



Multiple Imputation using R
Sarah R Haile, PD PhD (sarah.haile@uzh.ch)
Version 1.4 of April 16, 2026

Contents

1	Basic steps	2
2	Complete Cases analysis	2
3	A basic example	3
3.1	Impute the missing observations using <code>mice()</code>	3
3.2	Run regression model on imputed data using <code>with()</code>	5
3.3	Combine the results using <code>pool()</code>	5
4	Extended example	5
5	Cox proportional hazards regression	8
6	Multiple imputation with clustered or longitudinal data	9
7	Predictions / fitted values after MI	12
8	Frequently Asked Questions	13
8.1	Do I need to include imputed values in Table 1?	13
8.2	How do I report my MI analysis?	13
8.3	How many imputations?	13
8.4	Which variables should be included in the imputation step?	14
8.5	How can I use <code>mice()</code> when performing Cox PH regression?	14
8.6	Can I show odds ratios after logistic regression?	14
8.7	Which method do I use to impute the variables?	15
8.8	In my analysis, <code>pool</code> doesn't seem to work. Now what?	15
8.9	Can I pool quantities that are not regression coefficients?	17
8.10	Can I combine MI with cross-validation or bootstrap sampling?	17
8.11	Can I use parallel processing to run <code>mice</code> ?	17
9	References	18
A	Recent Changes	20
B	Code to make datasets used here	20

Multiple imputation (MI) is a common statistical method used to analyze datasets where some values are missing. In this document we describe multiple imputation briefly, and show how to perform the analysis in R. The main and extended examples show a dataset where the outcome is binary, and logistic regression is used. After that, we show shorter examples for linear regression and Cox proportional hazards regression.

For more information about MI, you might find the following other references helpful. [Sterne et al. \[2009\]](#) and [White et al. \[2011\]](#) provide overviews on MI, while [van Buuren and Groothuis-Oudshoorn \[2011\]](#) describes the `mice` package in R and gives examples. [White and Royston \[2009\]](#) specifically discusses MI when survival data are present. Several papers give good overviews of MI in epidemiology: [Perkins et al. \[2017\]](#), [Harel et al. \[2017\]](#), [Pedersen et al. \[2017\]](#). For further details on the underlying framework of the R package, `mice`, described here, see also the book *Flexible Imputation of Missing Data* (FIMD) [[van Buuren, 2018](#)] (full text available online), and the accompanying vignettes ([Vink and van Buuren \[2019\]](#), also linked from `?mice`).

1 Basic steps

The basic steps in R are as follows:

- `mice()` Impute the data. That is, make `m` copies of the original dataset and fill in the missing values.
- `with()` Analyze each of the completed datasets.
- `pool()` Combine the parameter estimates using Rubin's rules.

2 Complete Cases analysis

The `birthwt` dataset (from the R package `MASS`) has been adapted for the analysis described here (see code in Appendix). The variables are:

- `low` indicator of birth weight less than 2.5 kg.
- `bwt` birth weight in grams.
- `age` mother's age in years.
- `lwt` mother's weight in pounds at last menstrual period.
- `race` mother's race (1 = white, 2 = black, 3 = other).
- `smoke` smoking status during pregnancy.
- `pt1` number of previous premature labours.
- `ht` history of hypertension.
- `ui` presence of uterine irritability.
- `ftv` number of physician visits during the first trimester.

Using `md.pattern()` we can examine the pattern of missingness in the data. Each combination of missing variables is given a row in the output, with 1 in the row indicating observed and 0 indicating missing. Here we see that while most of the 189 subjects have complete observations [130, top row], many observations have a single missing variable (only one 0 in the row), and some observations have multiple variables missing (several 0s in the row). Other functions to visualize patterns in missing data are also available from the `VIM` package.

```
md.pattern(dat[, c("bwt", "age", "smoke", "race", "ftv")], plot = TRUE)
```

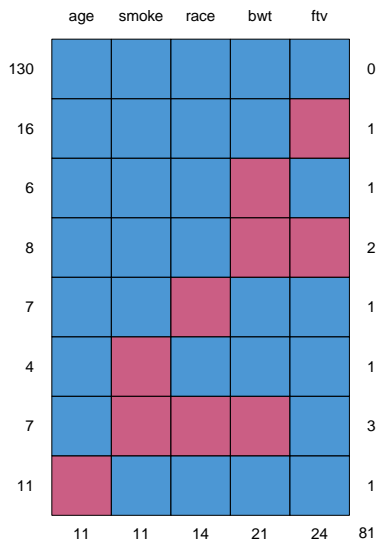


Figure 1: Pattern of missing values

First, we consider a regression model using only the complete cases: predictors for birthweight using linear regression.

```
library(tidyverse)
library(broom)
m1.cc <- lm(bwt ~ age + smoke, data = dat)
tidy(m1.cc, conf.int = TRUE)

## # A tibble: 3 x 7
##   term          estimate std.error statistic p.value  conf.low  conf.high
##   <chr>          <dbl>    <dbl>    <dbl> <chr>    <dbl>    <dbl>
## 1 (Intercept)  2927.    283.    10.3  <0.0001  2368.    3486.
## 2 age           2.16     11.9     0.181 0.856   -21.4    25.7
## 3 smoke        -238.    123.    -1.94 0.054   -480.     4.36
```

