



Imputation of Incomplete Multilevel Data with **mice**

Hanne Oberman
Utrecht University

Johanna Munoz Avila
University Medical Center Utrecht

Valentijn de Jong
University Medical Center Utrecht

Gerko Vink
Utrecht University

Thomas Debray
University Medical Center Utrecht

Abstract

Multilevel data is not spared the ubiquitous problem of missing information. This is a tutorial paper on imputing incomplete multilevel data with **mice**. Including methods for ignorable and non-ignorable missingness. Footnotes in the current version show work in progress/under construction. The last section is not part of the manuscript, but purely for reminders.

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

1. Introduction

1.1. Multilevel data

Research into any field with a hierarchical or clustered nature of observations may yield multilevel data. In the typical case, individuals are nested within groups, but there are many different types of multilevel data. In the medical field, clustering occurs at e.g., the hospital or center level in registry data, or at the study-level in meta-analyses (IPDMA¹). In the social sciences and official statistics we can find clustering e.g. at the country-level, or as imposed by a multi-stage sampling design. For the sake of legibility, we will refer to the grouping variable

¹remove or write in full?

as ‘cluster’, and the grouped variable as ‘(sample) unit’ throughout this paper.² And, for reasons of brevity, we only discuss clustering between units, not within units (e.g., timeseries or longitudinal data).³

Analyzing multilevel data requires special care, compared to ‘regular’, single level data. For instance, there may be cross-level interactions between unit-level variables and cluster-level variables. The cluster to which a unit belongs may influence the unit-level observations, and vice versa for the the units that make up the cluster (Hox, Moerbeek, and van de Schoot 2017). These relations can and should be taken into account when developing analysis models for multilevel data.⁴ Multilevel models typically include separate intercepts for each cluster, which relieves one restriction imposed by single-level models: equal group means across clusters. Additionally, there may be random predictor effects and/or random error terms (residual error variances), see e.g. Hox *et al.* (2017) and de Jong, Moons, Eijkemans, Riley, and Debray (2021).⁵ There are many names for models that take clustering into account. Some popular examples are ‘multilevel models’, ‘hierarchical models’, ‘mixed effect models’ and ‘random effect models’.

1.2. Missing data

Multilevel data is not spared the ubiquitous problem of missing information. Just as in single level data, missingness may occur at the unit level. But with multiple levels of data comes the potential for missingness at multiple levels. Missingness in multilevel data can be categorized into two general patterns: systematic missingness and sporadic missingness, see Resche-Rigon, White, Bartlett, Peters, Thompson, and Group (2013). In figure 1, we show a dataset with units in the rows and variables in the columns, there are 5 units nested within 3 clusters, and 3 variables of interest. Variable X1 is completely observed. Variable X2 is systematically missing, X3 is sporadically missing.⁶ Systematic missingness can be further subdivided into unobserved constants (i.e., the same value within clusters) and non-measured random variables (which may differ per unit within clusters). In Figure 1, the former implies that the unobserved values for units 3 and 4 on variable X2 would be equal. With the latter, the values may differ. Depending on the missing data pattern⁷, there are more or less optimal way of accommodating the missingness.

	X1	X2	X3
cluster1{	x_{11}	x_{12}	NA
	x_{21}	x_{22}	x_{32}
cluster2{	x_{31}	NA	x_{33}
	x_{41}	NA	NA
cluster3{	x_{51}	x_{52}	x_{53}

²Alternatives to unit: participant/response/record.

³Also, add that we’ll only discuss two levels, not more?

⁴Explain ICC here? The percentage of variance attributed to the cluster-level is expressed by the intra-class coefficient (ICC). The ICC can also be interpreted as the expected correlation between two randomly sampled units in same cluster. So if the ICC is high, a lot of variability in a variable is due to the clustering, which should be modeled accordingly.

⁵Add that heterogeneity refers to variability within clusters.

⁶Explain why.

⁷add missing data mechanisms here?

