



Imputation of Incomplete Multilevel Data with mice

Hanne I. Oberman

Utrecht University

Johanna Munoz

Julius Center for Health Sciences and Primary Care,
University Medical Center Utrecht, Utrecht University,
Utrecht, The Netherlands

Thomas P. A. Debray

Julius Center for Health Sciences and Primary Care

Gerko Vink

Utrecht University

Valentijn M. T. de Jong

Julius Center for Health Sciences and Primary Care,
University Medical Center Utrecht, Utrecht University,
Utrecht, The Netherlands

Abstract

This tutorial illustrates the imputation of incomplete multilevel data with the R package **mice**. Footnotes in the current version show work in progress/under construction. The last section is not part of the manuscript, but purely for reminders. See also all of the TODOs that need to be worked out. We aim to submit at JSS, so there is no word count limit (“There is no page limit, nor a limit on the number of figures or tables”). [Just adding some text to get a better guess of what the actual abstract will look like: Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.]

Keywords: missing data, multilevel, clustering, **mice**, R.

1. Introduction

Many datasets include individuals from multiple settings, geographic regions, or even different

studies. In the simplest case, individuals (e.g., students) are nested within so-called clusters (e.g., school classes). More complex clustered structures may occur when there are multiple hierarchical levels (e.g., patients within hospitals within regions or countries), or when the clustering is non-nested (e.g., electronic health record data from diverse settings and populations within large databases). In general, individuals from the same cluster tend to be more similar than individuals from other clusters. In statistical terms, this implies that observations from the same cluster are correlated. If this correlation is left unaddressed, estimates of p values, confidence intervals even model parameters are prone to bias (Localio, Berlin, Ten Have, and Kimmel 2001). [TODO: make a link to imputation methods, which require adequate handling and propagation of variance; we are not recommending the adoption of multilevel models for data analysis here, but rather for imputation. VMTdJ: “I would only make that link after introducing missing data, and missing data is not mentioned yet, so I would do that in the next section”] Statistical methods for clustered data typically adopt hierarchical models that explicitly describe the grouping of observations. These models are also known as ‘multilevel models’, ‘hierarchical models’, ‘mixed effect models’, ‘random effect models’, and in the context of time-to-event data as ‘frailty models’. Table ?? provides an overview of some key concepts in multilevel modeling.

1.1. Missingness in multilevel data

Like any other dataset, clustered datasets are prone to missing data. Several strategies can be used to handle missing data, including complete case analysis and imputation. We focus on the latter approach and discuss statistical methods for replacing the missing data with one or more plausible values. Afterwards, the completed data can be analyzed as if they were completely observed. In contrast to single imputation (where missing data are only replaced once), multiple imputation allows to preserve uncertainty due to missingness and is therefore recommended (c.f. Rubin 1976).

When clustered datasets are affected by missing values, we can distinguish between two types of missing data: sporadic missingness and systematic missingness (Resche-Rigon, White, Bartlett, Peters, and Thompson 2013). Sporadic missingness arises when variables are missing for some but not all of the units in a cluster (Van Buuren 2018; Jolani 2018). For example, it is possible that test results are missing for several students in one or more classes. [TODO: Provide an example for one of the case studies below.] When all observations are missing within one or more clusters, data are systematically missing. [TODO: Refer to Figure 1 and put interpretation in the figure caption. VMTdJ: “Mention in text that this is sporadic missingness?”]

2. Introduction

Many datasets include individuals that are clustered together, for example in geographic regions, or even different studies. In the simplest case, individuals (e.g., students) are nested within a single cluster (e.g., school classes). More complex clustered structures may occur when there are multiple hierarchical levels (e.g., students in different schools or patients within hospitals within regions across countries), or when the clustering is non-nested (e.g., electronic health record data from diverse settings and populations within large databases). With clustered data we generally assume that individuals from the same cluster tend to be more

similar than individuals from other clusters. In statistical terms, this implies that observations from the same cluster are not independent and may in fact be correlated. If this correlation is left unaddressed, estimates of p values, confidence intervals even model parameters are prone to bias (Localio et al. 2001). Statistical methods for clustered data typically adopt hierarchical models that explicitly describe the grouping of observations. These models are also known as ‘multilevel models’, ‘hierarchical models’, ‘mixed effect models’, ‘random effect models’, and in the context of time-to-event data as ‘frailty models’. Table ?? provides an overview of some key concepts in multilevel modeling.

Table 1: Concepts in multilevel methods

| Concept | Details |
|----------------------|---|
| Sample unit | Units of the population from which measurements are taken in a sample. |
| Hierarchical levels | Data are grouped into clusters at different levels, observations belonging to the same cluster are expected to share certain characteristics. |
| Fixed effect | Effects that are constant across all sample units, e.g. something that researchers control for and can repeat, such as the administration of a drug. |
| Random effect | Effects that are a source of random variation in the data, and whose levels are not fully sampled. e.g. individuals are drawn from a population of hospitals, here it is not possible to sample all hospitals but drug effects could vary between hospitals. |
| Mixed effect | Includes fixed and random effects, e.g. the fixed effect would be the treatment effect of a drug and the random effect would be the ID of the hospital where the patient is treated. Multilevel models typically accommodate for variability by including a separate group mean for each cluster e.g random intercept on hospitals. In addition to random intercepts, multilevel models can also include random coefficients and heterogeneous residual error variances across clusters (see e.g. Gelman and Hill 2006, Hox, Moerbeek, and van de Schoot (2017) and de Jong, Moons, Eijkemans, Riley, and Debray (2021)). |
| ICC | The variability due to clustering is often measured by means of the intraclass coefficient (ICC). The ICC can be seen as the percentage of variance that can be attributed to the cluster-level, where a high ICC would indicate that a lot of variability is due to the cluster structure. |
| Stratified intercept | |

2.1. Missingness in multilevel data

As with any other dataset, clustered datasets may be impacted by missingness in much the same way. Several strategies can be used to handle missing data, including complete case analysis and imputation. We focus on the latter approach and discuss statistical methods for replacing the missing data with one or more plausible values. Imputation separates the

| | cluster | X_1 | X_2 | X_3 | ... | X_p |
|-----|---------|-------|-------|-------|-----|-------|
| 1 | 1 | | | NA | | |
| 2 | 1 | | | | | |
| 3 | 2 | | NA | | | |
| 4 | 2 | | NA | NA | | |
| 5 | 3 | | | | | |
| ... | | | | | | |
| n | N | | | | | |

Figure 1: Missingness in multilevel data

Table 2: Concepts in missing data methods

| Concept | Details |
|---------|---|
| MCAR | Missing Completely At Random, where the probability to be missing is equal across all data entries |
| MAR | Missing At Random, where the probability to be missing depends on observed information |
| MNAR | Missing Not At Random (MNAR), where the probability to be missing depends on unrecorded information, making the missingness non-ignorable [rubi76; meng94]. [TODO: add congeniality, but maybe in-text?] |

missing data problem from the analysis and the completed data can be analyzed as if it were completely observed. It is generally recommended to impute the missing values more than once to preserve uncertainty due to missingness and to allow for valid inferences (c.f. Rubin 1976).

With incomplete clustered datasets we can distinguish between two types of missing data: sporadic missingness and systematic missingness (Resche-Rigon et al. 2013). Sporadic missingness arises when variables are missing for some but not all of the units in a cluster (Van Buuren 2018; Jolani 2018). For example, it is possible that test results are missing for several students in one or more classes. When all observations are missing within one or more clusters, data are said to be systematically missing.

Imputation of missing data requires consideration of the mechanism behind the missingness. Rubin proposed to distinguish between data that are missing completely at random (MCAR), data that are missing at random (MAR) and data that are missing not at random (MNAR; see Table 2). For each of these three missingness generating mechanisms, different imputation strategies are warranted Yucel (2008) and Hox, van Buuren, and Jolani (2015). We here consider the general case that data are MAR, and expand on certain MNAR situations.

The R package **mice** has become the de-facto standard for imputation by chained equations, which iteratively solves the missingness on a variable-by-variable basis. **mice** is known to yield valid inferences under many different missing data circumstances (Van Buuren 2018). However, commonly used imputation methods were not designed for use in clustered data and usually generate observations that are independent whereas multilevel data are dependent.

For this reason, we discuss how the R package **mice** can be used to impute multilevel data.

[TODO: clarify why clustering is relevant during imputation, and why this exposes the need for specialized imputation methods and more attention during their implementation (“thou shall not simply run `mice()` on any incomplete dataset”).] [TODO: Add that the more the random effects are of interest, the more you need multilevel imputation models.] [TODO: Add an overview of all possible predictor matrix values in manuscript or **ggmice** legend.]

2.2. Aim of this paper

This paper serves as a tutorial for imputing incomplete multilevel data with **mice** in R. We provide practical guidelines and code snippets for different missing data situations, including non-ignorable mechanisms. For reasons of brevity, we focus on multilevel imputation by chained equations with **mice** exclusively; other imputation methods and packages (e.g., **jomo** and **mdmb**) are outside the scope of this tutorial. [TODO: change this sentence because we do not assume familiarity: “Assumed knowledge includes basic familiarity with multilevel imputation (see e.g. [Audigier, White, Jolani, Debray, Quartagno, Carpenter, van Buuren, and Resche-Rigon 2018](#), and [Grund, Lüdtke, and Robitzsch \(2018\)](#)) and the **lme4** notation for multilevel models (see Table ??)”.]

We illustrate imputation of incomplete multilevel data using three case studies:

- **popmis** from the **mice** package (simulated data on perceived popularity, $n = 2,000$ pupils across $N = 100$ schools with data that are MAR, [van Buuren and Groothuis-Oudshoorn 2021](#));
- **impact** from the **metamisc** package (empirical data on traumatic brain injuries, $n = 11,022$ patients across $N = 15$ studies with data that are MAR, [Debray and de Jong 2021](#));
- **hiv** from the **GJRM** package (simulated data on HIV diagnoses, $n = 6,416$ patients across $N = 9$ regions with data that are MNAR, [Radice 2021](#)).

For each of these datasets, we discuss the nature of the missingness, choose one or more imputation models and evaluate the imputed data, but we will also highlight one specific aspect of the imputation workflow.