We would like to see if our results are affected by the missing values.

3 A basic example

Now, we use MI to guess what the missing variables "could have been", and then use our imputed datasets to estimate the regression coefficients again.

3.1 Impute the missing observations using `mice()`.

There are various methods used to impute missing values for multiple variables at once. We recommend using `mice()` ("Multivariate Imputation using Chained Equations"). **Before running** `mice()`,

check that your dataset contains only variables you will use in your analysis, or variables you think are related to missingness (see Section 8.4).

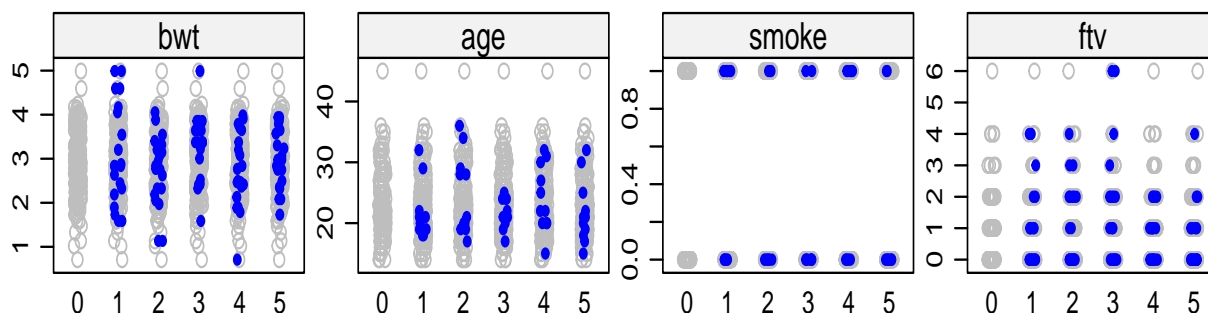
The model above considered 3 variables (bwt, age, smoke), and we think that race and ftv may be related to missingness, so we keep only those 5 variables. Then, we use the `mice()` command to impute the missing variables multiple times, using all other variables (including the outcome) as predictors in the imputation procedure (see Section 8.4). We also set a seed to ensure that we will get the same results any time we run the code. For more on which regression methods to use or how many datasets to impute see Sections 8.7 and 8.3.

```
library(mice)
small <- dat %>%
  select(bwt, age, smoke, race, ftv) %>%
  mutate(bwt = bwt / 1000) # convert g to kg
imp <- mice(small, m = 5, print = FALSE, seed = 12345)
imp

## Class: mids
## Number of multiple imputations: 5
## Imputation methods:
##      bwt      age      smoke      race      ftv
##      "pmm"    "pmm"    "pmm" "polyreg"  "pmm"
## PredictorMatrix:
##      bwt age smoke race ftv
## bwt   0  1   1   1   1
## age   1  0   1   1   1
## smoke 1  1   0   1   1
## race  1  1   1   0   1
## ftv   1  1   1   1   0
```

In this example, we see that `mice()` used `pmm` (predictive mean matching) to impute all the variables except race which used `polyreg`, a form of multinomial logistic regression. (Variables with no missing values will have no method ("") because no imputation method is needed.) Reading across the rows of the predictor matrix, we can also see that `mice()` used each of the other variables to impute both variables.

```
stripplot(imp, col=c("grey", "blue"), pch = c(1, 20))
```



Using `stripplot.mids()`¹ or `densityplot.mids()`, we can see that the imputed values are similar to the observed values for both age and bmi, indicating that imputation is probably appropriate for this analysis.