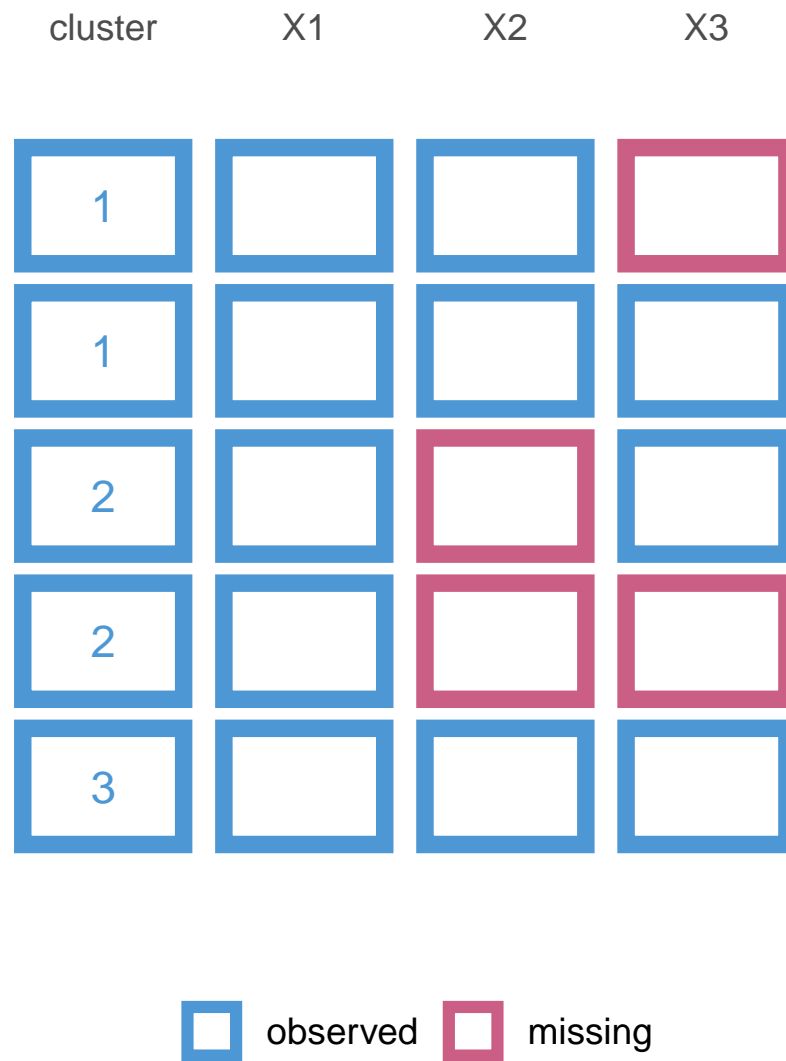


Figure 1: Missingness in multilevel data

Ignoring missing data in research endeavors is almost never a good idea. Complete case analysis (i.e., excluding all units with one or more missing entries) can introduce bias in statistical inference and lowers statistical power. Instead, the missingness should be accommodated *before* or *within* the analysis of scientific interest. Especially the former is very generic and popular. Imputing (i.e., filling in) the missing values splits the missing data problem from the scientific problem. The R package **mice** has become the de-facto standard for imputation by chained equations, which solves the missingness one variable at a time, iteratively. **mice** is known to yield valid inferences under many different missing data circumstances (Van Buuren 2018). In this paper, we'll discuss how to use **mice** in the context of multilevel data, under varying missing data mechanisms.⁸

1.3. Aim of this paper

This paper serves as a tutorial for imputing incomplete multilevel data with **mice**. We provide practical guidelines and code snippets for different missing data situations. For reasons of brevity, we focus on imputation by chained equations⁹. Other useful packages for incomplete multilevel data include **mitml**, **miceadds**, and **mdmb**.¹⁰

We structure this tutorial around three case studies:

- `mice::popmis` (simulated data on school kids, with MNAR/MAR mixture);
- `metamisc::impact` (real IPD on traumatic brain injuries, without NAs);
- `GJRM::hiv` (simulated patient data on HIV, without NAs)

For each case study we focus on a different aspect to illustrate how to impute incomplete multilevel data. In the `mice::popmis` data, we show the advantages of including the multilevel structure of the data into the imputation model. In the `metamisc::impact` data we'll show how to induce missingness and solve it in real-world data. In the `GJRM::hiv` we provide novel methodology¹¹ for imputing MNAR missingness according to the Heckman model. For all case studies we discuss the nature of the incomplete data, the imputation model(s), and evaluation of the imputed data.