This tutorial is dedicated to readers who are unfamiliar with multiple imputation. More experienced readers can skip the introduction (case study 1) and directly head to practical applications of multilevel imputation under MAR conditions (case study IMPACT) or under MNAR conditions (case study HIV).

Imputation of missing data requires consideration of the mechanism behind the missingness. Rubin proposed to distinguish between data that are missing completely at random (MCAR), data that are missing at random (MAR) and data that are missing not at random (MNAR; see Table 2). For each of these three missingness generating mechanisms, different imputation strategies are warranted ([Yucel \(2008\)](#) and [Hox et al. \(2015\)](#)). We here consider the general case that data are MAR, and expand on certain MNAR situations.

Table 2: Concepts in missing data methods

| Concept | Details |
|---------|--|
| MCAR | Missing Completely At Random, where the probability to be missing is equal across all data entries |
| MAR | Missing At Random, where the probability to be missing depends on observed information |
| MNAR | Missing Not At Random (MNAR), where the probability to be missing depends on unrecorded information, making the missingness non-ignorable (Rubin 1976; Meng 1994). |

2.3. Aim of this paper

This paper serves as a tutorial for imputing incomplete multilevel data with **mice** in R. **mice** has become the de-facto standard for imputation by chained equations, which iteratively solves the missingness on a variable-by-variable basis. **mice** is known to yield valid inferences under many different missing data circumstances (Van Buuren 2018).

We provide practical guidelines and code snippets for different missing data situations, including non-ignorable mechanisms. For reasons of brevity, we focus on multilevel imputation by chained equations with **mice** exclusively; other imputation methods and packages (see e.g. Audigier et al. 2018, and Grund et al. (2018)) are outside the scope of this tutorial. Assumed knowledge includes basic familiarity with the **lme4** notation for multilevel models (see Table ??).

We illustrate imputation of incomplete multilevel data using three case studies:

- **popmis** from the **mice** package (simulated data on perceived popularity, $n = 2,000$ pupils across $N = 100$ schools with data that are MAR, van Buuren and Groothuis-Oudshoorn 2021);
- **impact** from the **metamisc** package (empirical data on traumatic brain injuries, $n = 11,022$ patients across $N = 15$ studies with data that are MAR, Debray and de Jong 2021);
- **obesity** from the **miceheckman** package [simulated data on obesity, $n = 2,111$ patients across $N = 5$ regions with data that are MNAR].

For each of these datasets, we discuss the nature of the missingness, choose one or more imputation models and evaluate the imputed data, but we will also highlight one specific aspect of the imputation workflow.

This tutorial is dedicated to readers who are unfamiliar with multiple imputation. More experienced readers can skip the introduction (case study 1) and directly head to practical applications of multilevel imputation under MAR conditions (case study 2) or under MNAR conditions (case study 3).

TODO: explicit statement about not going into workings of the methods. Galimer 2l methods.

2.4. Setup

[TODO: Add environment info, seed and version number(s) somewhere.] Set up the R environment and load the necessary packages:

Install non-CRAN packages if necessary:

```
R> devtools::install_github("amices/ggmice")
R> devtools::install_github("hanneoberman/miceheckman")
```

Set up the R environment and load the necessary packages:

```
R> set.seed(2022)
R> library(mice)           # for imputation
R> library(miceadds)       # for imputation
R> library(ggmice)         # for visualization
R> library(ggplot2)        # for visualization
R> library(dplyr)          # for data wrangling
R> library(lme4)           # for multilevel modeling
R> library(mitml)          # for multilevel pooling
R> library(miceheckman)    # for imputation cf. heckman models
R> library(metamisc)       # for case study data
```

3. Case study I: popularity data

[TODO: explain case study]

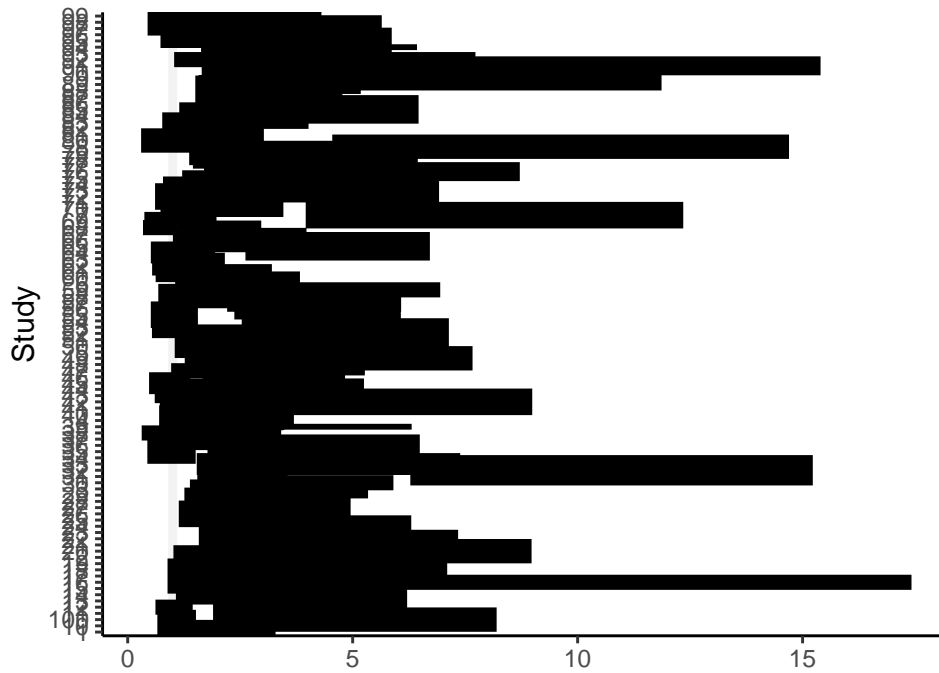
In this section we'll go over the different steps involved with imputing incomplete multilevel data with the R package `mice`. We consider the simulated `popmis` dataset, which included pupils ($n = 2000$) clustered within schools ($N = 100$). The following variables are of primary interest:

- `school`, school identification number (clustering variable);
- `popular`, pupil popularity (self-rating between 0 and 10; unit-level);
- `sex`, pupil sex (0=boy, 1=girl; unit-level);
- `texp`, teacher experience (in years; cluster-level).

The research objective of the `popmis` dataset is to predict the pupils' popularity based on their gender and the experience of the teacher. The analysis model corresponding to this dataset is multilevel regression with random intercepts, random slopes and a cross-level interaction. The outcome variable is `popular`, which is predicted from the unit-level variable `sex` and the cluster-level variable `texp`:

```
R> mod <- popular ~ 1 + sex + (1 | school)
```

The true effect is:



Load the data into the environment and select the relevant variables:

```
R> popmis <- popmis[, c("school", "popular", "sex")]
```

Plot the missing data pattern:

```
R> plot_pattern(popmis)
```

The missingness is univariate and sporadic, which is illustrated in the missing data pattern in Figure 2.

To develop the best imputation model for the incomplete variable `popular`, we need to know whether the observed values of `popular` are related to observed values of other variables. Plot the pair-wise complete correlations in the incomplete data:

```
R> plot_corr(popmis)
```

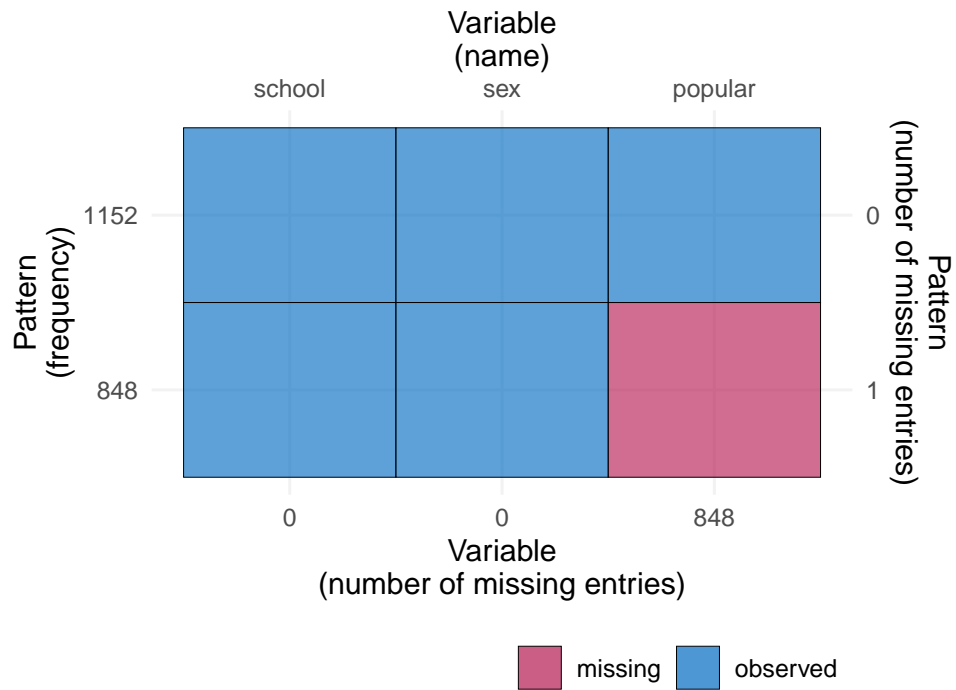
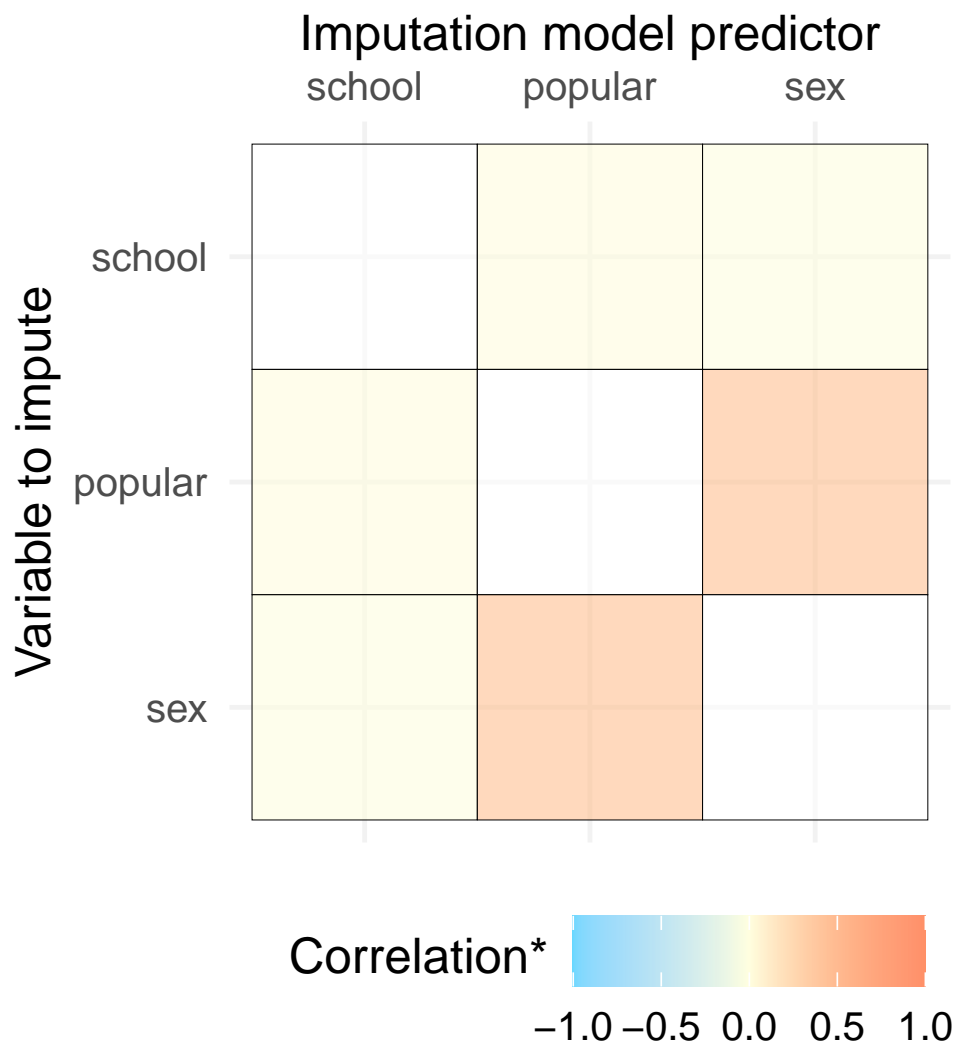