3.2 Run regression model on imputed data using `with()`.

To run a regression model with imputed data, we have to use `with()` (see `?with.mids`). Note that we use the same model formulation as above, but we leave out the `data` option. Note that we use the same imputed data to run several models, so there is no need to impute new data for every model of interest. In general we will not look at the results of `with()` directly, but instead `pool()` them first. We can however take a look at the analyses and get the results of each of the `m` fitted regression models.

```
m1.mi <- with(imp, lm(bwt ~ age + smoke))
sapply(m1.mi$analyses, coef)

##           [,1]      [,2]      [,3]      [,4]      [,5]
## (Intercept) 2.8151  2.8405  2.8423  2.8651  2.9905
## age         0.0092  0.0068  0.0089  0.0069  0.0013
## smoke      -0.1948 -0.2432 -0.2353 -0.3176 -0.2300
```

3.3 Combine the results using `pool()`.

```
summary(pool(m1.mi), conf.int = TRUE)

## # A tibble: 3 x 8
##   term      estimate std.error statistic    df p.value conf.low conf.high
##   <fct>      <dbl>    <dbl>    <dbl> <dbl> <chr>    <dbl>    <dbl>
## 1 (Intercept)  2.87      0.266    10.8   132. <0.0001  2.34     3.40
## 2 age         0.00663   0.0111    0.598  119.  0.551   -0.0153  0.0286
## 3 smoke      -0.244    0.121    -2.02  74.7  0.047   -0.486   -0.00280
```

The `pool()` function should work in most cases, for most types of models, as long as there is a tidy method for it in the `broom` or `broom.mixed` package (available methods are listed at broom.tidymodels.org). Note that, unlike with usual model results, there are not many additional options after MI has been performed and the results have been pooled. Typical functions like `coef()` and `predict()` are not available.

4 Extended example

Now let's suppose that we want to a) give each level of smoke a label instead of using 0/1 coding, b) consider age as categorical variable, and c) recode `ftv` by grouping together 2-6 visits, and include it in the new model.

Here, we could impute the missing data as above, and then calculate the variables we need for our regression models. The approach in this case would be as follows:

¹Changing the default colors using `col` and the shapes using `pch` in `stripplot()` makes these easier to read.

1. Recode smoke as a factor.
2. Impute age, bwt and ftv as before. The computed variables age_cut and ftv2 are not included in the dataset at this point.
3. Calculate age_cut and ftv2.

Therefore the following three options should be examined more critically before running `mice()` again:

m The number of imputed datasets (default `m = 5`). One good approach to determining `m` is to check the fraction of missing information (`fmi`) in model results, and use 100 times the highest value. For this data, the highest `fmi` is 0.29, indicating about 30 imputations. See Section 8.3.

```
check1 <- with(imp, lm(bwt ~ age + smoke + race + ftv))
pool(check1)$pooled
```

##	term	m	estimate	ubar	b	t	dfcom	df	riv	lambda	fmi
## 1	(Intercept)	5	3.4566	0.07576	0.0041691	0.08076	183	146	0.066	0.062	0.075
## 2	age	5	-0.0053	0.00011	0.0000079	0.00012	183	134	0.084	0.077	0.091
## 3	smoke	5	-0.4105	0.01271	0.0022498	0.01541	183	70	0.212	0.175	0.198
## 4	raceblack	5	-0.4728	0.02435	0.0009121	0.02545	183	160	0.045	0.043	0.055
## 5	raceother	5	-0.5265	0.01518	0.0019734	0.01755	183	91	0.156	0.135	0.153
## 6	ftv	5	0.0178	0.00240	0.0007079	0.00325	183	41	0.354	0.261	0.295

method Using the default settings of `mice()`, the imputation method for each of the variables except race will be predictive mean matching `pmm`, since all variables in the dataset are numeric. If we recode smokes as a factor, the default method would then be `logreg`, logistic regression. See also Section 8.7.

predictorMatrix By default, the imputation model for each variable includes all other variables. Reading across the rows below, 1 means the predictor in that column is included in the imputation model for that row, else 0 means it is not included. The default seems to be acceptable here, since there are no variables which are in the dataset twice (e.g. age and age_cut) and no variables which are calculated from other variables in the dataset (e.g. if height, weight and bmi were all present).

First we fix the smoking variable, and then get the default settings from `mice()` using the option `maxit = 0`.

```
dat2 <- dat[, c("bwt", "age", "smoke", "race", "ftv")]
dat2$smoke <- factor(dat2$smoke, 0:1, c("non-smoker", "smoker"))

init <- mice(dat2, maxit = 0)
```

If, for example, we want to change the imputation method for age, and we don't want to use `ftv` to predict race, we could add the following lines.

```

meth <- init$method
meth["age"] <- "norm"
meth

##      bwt      age      smoke      race      ftv
##      "pmm"    "norm"  "logreg" "polyreg"  "pmm"

pred <- init$predictorMatrix
pred["race", "ftv"] <- 0
pred

##      bwt age smoke race ftv
## bwt    0  1   1   1   1
## age    1  0   1   1   1
## smoke  1  1   0   1   1
## race   1  1   1   0   0
## ftv    1  1   1   1   0