2. Case Study I: Popularity

`popNCR` is a simulated dataset with pupils clustered in classes, where the number of units $n = 2000$, and the number of clusters $N = 100$, on 7 variables:

- `pupil` Pupil number within class,

⁸Discuss missingness mechanisms before this point, add references Yucel (2008) and Hox, van Buuren, and Jolani (2015).

⁹add that JOMO is available in **mice** as well?

¹⁰Rephrase: Some level of knowledge on multilevel models is assumed. We're providing an overview of implementations. It's up-to the reader to decide which multilevel strategy suits their data. So we won't go into detail for the different methods (and equations). Refer to Meng (1994), Audigier, White, Jolani, Debray, Quartagno, Carpenter, van Buuren, and Resche-Rigon (2018), and Grund, Lüdtke, and Robitzsch (2018). This paper is just a software tutorial. We'll keep it practical.

¹¹not really, the methods exist already, but how to show that this is something new and exciting?

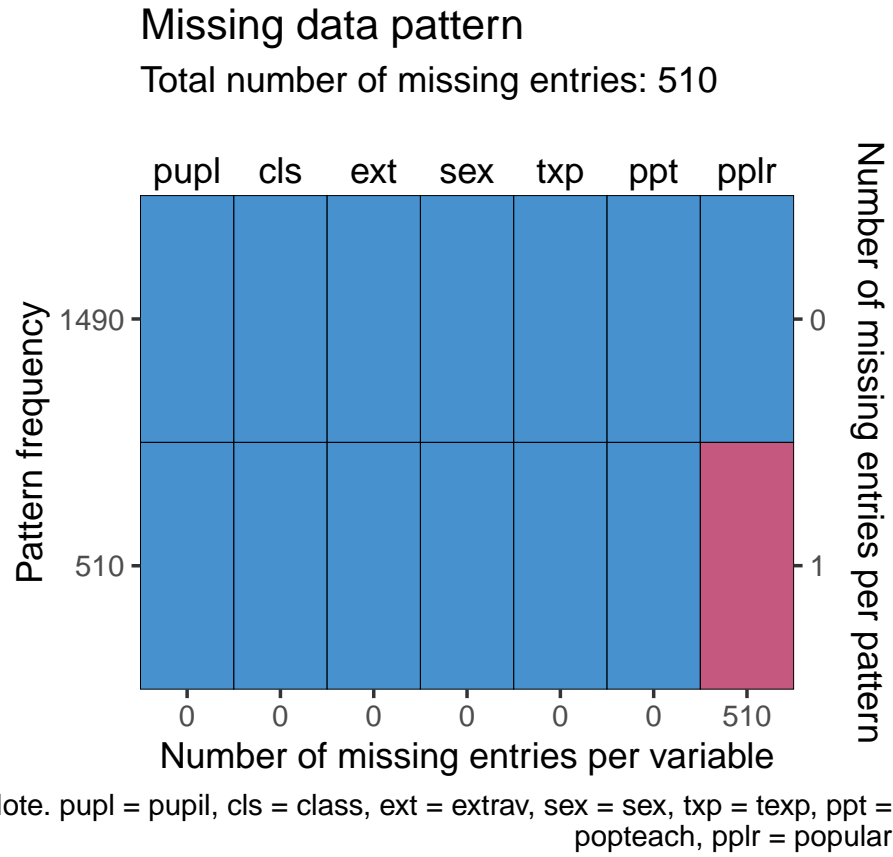


Figure 2: Missing data pattern in the popularity data

- **class** Class number,
- **extrav** Pupil extraversion,
- **sex** Pupil gender,
- **txp** Teacher experience (years),
- **popular** Pupil popularity,
- **popteach** Teacher popularity.

Incomplete data

The popularity data is created such that there are strong relations between the incomplete variables and the clustering variable **class**. We can express this using the intra-class correlation (ICC). For **popular** the ICC is 0.33. For **popteach** it is 0.34. It would thus be wise to use multilevel modeling.

The missingness in this dataset is induced conform MAR and MNAR mechanisms. The missing data pattern, Figure 2, shows that just one variable is incomplete **[the next part is not yet updated to reflect this]**.