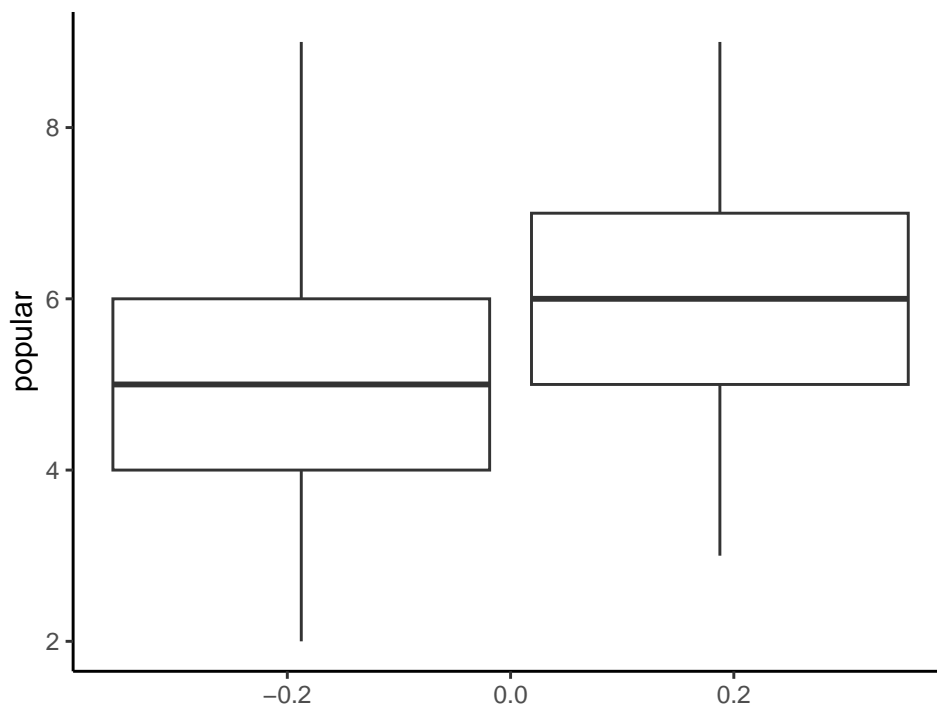
Figure 2: Missing data pattern in the popularity data



*pairwise complete observations

This shows us that `sex` may be a useful imputation model predictor. Moreover, the missingness in `popular` may depend on the observed values of other variables.

```
R> # ggplot(popmis, aes(sex)) +
R> #   geom_histogram(fill = "white") +
R> #   facet_grid(. ~ is.na(popular), scales = "free", labeller = label_both)
R>
R> ggplot(popmis, aes(y = popular, group = sex)) +
+   geom_boxplot() +
+   theme_classic()
```

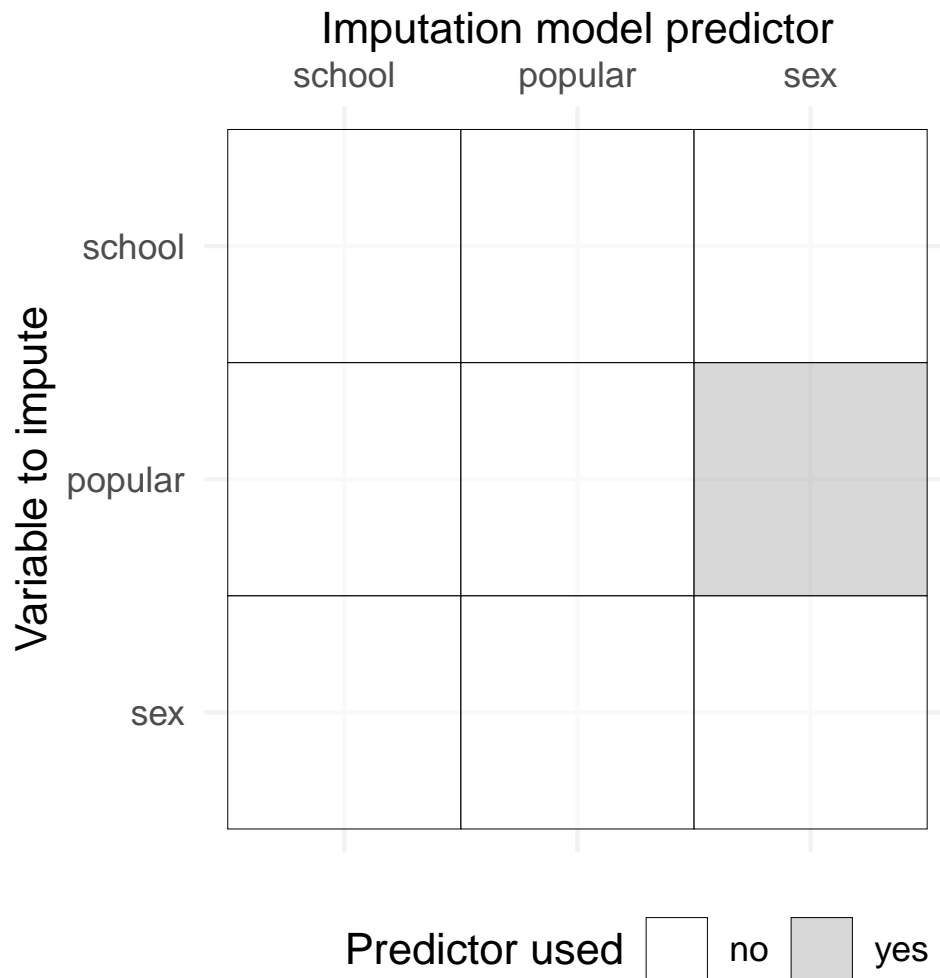


Imputation ignoring the cluster variable (not recommended)

The first imputation model that we'll use is likely to be invalid. We do not use the cluster identifier `school` as imputation model predictor. With this model, we ignore the multilevel structure of the data, despite the high ICC. This assumes exchangeability between units. We include it purely to illustrate the effects of ignoring the clustering in our imputation effort.

Create a methods vector and predictor matrix for `popular`, and make sure `school` is not included as predictor:

```
R> meth <- make.method(popmis) # methods vector
R> pred <- quickpred(popmis)   # predictor matrix
R> plot_pred(pred)
```



Impute the data, ignoring the cluster structure:

```
R> imp <- mice(popmis, pred = pred, print = FALSE)
```

Analyze the imputations:

```
R> fit <- with(imp,
+             lmer(popular ~ 1 + sex + (1 | school)))
```

Print the estimates:

```
R> testEstimates(as.mitml.result(fit), extra.pars = TRUE)
```

Call:

```
testEstimates(model = as.mitml.result(fit), extra.pars = TRUE)
```

Final parameter estimates and inferences obtained from 5 imputed data sets.

| | Estimate | Std.Error | t.value | df | P(> t) | RIV | FMI |
|-------------|----------|-----------|---------|-------|---------|--------|-------|
| (Intercept) | 5.012 | 0.295 | 16.994 | 4.362 | 0.000 | 22.587 | 0.969 |
| sex | 0.695 | 0.251 | 2.768 | 4.287 | 0.047 | 28.390 | 0.975 |

| | Estimate |
|-----------------------------|----------|
| Intercept~~Intercept school | 0.266 |
| Residual~~Residual | 1.035 |
| ICC school | 0.208 |