```

We then impute the missing observations, and then calculate the derived variables we want from age, birthweight, and ftv. This happens in 3 steps:

1. To get imputed data that look like a regular dataset, generate `complete()` imputed data,

```

imp2 <- mice(dat2, m = 30,
             method = meth, predict = pred,
             print = FALSE, seed = 123456)
impc <- complete(imp2, "long", include = TRUE)

```

2. derive any additional variables,

(a) age

```

impc <- impc %>%
  mutate(age_cut = cut(age, c(14, 20, 25, 30, 50),
                       include.lowest = TRUE, right = FALSE),
         age_cut = relevel(age_cut, "[20,25)"))

```

(b) number of physician visits

```

impc <- impc %>%
  mutate(ftv2 = case_when(ftv > 2 ~ 2,
                          .default = ftv),
         ftv2 = factor(ftv2, 2:0, c("2+ visits", "1 visit", "0 visits")))

```

3. and finally, declare the imputed data to be `mids` again using `as.mids()`. This is the format `mice` is expecting to use for any regression analyses.

```
impc <- as.mids(impc)
```

Such an approach could be used, for example, if height and weight measurements are available, but BMI should be included in the models.

We can then run any models with the derived variables as described above.

```
m3.mi <- with(impc, lm(bwt ~ age_cut + smoke + ftv2))
summary(pool(m3.mi), conf.int = TRUE)
```

```
## # A tibble: 7 x 8
##   term          estimate std.error statistic    df p.value conf.low conf.high
##   <fct>          <dbl>    <dbl>    <dbl> <dbl> <chr>    <dbl>    <dbl>
## 1 (Intercept)    2993.     159.     18.9   133. <0.0001  2679.    3307.
## 2 age_cut[14,20]  128.      146.      0.878  139. 0.382    -160.     416.
## 3 age_cut[25,30] -143.      158.     -0.901  122. 0.369    -456.     171.
## 4 age_cut[30,50]  180.      198.      0.910  122. 0.364    -212.     572.
## 5 smokesmoker   -245.     117.     -2.09   146. 0.038    -476.    -13.5
## 6 ftv21 visit    120.      169.      0.708  132. 0.480    -215.     454.
## 7 ftv20 visits  -70.2     153.     -0.459  131. 0.647    -372.     232.
```

There are other approaches to handling derived variables, for example, using what is often referred to as "passive imputation". That approach may be useful in more complicated situations. See [FIMD Section 6.4](#).

5 Cox proportional hazards regression

It is recommended to include two variables related to the survival endpoint in the imputation models, the Nelson-Aalen estimate of the cumulative hazard (`nelsonaalen()`) and the event indicator, in the imputation process. For more details, see [Section 8.5](#).

```
library(survival)
md.pattern(stanford2, plot = FALSE)
```

```
##   id time status age t5
## 157  1   1     1   1  1  0
##  27  1   1     1   1  0  1
##    0   0     0   0  27 27
```

```
stanford2$nelsonaalen <- nelsonaalen(stanford2, time, status)
```

```
imp.surv <- mice(stanford2, m = 20, print = FALSE)
m4.mi <- with(imp.surv, coxph(Surv(time, status) ~ t5 + age))
summary(pool(m4.mi), conf.int = TRUE, exponentiate = TRUE)
```

```
## # A tibble: 2 x 8
##   term estimate std.error statistic    df p.value conf.low conf.high
##   <fct>    <dbl>    <dbl>    <dbl> <dbl> <chr>    <dbl>    <dbl>
## 1 t5       1.18     0.180     0.929  98.8 0.3554    0.827    1.69
## 2 age     1.03     0.0106    2.72  109. 0.0077    1.01     1.05
```

Since we used the `exponentiate = TRUE` option, the column labelled `estimate` shows the hazard ratios.

6 Multiple imputation with clustered or longitudinal data

Multiple imputation on longitudinal or clustered data is a difficult problem. In the case of longitudinal data, the easiest approach is often to impute the data in wide format (one row per patient, with each timepoint in a different column), convert the data to long (one row per patient and timepoint), and then perform the regression analysis. An example of this type of analysis is shown below. For cluster randomized data where the random effects (for example, schools or child care centers) should be included in the imputation process, another approach should be used, see [Chapter 11 in FIMD](#). For more examples, see [Chapters 7 and 11](#). Below, we briefly show 2 approaches: 1) impute wide format, reshape and model in long format; and 2) impute long format and model in long format using "21" methods in mice package.

