` object using the functions such as `mice::as.mids()`, `miceadds::datalist2mids()`, `mitools::imputationList()`, ...

Tips & Tricks

In complex settings, variables may need to be **re-calculated** or **re-coded** after imputation:

- ▶ Use `complete()` to convert the imputed data from a `mids` object to a `data.frame`.
- ▶ Perform the necessary calculations.
- ▶ Convert the changed `data.frame` back to a `mids` object using the functions such as `mice::as.mids()`, `miceadds::datalist2mids()`, `mitools::imputationList()`, ...

Not just in imputation: Set a **seed value** to create reproducible results.

- ▶ in : `set.seed()`
- ▶ in `mice()`: argument `seed`

Imputation Methods

We have focussed on a few imputation methods that cover the most common types of data, but there are many more methods implemented.

Imputation methods implemented in the **mice** package:

<code>mice.impute.2l.bin</code>	<code>mice.impute.logreg</code>	<code>mice.impute.panImpute</code>
<code>mice.impute.2l.lmer</code>	<code>mice.impute.logreg.boot</code>	<code>mice.impute.passive</code>
<code>mice.impute.2l.norm</code>	<code>mice.impute.mean</code>	<code>mice.impute.pmm</code>
<code>mice.impute.2l.pan</code>	<code>mice.impute.midastouch</code>	<code>mice.impute.polr</code>
<code>mice.impute.2lonly.mean</code>	<code>mice.impute.mnar.logreg</code>	<code>mice.impute.polyreg</code>
<code>mice.impute.2lonly.norm</code>	<code>mice.impute.mnar.norm</code>	<code>mice.impute.quadratic</code>
<code>mice.impute.2lonly.pmm</code>	<code>mice.impute.norm</code>	<code>mice.impute.rf</code>
<code>mice.impute.cart</code>	<code>mice.impute.norm.boot</code>	<code>mice.impute.ri</code>
<code>mice.impute.jomolImpute</code>	<code>mice.impute.norm.nob</code>	<code>mice.impute.sample</code>
<code>mice.impute.lda</code>	<code>mice.impute.norm.predict</code>	NA

Note: That a method is implemented does not mean you need to / should use it.

Imputation Methods

Imputation methods implemented in the **miceadds** package:

mice.impute.2l.binary	mice.impute.hotDeck
mice.impute.2l.contextual.norm	mice.impute.lm
mice.impute.2l.contextual.pmm	mice.impute.lm_fun
mice.impute.2l.continuous	mice.impute.lqs
mice.impute.2l.groupmean	mice.impute.ml.lmer
mice.impute.2l.groupmean.elim	mice.impute.plausible.values
mice.impute.2l.latentgroupmean.mcmc	mice.impute.pls
mice.impute.2l.latentgroupmean.ml	mice.impute.pmm3
mice.impute.2l.plausible.values	mice.impute.pmm4
mice.impute.2l.pls	mice.impute.pmm5
mice.impute.2l.pls2	mice.impute.pmm6
mice.impute.2l.pmm	mice.impute.rlm
mice.impute.2lonly.function	mice.impute.smcfcfs
mice.impute.2lonly.norm2	mice.impute.tricube.pmm
mice.impute.2lonly.pmm2	mice.impute.tricube.pmm2
mice.impute.bygroup	mice.impute.weighted.norm
mice.impute.grouped	mice.impute.weighted.pmm

Imputation Methods


Imputation methods implemented in the **micemd** package:

mice.impute.2l.2stage.bin	mice.impute.2l.glm.bin
mice.impute.2l.2stage.norm	mice.impute.2l.glm.norm
mice.impute.2l.2stage.pmm	mice.impute.2l.glm.pois
mice.impute.2l.2stage.pois	mice.impute.2l.jomo

Other R Packages for Imputation

CRAN Task View on Missing Data:

<https://cran.r-project.org/web/views/MissingData.html>

- ▶ gives an overview on the available  packages for missing data / imputation
- ▶ good point to start when searching for a package with a particular functionality


Other R Packages for Imputation


Currently, there are **414 packages** available on CRAN that use the word **“missing”, “impute”, “imputation” or “incomplete”** in either the title or description.

Not all of these packages perform imputation or are useful for our purposes, but even if we excluded those packages, the number of useful packages for dealing with missing data would still be too large to mention them all.

➔ **The mice package is often a good option, but certainly not the only option to perform imputation!**

Imputation in other Software

In this second half of the course, we have focused on (multiple) imputation using .

Naturally,  is not the only statistical software that can perform multiple imputation.

- ▶ **Stata, SAS and MPLUS** provide packages/functions to perform multiple imputation and pool the results.
- ▶ There are macros and additional packages available, e.g., **smcfcs** is implemented for **Stata** as well
- ▶ **SPSS** provides some functionality to perform MI

Other Approaches to Handle Missing Values

Finally, we should not forget that **MICE is not the only method to handle missing values.**

Besides MICE, **multiple imputation** can be performed in a **joint model approach** (as for instance implemented in the R package **jomo**).

Furthermore,

- ▶ **direct likelihood methods,**
- ▶ **fully Bayesian methods** (as implemented in **JointAI**), or
- ▶ **weighted estimating equations**

are valid approaches and may **in certain settings be superior.**

References I

Graham, John W, Allison E Olchowski, and Tamika D Gilreath. 2007. "How Many Imputations Are Really Needed? Some Practical Clarifications of Multiple Imputation Theory." *Prevention Science* 8 (3): 206–13.

Rubin, Donald B. 2004. "The Design of a General and Flexible System for Handling Nonresponse in Sample Surveys." *The American Statistician* 58 (4): 298–302. <https://doi.org/10.1198/000313004X6355>.

Schafer, Joseph L. 1997. *Analysis of Incomplete Multivariate Data*. CRC press.

Van Buuren, Stef. 2012. *Flexible Imputation of Missing Data*. Chapman & Hall/Crc Interdisciplinary Statistics. Taylor & Francis.
<https://stefvanbuuren.name/fimd/>.

References II

White, Ian R, Patrick Royston, and Angela M Wood. 2011. "Multiple Imputation Using Chained Equations: Issues and Guidance for Practice." *Statistics in Medicine* 30 (4): 377–99.