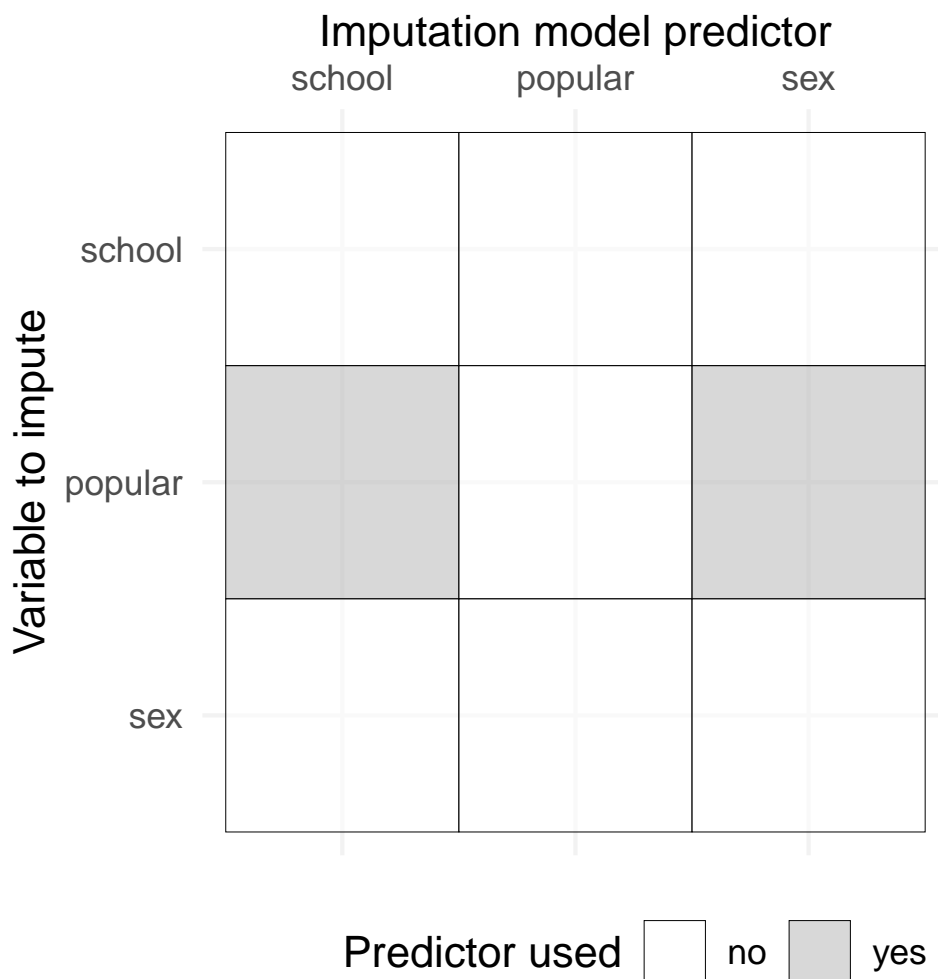
Unadjusted hypothesis test as appropriate in larger samples.

Imputation with the cluster variable as predictor (not recommended)

We'll now use `school` as a predictor to impute all other variables. This is still not recommended practice, since it only works under certain circumstances and results may be biased (Drechsler 2015; Enders, Mistler, and Keller 2016). But at least, it includes some multilevel aspect. This method is also called 'fixed cluster imputation', and uses N-1 indicator variables representing allocation of N clusters as a fixed factor in the model (Reiter, Raghunathan, and Kinney 2006; Enders et al. 2016). Colloquially, this is 'multilevel imputation for dummies'.

[TODO: Add that it doesn't work with systematic missingness (only with sporadic). There's some pros and cons, and it may not even differ much if the number of clusters is low.]

```
R> # adjust the predictor matrix
R> pred["popular", "school"] <- 1
R> plot_pred(pred)
```



```
R> # impute the data, cluster as predictor
R> imp <- mice(popmis, pred = pred, print = FALSE)
```

Analyze the imputations:

```
R> fit <- with(imp,
+             lmer(popular ~ 1 + sex + (1 | school)))
```

Print the estimates:

```
R> testEstimates(as.mitml.result(fit), extra.pars = TRUE)
```

Call:

```
testEstimates(model = as.mitml.result(fit), extra.pars = TRUE)
```

Final parameter estimates and inferences obtained from 5 imputed data sets.

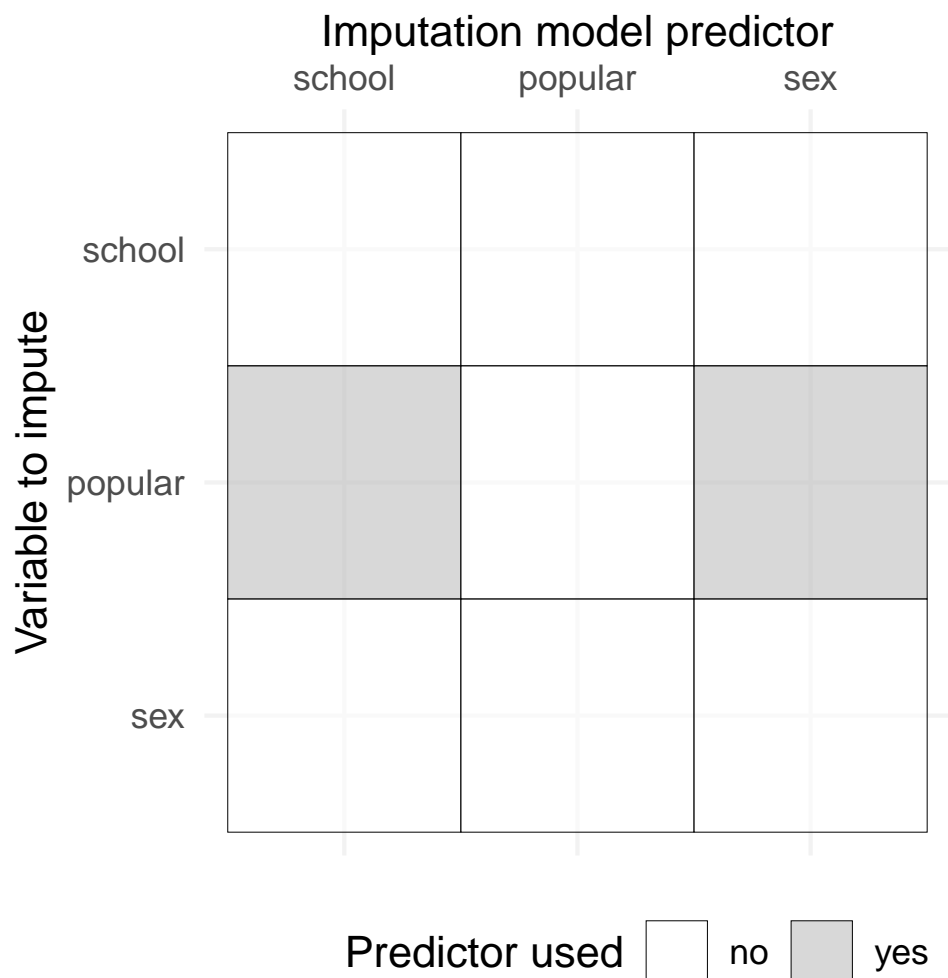
| | Estimate | Std.Error | t.value | df | P(> t) | RIV | FMI |
|-------------|----------|-----------|---------|-------|---------|--------|-------|
| (Intercept) | 4.915 | 0.217 | 22.642 | 4.926 | 0.000 | 9.110 | 0.926 |
| sex | 0.975 | 0.283 | 3.444 | 4.250 | 0.024 | 32.504 | 0.978 |

| | Estimate |
|-----------------------------|----------|
| Intercept~~Intercept school | 0.351 |
| Residual~~Residual | 1.153 |
| ICC school | 0.233 |

Unadjusted hypothesis test as appropriate in larger samples.

Imputation with multilevel model

```
R> # adjust the predictor matrix
R> pred["popular", "school"] <- -2
R> plot_pred(pred)
```



```
R> # impute the data, cluster as predictor
R> imp <- mice(popmis, pred = pred, print = FALSE)
```

Analyze the imputations:

```
R> fit <- with(imp,
+             lmer(popular ~ 1 + sex + (1 | school)))
```

Print the estimates:

```
R> testEstimates(as.mitml.result(fit), extra.pars = TRUE)
```

Call:

```
testEstimates(model = as.mitml.result(fit), extra.pars = TRUE)
```

Final parameter estimates and inferences obtained from 5 imputed data sets.

| | Estimate | Std.Error | t.value | df | P(> t) | RIV | FMI |
|-------------|----------|-----------|---------|-------|---------|--------|-------|
| (Intercept) | 5.011 | 0.410 | 12.222 | 4.226 | 0.000 | 35.955 | 0.980 |
| sex | 0.928 | 0.381 | 2.434 | 4.168 | 0.069 | 48.221 | 0.985 |

| | Estimate |
|-----------------------------|----------|
| Intercept~~Intercept school | 0.313 |
| Residual~~Residual | 1.428 |
| ICC school | 0.188 |

Unadjusted hypothesis test as appropriate in larger samples.

4. Case study II: IMPACT data (syst missingness, pred matrix)

[TODO: check if there is systematic missingness in this dataset, if not make Marshall Computerized Tomography classification (ct) systematically missing.]

We illustrate how to impute incomplete multilevel data by means of a case study: **impact** from the **metamisc** package (empirical data on traumatic brain injuries, $n = 11,022$ units across $N = 15$ clusters, [Debray and de Jong 2021](#)). [TODO: add more info about the complete data.] The **impact** data set contains traumatic brain injury data on $n = 11022$ patients clustered in $N = 15$ studies with the following 11 variables:

- **name** Name of the study,
- **type** Type of study (RCT: randomized controlled trial, OBS: observational cohort),
- **age** Age of the patient,
- **motor_score** Glasgow Coma Scale motor score,
- **pupil** Pupillary reactivity,

- **ct** Marshall Computerized Tomography classification, [TODO: make this one var? also shows that you don't always need random effects everywhere?]
- **hypox** Hypoxia (0=no, 1=yes),
- **hypots** Hypotension (0=no, 1=yes),
- **tsah** Traumatic subarachnoid hemorrhage (0=no, 1=yes),
- **edh** Epidural hematoma (0=no, 1=yes),
- **mort** 6-month mortality (0=alive, 1=dead).