Potthoff-Roy data: approach 1 (impute wide, reshape, analyse) In the data from Potthoff and Roy (1964), we have repeated measures from 4 timepoints, in wide format, to which we have added some missing values. In the usual complete cases analysis, we first reshape the data to long, plot the observations, and run a linear mixed model using `lmer()`². The outcome, distance, appears to increase in a statistically significant fashion over time.

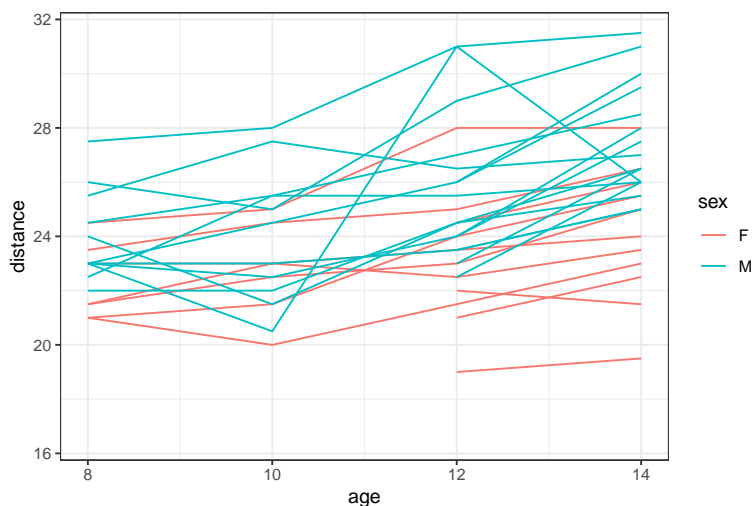
```
library(lme4)
library(broom.mixed)
head(phr, 3)

##   id sex d8 d10 d12 d14
## 1  1  F 21  20  22  23
## 2  2  F 21  22  24  26
## 3  3  F 20  NA  24  26

phr_long <- phr %>%
  pivot_longer(d8:d14,
               values_to = "distance",
               names_to = "age",
               names_pattern = "d(.*)",
               names_transform = as.numeric)

ggplot(aes(age, distance, group = id,
           color = sex), data = phr_long) +
  geom_line() +
  theme_bw()
```

²Using the R package `lmerTest` adds *p*-values to the usual `lme4` results, but causes a problem with `pool()`.



```
m5.cc <- lmer(distance ~ age + sex + (1 | id), data = phr_long)
broom.mixed::tidy(m5.cc, conf.int = TRUE)
```

```
## # A tibble: 5 x 8
##   effect  group  term          estimate std.error statistic  conf.low  conf.high
##   <chr>   <chr>  <chr>         <dbl>    <dbl>    <dbl>    <dbl>    <dbl>
## 1 fixed   <NA>   (Intercept)  15.3     0.936     16.3     13.4     17.1
## 2 fixed   <NA>   age           0.666    0.0646    10.3     0.539    0.792
## 3 fixed   <NA>   sexM          2.39     0.779     3.07     0.863    3.92
## 4 ran_pars id      sd__(Interc~ 1.83     NA        NA       NA       NA
## 5 ran_pars Residual sd__Observa~ 1.48     NA        NA       NA       NA
```

We follow the same basic approach with multiple imputation. This is the same approach used in previous examples, but we add an extra step to reshape the data between the imputation step and the analysis step.

```
mice(phr, m = 10, print = FALSE) %>%
  mice::complete("all") %>%
  map(~ pivot_longer(., d8:d14,
    values_to = "distance",
    names_to = "age",
    names_pattern = "d(.*)",
    names_transform = as.numeric)) %>%
  map(lmer, formula = distance ~ age + sex + (1 | id)) %>%
  pool() %>%
  summary(conf.int = TRUE)
```

```
## # A tibble: 3 x 8
##   term          estimate std.error statistic    df p.value  conf.low  conf.high
##   <fct>         <dbl>    <dbl>    <dbl> <dbl> <chr>    <dbl>    <dbl>
## 1 (Intercept)  15.3     0.922     16.6  100. <0.0001  13.5     17.1
## 2 age           0.668    0.0648    10.3  100. <0.0001  0.539    0.796
## 3 sexM          2.26     0.757     2.98  100. 0.0036  0.755    3.76
```