To develop the best imputation model, we need to know whether the missingness in one

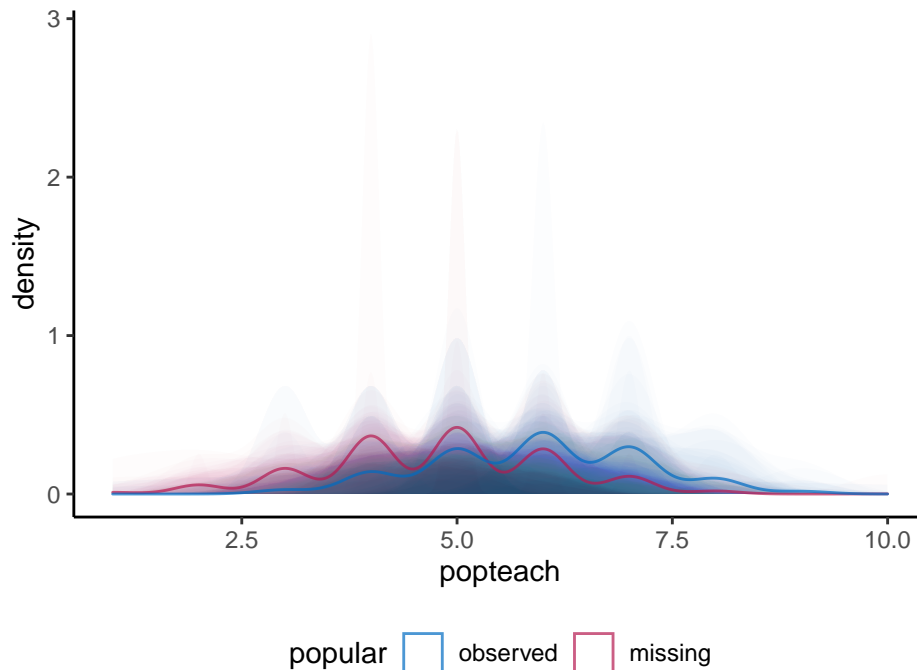


Figure 3: Conditional distributions in the popularity data

variable depends on the observed values of other variables. Visual inspection usually suffices. We'll highlight only two variables to illustrate, but ideally one would inspect all relations. The questions we'll ask are: 'Does the missing data of **popular** depend on **popteacher**?' and 'Does the missingness in teacher popularity depend on pupil popularity?' We'll evaluate this by making a histogram of **popteacher** separately for the pupils with known popularity and missing popularity, and the other way around.

In Figure 3 we see that the distribution for the missing **popular** is further to the right than the distribution for observed **popular**. This would indicate a right-tailed MAR missingness. In fact, this is exactly what happens, because the missingness in these data was created manually. Now, we've made it observable by examining the relations between the missingness in **popular** and the observed data in **popteacher**. There is also a dependency between the missingness in teacher popularity and pupil popularity. The relation seems to be right-tailed as well.

Imputation model

The first imputation model that we'll use is likely to be invalid. In this model, we ignore the multilevel structure of the data, despite the high ICCs. This is purely to illustrate the effects of ignoring the clustering in our imputation effort.

We'll use predictive mean matching to impute the continuous variables and logistic regression to impute the binary variable **sex**. We do not use the observation identifier **pupil** or cluster identifier **class** as predictors to impute other variables.