The analysis model for this dataset is a prediction model with **mort** as the outcome. In this tutorial we'll estimate the adjusted prognostic effect of **ct** on mortality outcomes. The estimand is the adjusted odds ratio for **ct**, after including **type**, **age**, **motor_score** and **pupil** into the analysis model:

```
R> mod <- mort ~ 1 + type + age + motor_score + pupil + ct + (1 | name)
```

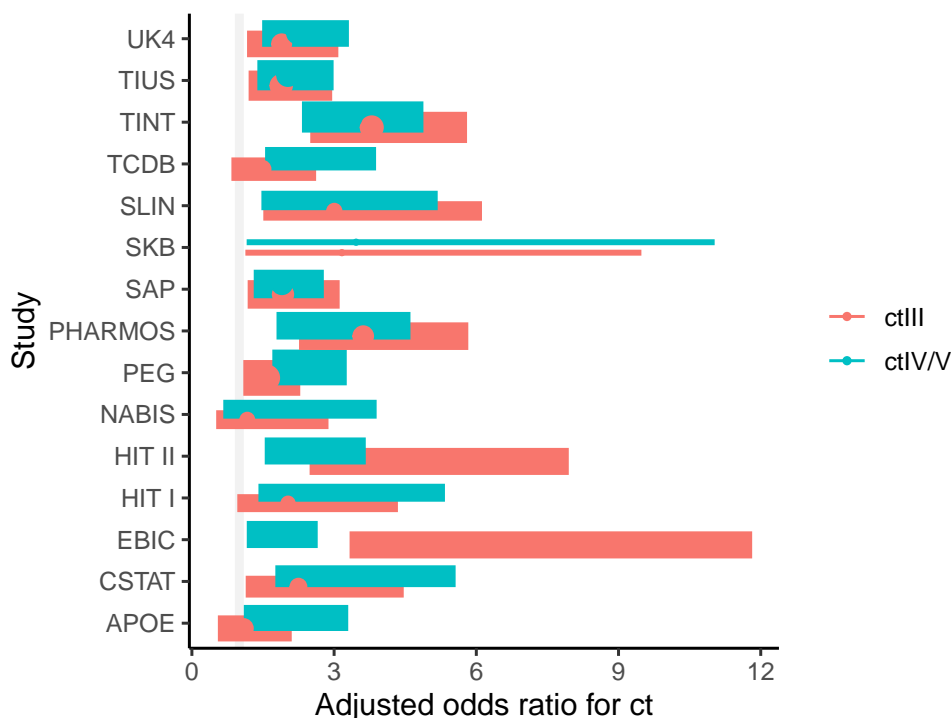
Note that variables **hypots**, **hypox**, **tsah** and **edh** are not part of the analysis model, and may thus serve as auxiliary variables for imputation.

The **impact** data included in the **metamisc** package is a complete data set. The original data has already been imputed once (Steyerberg et al, 2008). For the purpose of this tutorial we have induced missingness (mimicking the missing data in the original data set before imputation). The resulting incomplete data can be accessed from [zenodo link to be created](#).

Load the complete and incomplete data into the R workspace:

```
R> data("impact", package = "metamisc")      # complete data
R> dat <- read.table("link/to/the/data.txt") # incomplete data
```

The estimated effects in the complete data are visualized in Figure ??.




```
R> # fit <- glmer(mod, family = "binomial", data = impact) # fit the model
R> # tidy(fit, conf.int = TRUE, exponentiate = TRUE)      # print estimates
```

[TODO: show how much variance there is after different methods]

[TODO: add ICC before/after imputation and interpret: This tells us that the multilevel structure of the data should probably be taken into account. If we don't, we'll may end up with incorrect imputations, biasing the effect of the clusters towards zero.]

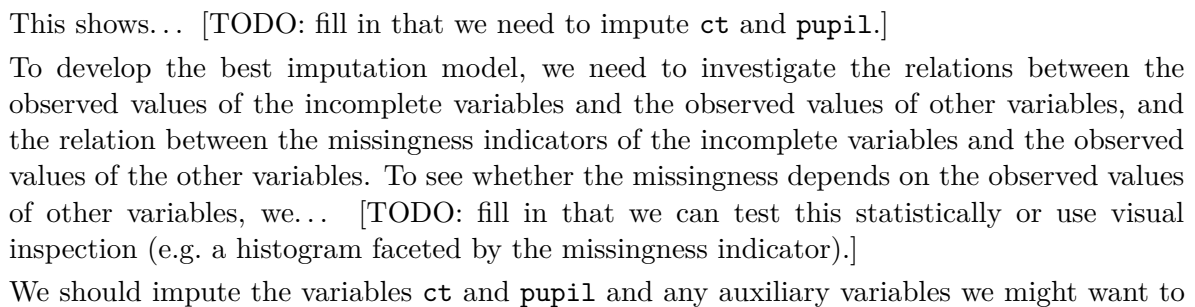
[TODO: add descriptive statistics of the complete and incomplete data.]

4.1. Missingness

To explore the missingness, it is wise to look at the missing data pattern. The ten most frequent missingness patterns are shown:

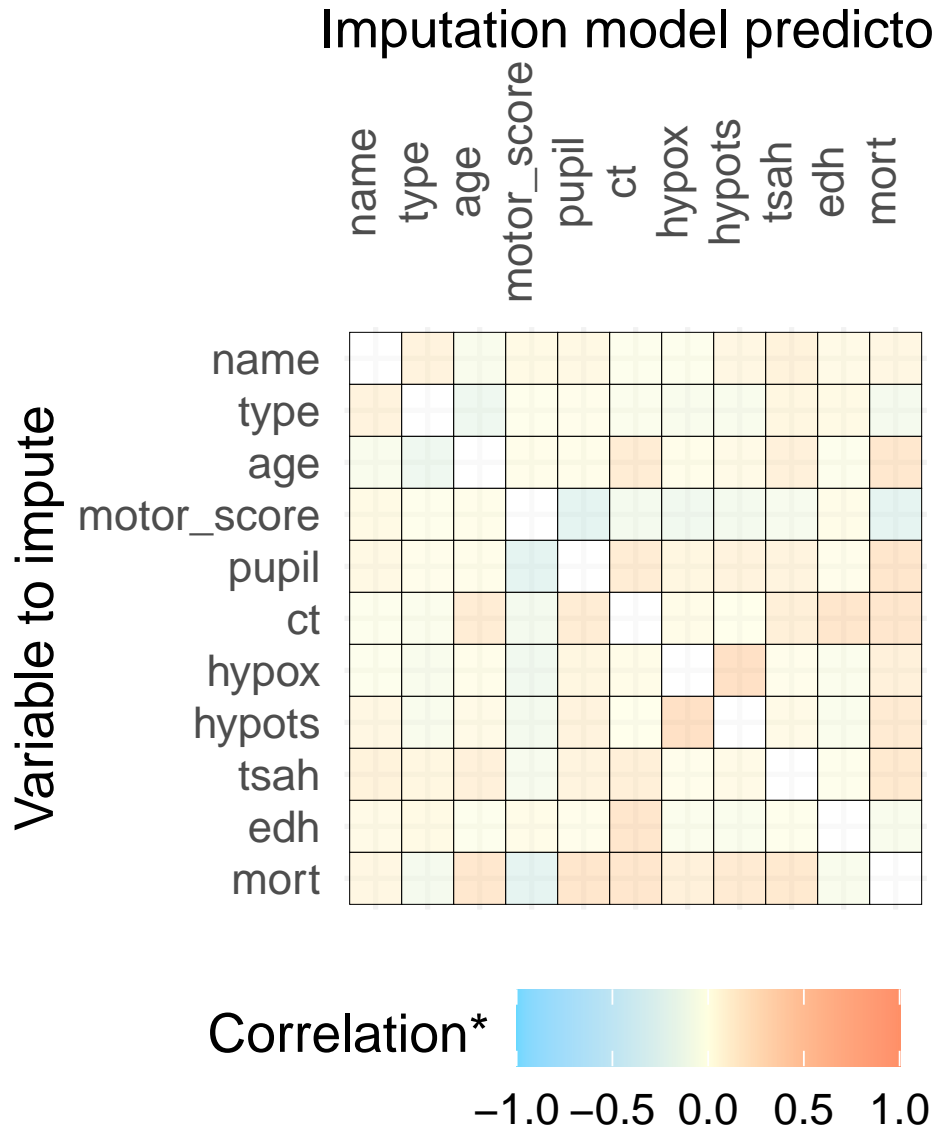
To explore the missingness, it is wise to look at the missing data pattern:

```
R> plot_pattern(dat, rotate = TRUE) # plot missingness pattern
```



use to impute these incomplete analysis model variables. We can evaluate which variables may be useful auxiliaries by plotting the pairwise complete correlations:

```
R> plot_corr(dat, rotate = TRUE) # plot correlations
```



*pairwise complete observations

This shows us that `hypox` and `hypot` would not be useful auxiliary variables for imputing `ct`. Depending on the minimum required correlation, `tsah` could be useful, while `edh` has the strongest correlation with `ct` out of all the variables in the data and should definitely be included in the imputation model. For the imputation of `pupil`, none of the potential auxiliary variables has a very strong relation, but `hypots` could be used. We conclude that we can exclude `hypox` from the data, since this is neither an analysis model variable nor an auxiliary variable for imputation:

```
R> dat <- select(dat, !hypox) # remove variable
```

4.2. Complete case analysis [TODO: remove this?]

As previously stated, complete case analysis lowers statistical power and may bias results. The complete case analysis estimates are:

```
R> fit <- glmer(mod, family = "binomial", data = na.omit(dat)) # fit the model
R> tidy(fit, conf.int = TRUE, exponentiate = TRUE) # print estimates
```

```
# A tibble: 11 x 9
  effect group term estimate std.error statistic p.value conf.low conf.high
  <chr>   <chr> <chr>   <dbl>     <dbl>     <dbl>     <dbl>     <dbl>     <dbl>
1 fixed  <NA> (Int~  0.0863  0.0182    -11.6  3.00e-31  0.0571  0.130
2 fixed  <NA> type~  0.757   0.137     -1.54  1.22e- 1  0.531  1.08
3 fixed  <NA> age   1.03   0.00265    12.9  7.40e-38  1.03   1.04
4 fixed  <NA> moto~ 0.651   0.0732     -3.82  1.34e- 4  0.522  0.811
5 fixed  <NA> moto~ 0.489   0.0555     -6.30  2.97e-10  0.391  0.611
6 fixed  <NA> moto~ 0.274   0.0321    -11.0  2.28e-28  0.218  0.345
7 fixed  <NA> pupi~ 3.20    0.317     11.7  8.18e-32  2.63   3.88
8 fixed  <NA> pupi~ 1.75    0.195      5.06  4.27e- 7  1.41   2.18
9 fixed  <NA> ctIII 2.41    0.268      7.89  3.05e-15  1.94   2.99
10 fixed <NA> ctIV~ 2.30    0.214      8.95  3.56e-19  1.92   2.76
11 ran_pa~ name sd__~ 0.230   NA        NA     NA     NA     NA     NA
```

As we can see... [TODO: fill in: This means that a higher ct (Marshall Computerized Tomography classification) is associated with a lower odds of 6-month mortality, given by the odds ratio $\exp(0.42)$, CI ... to ..., when controlling for ...]