Pothoff-Roy data: approach 2 (impute long, analyse) In an alternate approach, we try to impute and analyze the data in long format, using random effects. The method is then 2l.pmm (PMM for 2-level data, from the miceadds package [Robitzsch et al., 2019]), and we mark id as the cluster variable with -2 in the predictor matrix. This would also be an appropriate way to include centers in the imputation procedure for multicenter RCTs.

```
library(miceadds)
imp0 <- mice(phr_long, maxit = 0)

meth <- imp0$meth
meth["distance"] <- "2l.pmm"
meth

##      id      sex      age distance
##      ""      ""      ""  "2l.pmm"

pred <- imp0$pred
pred[, "id"] <- -2
pred

##      id sex age distance
## id      -2  1  1         1
## sex      -2  0  1         1
## age      -2  1  0         1
## distance -2  1  1         0

mice(phr_long, m = 10,
      predictorMatrix = pred, method = meth,
      print = FALSE) %>%
  complete("all") %>%
  map(lmer, formula = distance ~ age + sex + (1 | id)) %>%
  pool() %>%
  summary(conf.int = TRUE)

## # A tibble: 3 x 8
##   term      estimate std.error statistic    df p.value  conf.low  conf.high
##   <fct>      <dbl>    <dbl>    <dbl> <dbl> <chr>    <dbl>    <dbl>
## 1 (Intercept)  15.3      0.926     16.5  98.4 <0.0001  13.4     17.1
## 2 age          0.668    0.0637    10.5  99.7 <0.0001  0.542    0.794
## 3 sexM         2.32     0.778      2.99  99.6 0.0036  0.779    3.87
```

In some cluster-randomized trials, it may be that the clusters are relatively unimportant, and may be reasonably ignored in the imputation procedure (see for example, [FIMD Section 7.3.2](#)). In the case of truly multilevel data, the micemd package may be helpful. It is based on the work of [Audigier et al.](#), and discusses which methods should be used depending on the data structure.