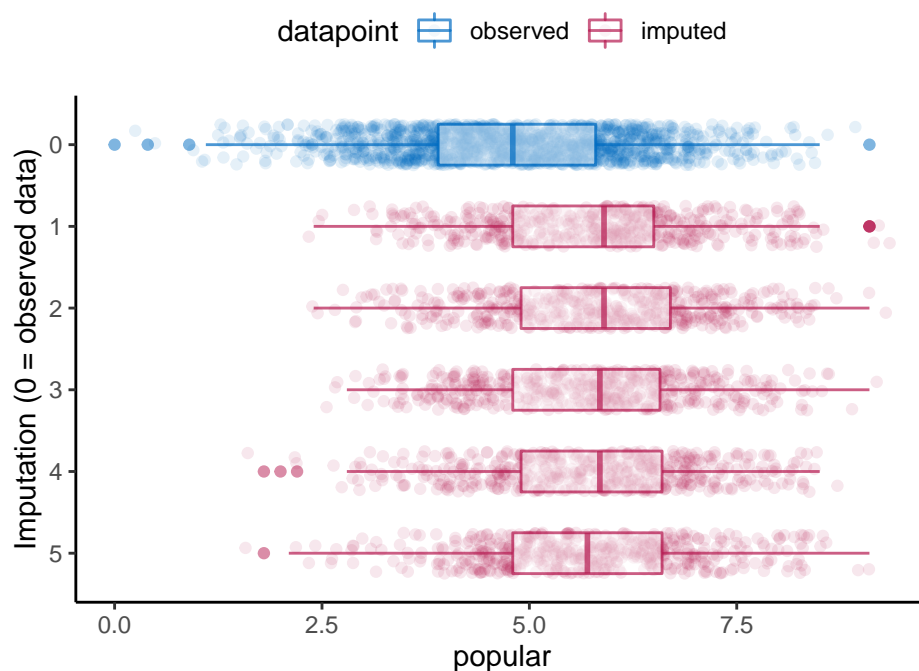
```
R> # dry run to get imputation parameters
```

```

R> ini <- mice(pop, maxit = 0)
R>
R> # extract predictor matrix and adjust
R> pred <- ini$pred
R> pred[, c("class", "pupil")] <- 0
R>
R> # impute the data, ignoring the cluster structure
R> imp_ignored <- mice(pop, maxit = 1, pred = pred, print = FALSE)

```

Imputed data



	vars	incomplete	ignored
1	popular	0.3280070	0.3338626
2	popteach	0.3414766	0.3414766
3	texp	1.0000000	1.0000000

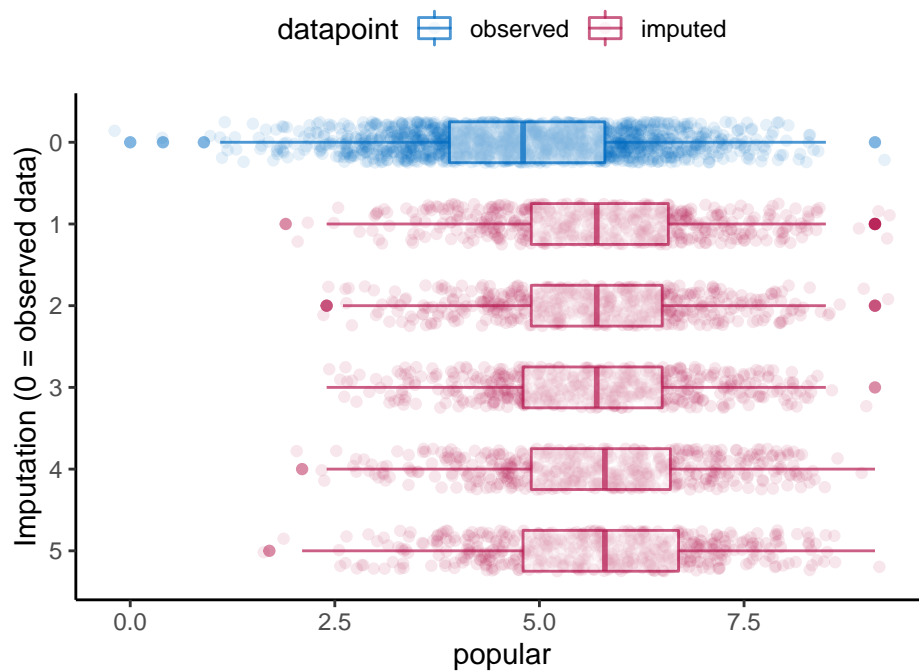
As the original ICCs show, 100% of the variance in `texp` can be attributed to the clustering variable `class`. This tells us that the multilevel structure of the data should be taken into account. If we don't, we'll end up with incorrect imputations, biasing the effect of the clusters towards zero.

We can also observe that the teacher experience increases slightly after imputation. This is due to the MNAR missingness in `texp`. Higher values for `texp` have a larger probability to be missing. This may not be a problem, however, if at least one pupil in each class has teacher experience recorded, we can deductively impute the correct (i.e. true) value for every pupil in the class.