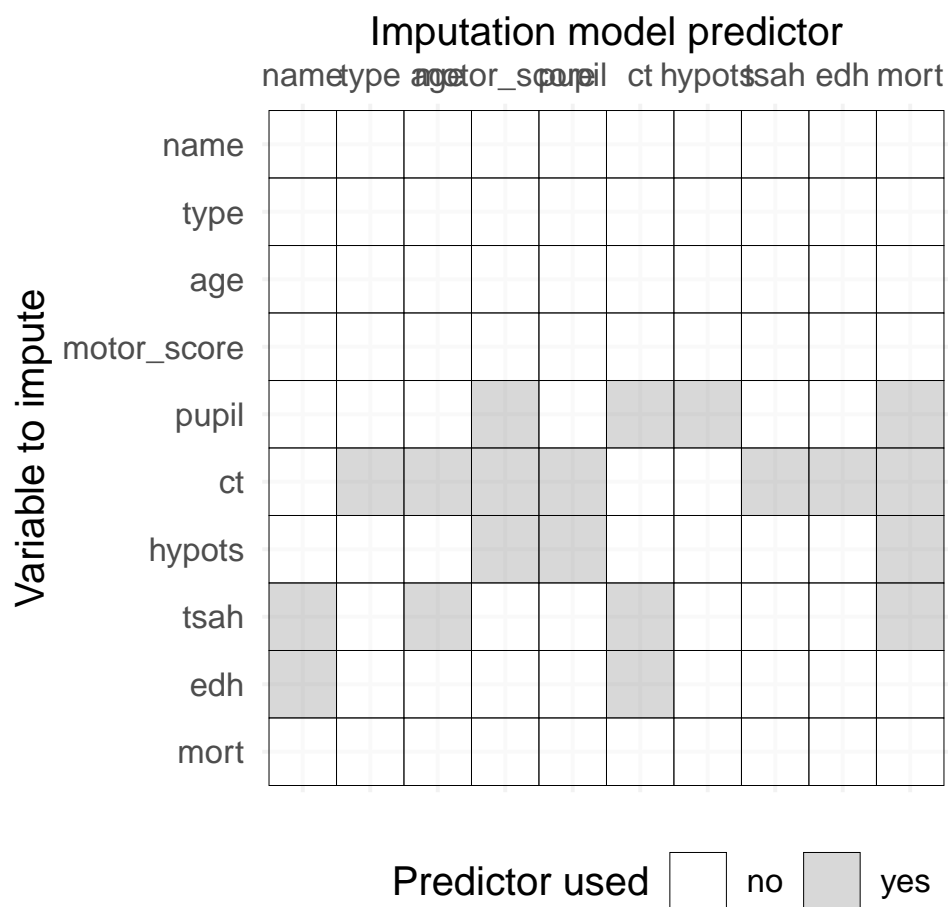
4.3. Imputation model

Mutate data to get the right data types for imputation (e.g. integer for clustering variable).

```
R> dat <- dat %>% mutate(across(everything(), as.integer))
```

Create a methods vector and predictor matrix, and make sure **name** is not included as predictor, but as clustering variable:

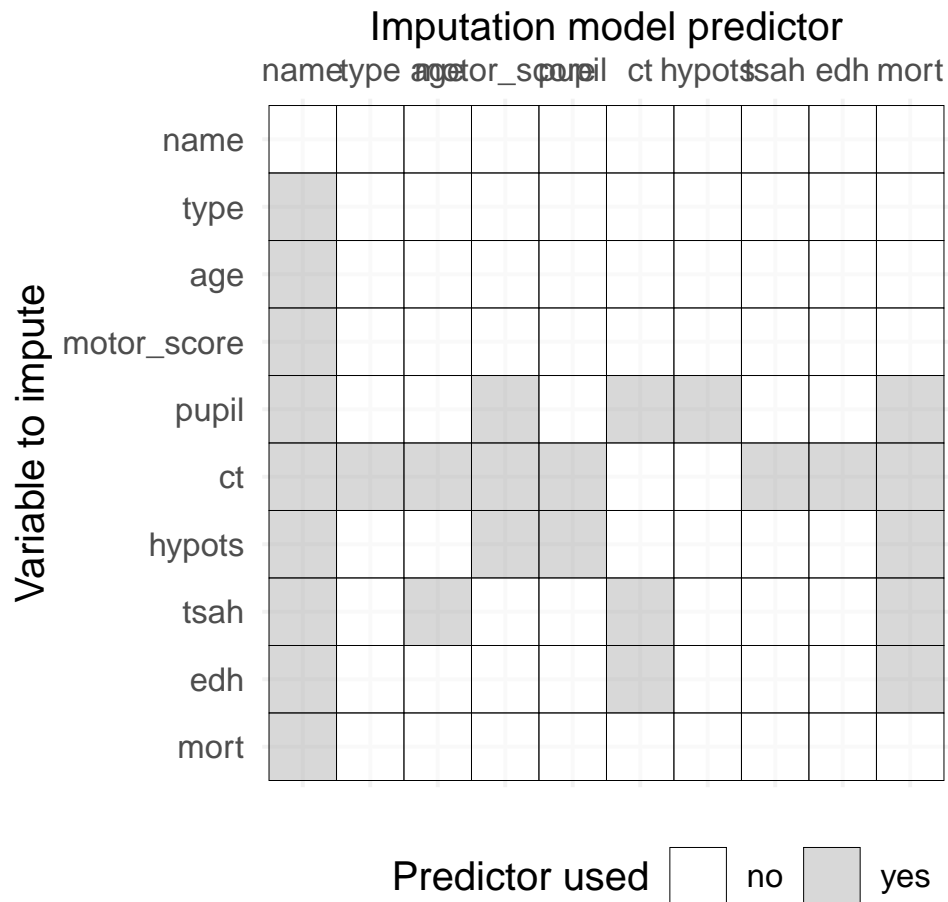
```
R> meth <- make.method(dat) # methods vector
R> pred <- quickpred(dat)   # predictor matrix
R> plot_pred(pred)
```



```

R> pred[pred == 1] <- 2
R> pred["mort", ] <- 2
R> pred[, "mort"] <- 2
R> pred[c("name", "type", "age", "motor_score", "mort"), ] <- 0
R> pred[, "name"] <- -2
R> diag(pred) <- 0
R> plot_pred(pred)

```



```
R> meth <- make.method(dat)
R> meth
```

| name | type | age | motor_score | pupil | ct |
|--------|-------|-------|-------------|-------|-------|
| "" | "" | "" | "" | "pmm" | "pmm" |
| hypots | tsah | edh | mort | | |
| "pmm" | "pmm" | "pmm" | "" | | |

Impute the incomplete data

```
R> imp <- mice(dat, method = meth, predictorMatrix = pred, printFlag = FALSE)
```

```
R> fit <- imp %>%
+   with(glmer(mort ~ type + age + as.factor(motor_score) + pupil + ct + (1 | name), famil
R> tidy(pool(fit))
```

| | term | estimate | std.error | statistic | p.value |
|---|-------------|-------------|-------------|-----------|--------------|
| 1 | (Intercept) | -2.35203726 | 0.340181747 | -6.914061 | 4.994037e-12 |
| 2 | type | -0.41265892 | 0.180274846 | -2.289054 | 2.209524e-02 |
| 3 | age | 0.03049023 | 0.001570162 | 19.418521 | 1.238416e-81 |

```

4 as.factor(motor_score)2 -0.66764920 0.068737865 -9.712975 3.480413e-22
5 as.factor(motor_score)3 -1.05520001 0.070218940 -15.027285 2.540225e-50
6 as.factor(motor_score)4 -1.51238349 0.072304262 -20.916934 1.850073e-90
7
8 pupil 0.48421447 0.038982800 12.421234 6.772069e-17
8 ct 0.43474621 0.029968474 14.506785 2.342253e-36
      b      df dfcom      fmi      lambda m      riv
1 5.281599e-04 10119.70041 11013 0.005673265 0.005476771 5 0.005506932
2 4.881699e-05 10893.89378 11013 0.001985736 0.001802528 5 0.001805783
3 3.335358e-08 6320.89582 11013 0.016545467 0.016234340 5 0.016502243
4 4.188703e-05 8327.24771 11013 0.010875748 0.010638213 5 0.010752602
5 5.135721e-05 7632.20045 11013 0.012757637 0.012498967 5 0.012657168
6 1.403058e-04 2831.75462 11013 0.032888241 0.032205434 5 0.033277139
7 3.571497e-04 49.97295 11013 0.309130930 0.282023647 5 0.392803531
8 8.586645e-05 294.69765 11013 0.120677044 0.114729598 5 0.129598366
      ubar
1 1.150898e-01
2 3.244044e-02
3 2.425385e-06
4 4.674630e-03
5 4.869071e-03
6 5.059539e-03
7 1.091079e-03
8 7.950697e-04