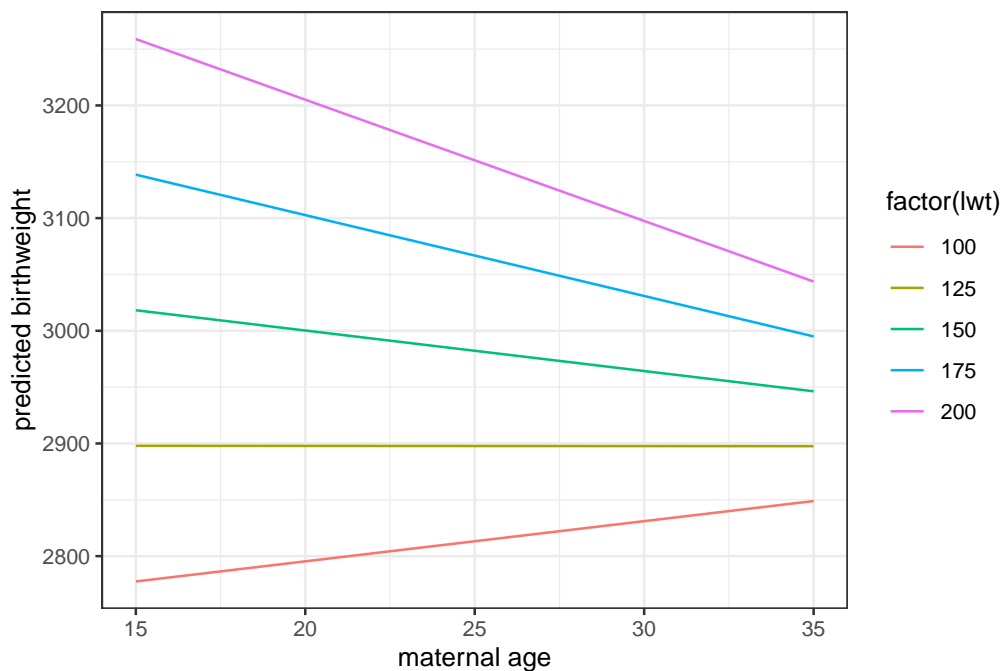
7 Predictions / fitted values after MI

Can I get fitted values from my pooled results? There are differing approaches to how to compute the fitted values. Do we 1) pool model coefficients and then get fitted values, or 2) get fitted values first and then pool them? Miles [2016] suggests that the order does not matter, and that the second approach is easier. Since mice version 3.19, there is a new function `predict_mi()` to obtain such predictions which appears to use the 2nd approach. By default, it pools the predictions, but the option `pool = FALSE` can be used to obtain predictions for each imputation.

```
imp <- mice(dat, m = 3, print = FALSE)
# 1. get model results
mod <- with(imp, lm(bwt ~ I(age - 20) * I(lwt - 125)))
# 2. create new data
nd <- crossing(bwt = 3000, # the outcome can be any non-missing value
              age = seq(15, 35, 5),
              lwt = seq(100, 200, 25))
# 3. use predict_mi() to get predictions
pred <- predict_mi(mod, nd, interval = "confidence")
nd$pred <- pred[, "fit"]
nd$lb <- pred[, "lwr"]
nd$sub <- pred[, "upr"]
nd

## # A tibble: 25 x 6
##   bwt   age   lwt  pred    lb    ub
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 3000    15   100 2778. 2212. 3343.
## 2 3000    15   125 2898. 2661. 3135.
## 3 3000    15   150 3018. 2746. 3291.
## 4 3000    15   175 3139. 2566. 3712.
## 5 3000    15   200 3259. 2324. 4193.
## 6 3000    20   100 2795. 2557. 3034.
## 7 3000    20   125 2898. 2775. 3020.
## 8 3000    20   150 3000. 2800. 3200.
## 9 3000    20   175 3103. 2726. 3479.
## 10 3000    20   200 3205. 2635. 3775.
## # i 15 more rows

ggplot(aes(age, pred, color = factor(lwt)), data = nd) +
  geom_line() +
  xlab("maternal age") + ylab("predicted birthweight") +
  theme_bw()
```



8 Frequently Asked Questions

8.1 Do I need to include imputed values in Table 1?

Multiple imputation is only for model fitting, and should not be used for tables of patient characteristics (which ideally should show how many missing values there are, see example below). MI is also not appropriate for use in model choice, in most cases at least, as measures of goodness-of-fit like AIC cannot be pooled. Use it only when you have a final model or a small group of final models, if you would like a "sensitivity analysis" to see how robust the coefficients in the complete cases analysis are to the missing measurements or observations.

```
library(gtsummary)
small %>%
  tbl_summary(missing = "ifany")
```

8.2 How do I report my MI analysis?

Beyond what can be found in the relevant **CONSORT** (randomized trials) or **STROBE** (observational studies) guidelines, Box 2 in [Sterne et al. \[2009\]](#) has a good description of how results should be reported if MI was used in the analysis.

8.3 How many imputations?

[Madley-Dowd et al. \[2019\]](#) has suggested that considering fraction of missing information (fmi) is more appropriate than proportion of incomplete cases (as suggested by [White et al. \[2011\]](#), Section 7.3) among others) when determining the number of imputations. However, fmi cannot be checked until at least some imputations have been generated and a model has been analyzed. Given that, the following procedure could be employed:

1. Set up your imputation procedure, leaving $m = 5$.
2. Check that the imputation runs through without errors.
3. Run a model which includes all the variables you want to include, and pool the results.
4. Print `summary()` of the pooled model results, and inspect the `fmi` column.
5. Round the highest `fmi` up to the nearest 0.05 (e.g., round 0.236 to 0.25) and multiply that value by 100 to get a new value of m (e.g., $m = 25$).
6. Run `mice()` again and continue with your analysis.

A two-stage approach for determining the number of imputations was proposed by [von Hippel \[2018\]](#) imputation, based on `fmi` and between imputation variability (implemented in the R package `howManyImputations`). For the previous example, the function can be used as follows:

```
#install.packages("howManyImputations")
library(howManyImputations)
how_many_imputations(mod, cv = 0.05)

## [1] 179
```

8.4 Which variables should be included in the imputation step?

[White et al. \[2011, Section 5\]](#) recommend using covariates and the outcome from the analysis models you want to run, as well as predictors of the incomplete variable(s). In general, it's good practice to create a small analysis dataset containing only the variables you need for the analysis, and use that in `mice()`, rather than the full dataset. If you plan on using any interactions in your analysis model(s), these should also be included in the imputation model (see [FIMD Section 6.4.2](#)).

8.5 How can I use `mice()` when performing Cox PH regression?

[White et al. \[2011, Section 5\]](#) recommend using covariates and the outcome from the analysis models, as well as predictors of the incomplete variable. Further, [White and Royston \[2009\]](#) recommend using the

- the Nelson-Aalen estimate of the cumulative hazard (computed using `nelsonaalen()`), and
- the event indicator (for example, `died` as a 0/1 variable).

See [Section 5](#).