Imputation model

We'll now use `class` 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. But at least, it includes some multilevel aspect. Colloquially, this is 'multilevel imputation for dummies'.

```
R> # adjust the predictor matrix
R> pred <- ini$pred
R> pred[, "pupil"] <- 0
R>
R> # impute the data, cluster as predictor
R> imp_predictor <- mice(pop, maxit = 1, pred = pred, print = FALSE)
```

Imputed data

	vars	incomplete	ignored	predictor
1	popular	0.3280070	0.3338626	0.3871020
2	popteach	0.3414766	0.3414766	0.3414766
3	texp	1.0000000	1.0000000	1.0000000

Now, we can clearly see that the imputed values of `texp` are higher than the observed values, which is in line with right-tailed MNAR.

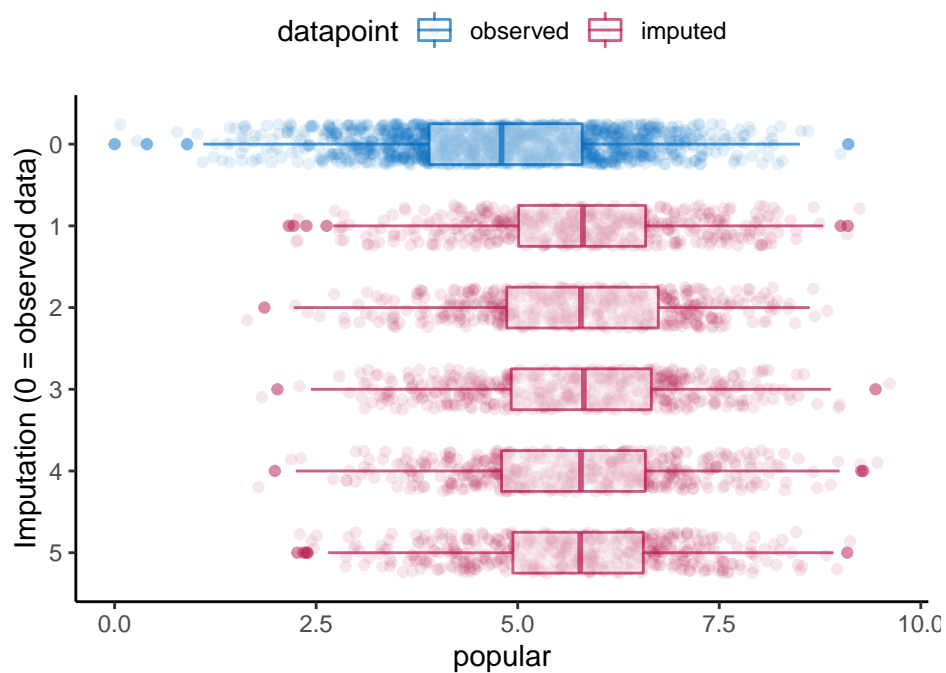
The ICCs are way more in line with the ICCs in the incomplete data. But this is a quick and dirty way of imputing multilevel data. We *should* be using a multilevel model.