```

```
R> as.mitml.result(fit)
```

```

[[1]]
Generalized linear mixed model fit by maximum likelihood (Laplace
Approximation) [glmerMod]
Family: binomial ( logit )
Formula: mort ~ type + age + as.factor(motor_score) + pupil + ct + (1 |
name)
      AIC      BIC    logLik deviance df.resid
10495.423 10561.192 -5238.712 10477.423     11013
Random effects:
Groups Name      Std.Dev.
name (Intercept) 0.2843
Number of obs: 11022, groups: name, 15
Fixed Effects:
      (Intercept)                type                age
      -2.37195                -0.41014                0.03052
as.factor(motor_score)2 as.factor(motor_score)3 as.factor(motor_score)4
      -0.65802                -1.04611                -1.51245
      pupil                ct
      0.50405                0.42496
optimizer (Nelder_Mead) convergence code: 0 (OK) ; 0 optimizer warnings; 1 lme4 warnings

```

```
[[2]]
Generalized linear mixed model fit by maximum likelihood (Laplace
  Approximation) [glmerMod]
Family: binomial ( logit )
Formula: mort ~ type + age + as.factor(motor_score) + pupil + ct + (1 |
  name)
      AIC      BIC    logLik deviance df.resid
10500.88 10566.65 -5241.44 10482.88    11013
Random effects:
Groups Name      Std.Dev.
name (Intercept) 0.2917
Number of obs: 11022, groups:  name, 15
Fixed Effects:
              (Intercept)                  type                  age
                -2.37718                  -0.41511                  0.03067
as.factor(motor_score)2 as.factor(motor_score)3 as.factor(motor_score)4
                -0.66935                  -1.05211                  -1.49429
                pupil                      ct
                 0.49013                  0.43835
optimizer (Nelder_Mead) convergence code: 0 (OK) ; 0 optimizer warnings; 1 lme4 warnings
```

```
[[3]]
Generalized linear mixed model fit by maximum likelihood (Laplace
  Approximation) [glmerMod]
Family: binomial ( logit )
Formula: mort ~ type + age + as.factor(motor_score) + pupil + ct + (1 |
  name)
      AIC      BIC    logLik deviance df.resid
10505.026 10570.795 -5243.513 10487.026    11013
Random effects:
Groups Name      Std.Dev.
name (Intercept) 0.2908
Number of obs: 11022, groups:  name, 15
Fixed Effects:
              (Intercept)                  type                  age
                -2.32339                  -0.42359                  0.03023
as.factor(motor_score)2 as.factor(motor_score)3 as.factor(motor_score)4
                -0.67142                  -1.05776                  -1.51038
                pupil                      ct
                 0.49756                  0.42474
optimizer (Nelder_Mead) convergence code: 0 (OK) ; 0 optimizer warnings; 1 lme4 warnings
```

```
[[4]]
Generalized linear mixed model fit by maximum likelihood (Laplace
  Approximation) [glmerMod]
Family: binomial ( logit )
Formula: mort ~ type + age + as.factor(motor_score) + pupil + ct + (1 |
```



```

      name)
      AIC      BIC    logLik deviance df.resid
10519.511 10585.280 -5250.755 10501.511    11013
Random effects:
  Groups Name      Std.Dev.
  name      (Intercept) 0.2961
Number of obs: 11022, groups:  name, 15
Fixed Effects:
              (Intercept)                  type                  age
              -2.33581                  -0.40871                  0.03039
as.factor(motor_score)2 as.factor(motor_score)3 as.factor(motor_score)4
              -0.66477                  -1.05451                  -1.51858
              pupil                  ct
              0.45928                  0.44419
optimizer (Nelder_Mead) convergence code: 0 (OK) ; 0 optimizer warnings; 2 lme4 warnings

[[5]]
Generalized linear mixed model fit by maximum likelihood (Laplace
Approximation) [glmerMod]
Family: binomial ( logit )
Formula: mort ~ type + age + as.factor(motor_score) + pupil + ct + (1 |
      name)
      AIC      BIC    logLik deviance df.resid
10522.038 10587.807 -5252.019 10504.038    11013
Random effects:
  Groups Name      Std.Dev.
  name      (Intercept) 0.2955
Number of obs: 11022, groups:  name, 15
Fixed Effects:
              (Intercept)                  type                  age
              -2.35187                  -0.40574                  0.03064
as.factor(motor_score)2 as.factor(motor_score)3 as.factor(motor_score)4
              -0.67468                  -1.06551                  -1.52622
              pupil                  ct
              0.47006                  0.44148
optimizer (Nelder_Mead) convergence code: 0 (OK) ; 0 optimizer warnings; 1 lme4 warnings

attr(,"class")
[1] "mitml.result" "list"

R> # testEstimates(as.mitml.result(fit))

```

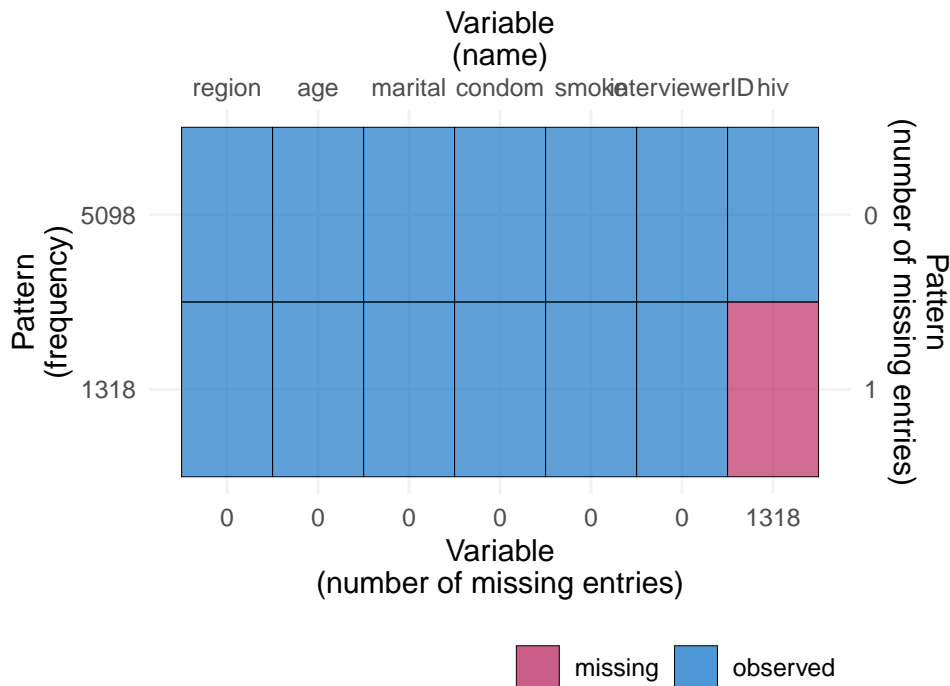
5. Case study III: HIV data

Data are simulated and included in the `miceheckman` package. We will use the following variables:

Data are simulated and included in the `GJRM` package. We will use the following variables:

- `region` Cluster variable,
- `hiv` HIV diagnosis (0=no, 1=yes),
- `age` Age of the patient,
- `marital` Marital status,
- `condom` Condom use during last intercourse,
- `smoke` Smoker (levels; inclusion restriction variable).

The imputation of these data is based on the toy example from [IPDMA Heckman Github repo](#).



- `region` Cluster variable,
- `gender` Gender (0=male, 1=female),
- `age` Age of the patient,
- `height` Height in meters,
- `weight` Weight in kilograms,
- `time` Response time in minutes (exclusion restriction variable).

The imputation of these data is based on the [IPDMA Heckman Github repo](#)
load the data

```
R> data("obesity", package = "miceheckman")
```

We obtain here the random effects for each interviewer, this is an approximation of the interviewer's skill which will be used as an exclusion constraint. Here, since the location of

the interviewer was not randomly assigned to the subjects, the assignment was corrected for region and language.

Set the Heckman model as imputation method

```
R> # Set prediction matrix and methods
R> #ini <- mice(obesity, maxit = 0)
R> #meth<-ini$method
R> #meth["weight"]<-"2l.heckman"
R> #pred <- ini$pred
R> #pred["weight","cluster"]<- -2
R> #pred["weight","rt"] <- -3
```

Select relevant variables:

```
R> #dat <- select(obesity, ~bmi)
```

Visualize missing data pattern:

```
R> #plot_pattern(dat)
```

Create predictor matrix:

```
R> #pred <- quickpred(dat) # predictor matrix
R> #pred["weight","cluster"]<- -2 # clustering variable
R> #pred["weight","rt"] <- -3 #inclusion-restriction variable
R> #plot_pred(pred)
```

Set the Heckman model as imputation method:

```
R> #meth <- make.method(dat) # methods vector
R> #meth["weight"]<-"2l.heckman"
```

Impute the missingness:

```
R> #imp <- mice(obesity, # dataset with missing values
R> #           m = 10, # number of imputations
R> #           maxit = 1,
R> #           seed = 1234, #seed attached to the dataID
R> #           meth = meth, #imputation method vector
R> #           pred = pred, #imputation predictors matrix
R> #           print = T,
R> #           meta_method="reml",
R> #           pmm=FALSE)
```

6. Discussion

- JOMO in **mice** -> on the side for now
- Additional levels of clustering
- More complex data types: timeseries and polynomial relationship in the clustering.
- FIML vs MI

An alternative approach to missing data is to use Full Information Maximum Likelihood (FIML). This method does not require the imputation of any missing values. Whereas MI consists of imputation, analyses and pooling steps, FIML analyses the data in a single step. When the assumptions are met the two approaches should produce equivalent results. [REF] As FIML requires specialised software, not all analyses can be performed with standard software. [REF]

- Survival / TTE, this could be put in the paragraph on congeniality

When the outcome is time-to-event, the Nelson-Aalen estimate of the time to event should be included as a covariate in the imputation model [REF]

- Additional levels of clustering
- More complex data types: timeseries and polynomial relationship in the clustering.

7. Think about

- Adding evaluations of the imputations such as convergence checks
- Adding some kind of help function to mice that suggests a suitable predictor matrix to the user, given a certain analysis model.
- Adding a `multilevel_ampute()` wrapper function in mice.
- Exporting `mids` objects to other packages like `lme4` or `coxme`?
- Adding a ICC=0 dataset to show that even if there is no clustering it doesn't hurt.
- Show use case for deductive imputation for cluster level variables?
- env dump in repo
- I don't know if in your article you cover something about model complexity, for example sometimes I have to switch from 2l. methods to 1l. methods just because the model didn't converge.. this is due to the considered imputation model is very complex regarding the amount of information counted... I know that a solution for an imputation model with many predictors is to check correlation plots as you did or use `quickpred()`.. (maybe you can add this somewhere after the correlation plots)....But as for cluster specification, I don't know if besides plots of distribution per cluster there is something else can be done to see if i have to use 1l. or 2l. for a given variable, also I have no idea.. how to test which is better between 2l.norm or 2l.2stage.norm.

In hierarchical datasets, clustering is a concern because the homoscedasticity in the error terms cannot be assumed across clusters and the relationship among variables may vary at different hierarchical levels. When multiple imputation is used to deal with missing data, as the imputation and analysis process is performed separately, it is necessary that imputation model being congenial with the main analysis model (Meng, 1994), e.g. if the main model accounts for the hierarchical structure also imputation model should do it (Audigier, 2021). Not including clustering into the imputation process may lead to effect estimates with smaller standard errors and inflated type I error.

There are different strategies that can be adopted in the imputation process that account for clustering: inclusion of cluster indicator variable, performing a separate imputation process for each cluster, or performing a simultaneous imputation process by using an imputation method that accounts for clustering. (Stata: <https://www.stata.com/support/faqs/statistics/clustering-and-mi-impute/>)

The selection of each strategy depends mainly on the assumptions in the main analysis and also on the restriction of the analyzed data.

Regarding the restrictions imposed by the data, for instance, the use of cluster indicator variables is restricted in datasets where there are not many clusters and many observations per cluster (Graham, 2009). The last restriction is also required when imputations are performed on each cluster separately. When this restriction cannot be achieved, one can use an imputation model that simultaneously imputes all clusters using a hierarchical model (Allison 2002).

Under this hierarchical imputation model, observations within clusters are correlated and this correlation is modelled by a random effect so the hierarchical model can be estimated even when there are few observations per cluster. However, this strategy is best suited for balanced data (Grund, 2017) and when random effects model is appropriated, i.e. the number of clusters is adequate. (Austin, 2018).

Here it is important to evaluate the assumptions imposed by the main model, for instance by using the cluster indicator strategy may lead to bias estimates when the model is based on a hierarchical model (<https://onlinelibrary.wiley.com/doi/full/10.1002/bimj.201900051>). Even when an imputation strategy congenial with the main model is preferred, it is important to consider whether it is appropriate for the data as a less complex imputation strategies may also lead to unbiased estimates in certain scenarios (Bailey 2020). For instance, in causal effect analysis, separately imputation may lead to smaller bias when the size of the smaller exposure cluster is large, compared with an imputation model that includes exposure-confounder interactions. (Zhang, 2023).

8. Funding

This project has received funding from the European Union's Horizon 2020 research and innovation programme under ReCoDID grant agreement No 825746.

The views expressed in this paper are the personal views of the authors and may not be understood or quoted as being made on behalf of or reflecting the position of the regulatory agency/agencies or organizations with which the authors are employed/affiliated.

9. References

10. Appendix

Table 3: Notation based on <https://datalorax.github.io/equationomatic/articles/lmer4-lmer.html> <https://cran.r-project.org/web/packages/lme4/vignettes/lmer.pdf>

| Formula lme4 | Details |
|---|--|
| $y \sim x1 + (1 g1)$ | Fixed $x1$ predictor with random intercept varying among $g1$ |
| $y \sim x1 + (x1 g1)$ or $1 + x1 + (1+x1 g1)$ | Fixed $x1$ with correlated random intercept and random slope of $x1$ |
| $y \sim x1 + (x1 g1)$ or $1 + x1 + (1 g1) + (0 + x1 g1)$ | Fixed $x1$ with uncorrelated random intercept and random slope of $x1$ |
| $y \sim x1 + x2 + (x1 g1) + (x2 g1)$ | variance-covariance matrix estimated separately for each fixed predictor, i.e, one for intercept and $x1$ and another for intercept and $x2$ |
| $y \sim (1 g1) + (1 g2)$ | Random intercept varying among $g1$ and among $g2$ |

Table 4: Imputation methods for hierarchical data, based on Audigier 2018 study comparison.

| Type approach | Imputation method | Type of variables | Systematically missing data | Missing mechanism | Heteroscedasticity | Details |
|---------------|---|-------------------|-----------------------------|-------------------|--------------------|--|
| JM | Pan | C-B | Yes | MAR | | |
| | Jomo (Quartagno and Carpenter, 2016) | C-B | Yes | MAR | Y | Time consuming, recommended for large clusters when proportion of binary is high |
| FCS | Pan | C | Yes | MAR | | |
| | 2lnorm | C | Yes | MAR | | |
| | GLM (Jolani, 2017) | C-B | No | MAR | N | Recommended for small clusters |
| | 2l.stage (Resche-Rigon and White, 2016, Audigier, 2018) | C-B | Yes | MAR/MNAR | | Combine estimates with random intercept model, recommended for large studies |

References

- Audigier V, White IR, Jolani S, Debray TPA, Quartagno M, Carpenter J, van Buuren S, Resche-Rigon M (2018). “Multiple Imputation for Multilevel Data with Continuous and Binary Variables.” *Statistical Science*, **33**(2), 160–183. ISSN 0883-4237, 2168-8745. doi:[10.1214/18-STS646](https://doi.org/10.1214/18-STS646). [1702.00971](https://doi.org/10.1214/18-STS646).
- de Jong VMT, Moons KGM, Eijkemans MJC, Riley RD, Debray TPA (2021). “Developing More Generalizable Prediction Models from Pooled Studies and Large Clustered Data Sets.” *Statistics in Medicine*, **40**(15), 3533–3559. ISSN 1097-0258. doi:[10.1002/sim.8981](https://doi.org/10.1002/sim.8981).
- Debray T, de Jong V (2021). “Metamisc: Meta-Analysis of Diagnosis and Prognosis Research Studies.”
- Drechsler J (2015). “Multiple Imputation of Multilevel Missing Data—Rigor Versus Simplicity.” *Journal of Educational and Behavioral Statistics*, **40**(1), 69–95. ISSN 1076-9986. doi:[10.3102/1076998614563393](https://doi.org/10.3102/1076998614563393).
- Enders CK, Mistler SA, Keller BT (2016). “Multilevel Multiple Imputation: A Review and Evaluation of Joint Modeling and Chained Equations Imputation.” *Psychological Methods*, **21**(2), 222–240. ISSN 1939-1463. doi:[10.1037/met0000063](https://doi.org/10.1037/met0000063).
- Gelman A, Hill J (2006). *Data Analysis Using Regression and Multilevel/Hierarchical Models*. Cambridge University Press. ISBN 978-1-139-46093-4.
- Grund S, Lüdtke O, Robitzsch A (2018). “Multiple Imputation of Missing Data for Multilevel Models: Simulations and Recommendations.” *Organizational Research Methods*, **21**(1), 111–149. ISSN 1094-4281. doi:[10.1177/1094428117703686](https://doi.org/10.1177/1094428117703686).
- Hox J, van Buuren S, Jolani S (2015). “Incomplete Multilevel Data: Problems and Solutions.” In J Harring, L Stapleton, S Beretvas (eds.), *Advances in Multilevel Modeling for Educational Research: Addressing Practical Issues Found in Real-World Applications*, CILVR Series on Latent Variable Methodology, pp. 39–62. Information Age Publishing Inc., Charlotte, NC. ISBN 978-1-68123-328-4.
- Hox JJ, Moerbeek M, van de Schoot R (2017). *Multilevel Analysis: Techniques and Applications, Third Edition*. Routledge. ISBN 978-1-317-30868-3.
- Jolani S (2018). “Hierarchical Imputation of Systematically and Sporadically Missing Data: An Approximate Bayesian Approach Using Chained Equations.” *Biometrical Journal. Biometrische Zeitschrift*, **60**(2), 333–351. ISSN 1521-4036. doi:[10.1002/bimj.201600220](https://doi.org/10.1002/bimj.201600220).
- Localio AR, Berlin JA, Ten Have TR, Kimmel SE (2001). “Adjustments for Center in Multicenter Studies: An Overview.” *Annals of Internal Medicine*, **135**(2), 112–123. ISSN 0003-4819. doi:[10.7326/0003-4819-135-2-200107170-00012](https://doi.org/10.7326/0003-4819-135-2-200107170-00012).
- Meng XL (1994). “Multiple-Imputation Inferences with Uncongenial Sources of Input.” *Statistical Science*, **9**(4), 538–558. ISSN 0883-4237, 2168-8745. doi:[10.1214/ss/1177010269](https://doi.org/10.1214/ss/1177010269).
- Radice GMar (2021). “GJRM: Generalised Joint Regression Modelling.”
- Reiter JP, Raghunathan T, Kinney SK (2006). “The Importance of Modeling the Sampling Design in Multiple Imputation for Missing Data.” *undefined*.

- Resche-Rigon M, White IR, Bartlett JW, Peters SAE, Thompson SG (2013). “Multiple Imputation for Handling Systematically Missing Confounders in Meta-Analysis of Individual Participant Data.” *Statistics in medicine*, **32**(28), 4890–4905. ISSN 1097-0258 0277-6715. doi:10.1002/sim.5894.
- Rubin DB (1976). “Inference and Missing Data.” *Biometrika*, **63**(3), 581–592. doi:10.2307/2335739.
- Van Buuren S (2018). *Flexible Imputation of Missing Data*. Chapman and Hall/CRC.
- van Buuren S, Groothuis-Oudshoorn K (2021). “Mice: Multivariate Imputation by Chained Equations.”
- Yucel RM (2008). “Multiple Imputation Inference for Multivariate Multilevel Continuous Data with Ignorable Non-Response.” *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, **366**(1874), 2389–2403. doi:10.1098/rsta.2008.0038.

Affiliation:

Hanne I. Oberman
 Utrecht University
 Padualaan 14
 3584 CH Utrecht
 E-mail: h.i.oberman@uu.nl
 URL: <https://hanneoberman.github.io/>

Thomas P. A. Debray
 Julius Center for Health Sciences and Primary Care
 University Medical Center Utrecht, Utrecht University,
 Utrecht, The Netherlands