Imputation model

To include...

```
R> pred <- ini$pred
R> pred["popular", ] <- c(0, -2, 2, 2, 2, 0, 2)
R> #-2 for the cluster variable, 2 for random effects
R> meth <- ini$meth
R> meth <- c("", "", "", "", "", "2l.norm", "")
R> imp_norm_2l <-
+   mice(
+     pop %>% mutate(class = as.integer(class)),
+     pred = pred,
+     meth = meth,
+     maxit = 1,
+     print = FALSE
+   )

R> # plot(imp_norm)
R> plot_box(imp_norm_2l, x = "popular", strip = TRUE)
```



```
R> ICCs <- ICCs %>% mutate(
+   norm = c(icc(popular ~ as.factor(class), complete(imp_norm_2l)),
+           icc(popteach ~ as.factor(class), complete(imp_norm_2l)),
+           icc(teexp ~ as.factor(class), complete(imp_norm_2l)))
+ )
R> ICCs
```

```

      vars incomplete  ignored predictor      norm
1 popular  0.3280070 0.3338626 0.3871020 0.3612830
2 popteach 0.3414766 0.3414766 0.3414766 0.3414766
3      texp  1.0000000 1.0000000 1.0000000 1.0000000

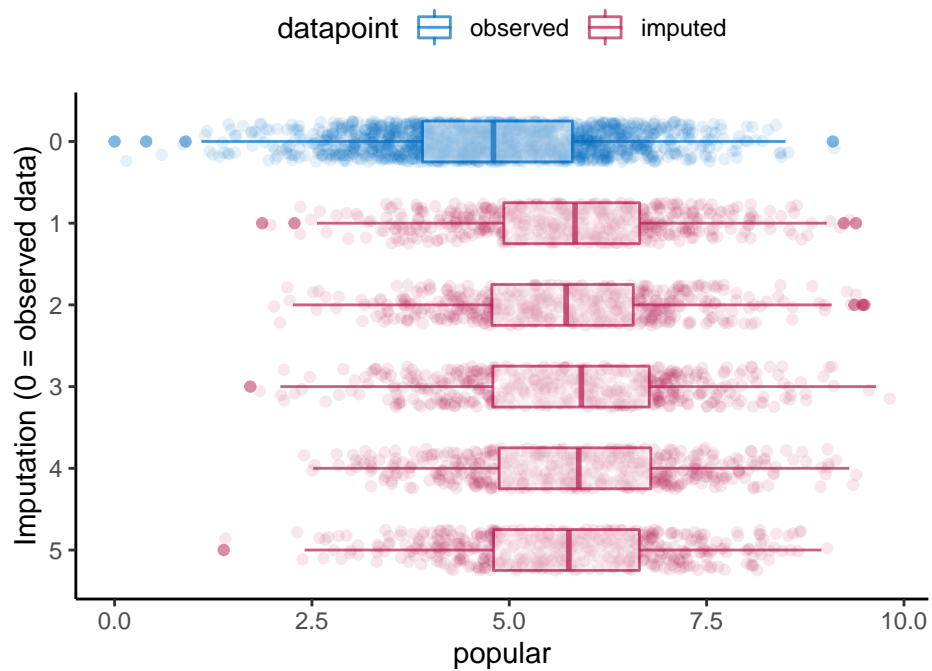
```

```

R> pred["popular", ] <- c(0, -2, 2, 2, 1, 0, 2)
R> meth <- c("", "", "", "", "", "2l.pan", "")
R> imp_pan_2l <-
+   mice(
+     pop %>% mutate(class = as.integer(class)),
+     pred = pred,
+     meth = meth,
+     maxit = 1,
+     print = FALSE
+   )

R> # plot(imp_pan)
R> plot_box(imp_pan_2l, x = "popular", strip = TRUE)

```



```

R> ICCs <- ICCs %>% mutate(
+   pan = c(icc(popular ~ as.factor(class), complete(imp_pan_2l)),
+           icc(popteach ~ as.factor(class), complete(imp_pan_2l)),
+           icc(texp ~ as.factor(class), complete(imp_pan_2l)))
+ )
R> ICCs

```

	vars	incomplete	ignored	predictor	norm	pan
1	popular	0.3280070	0.3338626	0.3871020	0.3612830	0.3745108
2	popteach	0.3414766	0.3414766	0.3414766	0.3414766	0.3414766
3	texp	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000

3. Case study II: IMPACT

impact is traumatic brain injury data with patients, $n = 11022$, clustered in studies, $N = 15$. 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,
- **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 data is already imputed (Steyerberg et al, 2008), so we'll induce missingness ourselves. For example, MAR missingness varying by cluster.¹²

4. Case study III: HIV

Toy example from [Heckman Github repo](#).

5. Discussion

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

6. Think about

- Adding some kind of help function to **mice** that suggests a suitable predictor matrix to the user, given a certain analysis model.

¹²Observed data pattern should differ per cluster. So, in cluster 1, the missingness would depend on age, but not in cluster two. Split the dataframe and run **ampute()** on each cluster.

- 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.

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/10.1702.00971).
- 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).
- 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.
- 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).
- Resche-Rigon M, White IR, Bartlett JW, Peters SA, Thompson SG, Group obotPIS (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. doi:[10.1002/sim.5894](https://doi.org/10.1002/sim.5894).
- Van Buuren S (2018). *Flexible Imputation of Missing Data*. Chapman and Hall/CRC.
- 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](https://doi.org/10.1098/rsta.2008.0038).

Affiliation:

Hanne Oberman

Utrecht University

Padualaan 14

3584 CH Utrecht

E-mail: h.i.oberman@uu.nl

URL: <https://hanneoberman.github.io/>