8.6 Can I show odds ratios after logistic regression?

Yes, try the option `exponentiate = TRUE`.

```
imp %>%
  complete("all") %>%
  map(glm, formula = I(bwt < 2500) ~ age + race + smoke,
      family = binomial) %>%
  pool() %>%
  summary(exponentiate = TRUE, conf.int = TRUE)
```

```
## # A tibble: 5 x 8
##   term          estimate std.error statistic    df p.value conf.low conf.high
##   <fct>         <dbl>     <dbl>    <dbl> <dbl> <chr>      <dbl>    <dbl>
## 1 (Intercept)    0.149     1.02     -1.86   28.6 0.0729    0.0183    1.21
## 2 age            1.00     0.0375    0.0381  32.5 0.9699    0.928     1.08
## 3 raceblack      3.00     0.492     2.23   177. 0.0267    1.14     7.92
## 4 raceother      3.48     0.456     2.73   42.3 0.0091    1.39     8.73
## 5 smoke          3.15     0.424     2.71   30.9 0.0110    1.33     7.48
```

8.7 Which method do I use to impute the variables?

By default, `mice()` uses the following default methods (option `defaultMethod`): **predictive mean matching** (`pmm`) for numeric data, logistic regression for factors with 2 levels, and multinomial logistic regression for factors with 3 levels. Note that binary variables that are still coded numerically (0/1, 1/2, etc) will have the default method "pmm", unless you recode it as a factor. There are many more available methods. See also multilevel methods available from the `micemd` package, and [Chapter 3 in FIMD](#) for discussion of the various methods..

```
apropos("mice.impute")
```

```
## [1] "mice.impute.2l.bin"           "mice.impute.2l.lmer"
## [3] "mice.impute.2l.norm"         "mice.impute.2l.pan"
## [5] "mice.impute.2lonly.mean"     "mice.impute.2lonly.norm"
## [7] "mice.impute.2lonly.pmm"      "mice.impute.cart"
## [9] "mice.impute.jomoImpute"      "mice.impute.lasso.logreg"
## [11] "mice.impute.lasso.norm"      "mice.impute.lasso.select.logreg"
## [13] "mice.impute.lasso.select.norm" "mice.impute.lda"
## [15] "mice.impute.logreg"          "mice.impute.logreg.boot"
## [17] "mice.impute.mean"            "mice.impute.midastouch"
## [19] "mice.impute.mnar.logreg"     "mice.impute.mnar.norm"
## [21] "mice.impute.mpmm"            "mice.impute.norm"
## [23] "mice.impute.norm.boot"       "mice.impute.norm.nob"
## [25] "mice.impute.norm.predict"    "mice.impute.panImpute"
## [27] "mice.impute.passive"         "mice.impute.pmm"
## [29] "mice.impute.polr"            "mice.impute.polyreg"
## [31] "mice.impute.quadratic"       "mice.impute.rf"
## [33] "mice.impute.ri"              "mice.impute.sample"
```

8.8 In my analysis, pool doesn't seem to work. Now what?

The `mice` package should be able to pool any model results for which a tidy method is available (`?generics::tidy`). If there is none, `?pool` suggests using `pool.scalar` or writing your own tidy method (see also [this overview on adding tidiers](#)).

```

imp <- mice(dat, m = 25, print = FALSE, seed = 12345)
nrow(dat) # n = 189

## [1] 189

mod <- with(imp, lm(bwt ~ smoke + age))
summary(pool(mod), conf.int = TRUE)

## # A tibble: 3 x 8
##   term          estimate std.error statistic    df p.value  conf.low  conf.high
##   <fct>          <dbl>    <dbl>    <dbl> <dbl> <chr>    <dbl>    <dbl>
## 1 (Intercept)  2876.    269.    10.7   141. <0.0001  2344.    3408.
## 2 smoke       -255.    114.    -2.24  155. 0.026   -480.    -30.5
## 3 age          6.07    11.4     0.531  122. 0.597   -16.6    28.7

coef(mod$analyses[[1]])

## (Intercept)      smoke      age
##      2791.6      -248.6       8.8

length(coef(mod$analyses[[1]])) # k = 3

## [1] 3

mod$analyses %>%
  map(tidy) %>%
  bind_rows() %>%
  mutate(variance = std.error ^ 2) %>%
  group_by(term) %>%
  nest() %>%
  mutate(out = map(data, ~ pool.scalar(.estimate, .variance,
                                     n = 189, k = 3)),
         estimate = map_dbl(out, ~ .$qbar),
         std.error = map_dbl(out, ~ sqrt(.t)),
         statistic = estimate / std.error,
         df = map_dbl(out, ~ .$df),
         p.value = 2 * (pt(abs(statistic), df, lower.tail = FALSE)),
         zalpha = qt(1 - 0.05 / 2, df),
         conf.low = estimate - zalpha * std.error,
         conf.high = estimate + zalpha * std.error ) %>%
  select(-data, -out) %>%
  mutate(p.value = format.pval(p.value, eps = 0.0001,
                              sci = FALSE, digits = 2))

## # A tibble: 3 x 9
## # Groups:   term [3]
##   term          estimate std.error statistic    df p.value  zalpha  conf.low  conf.high
##   <chr>          <dbl>    <dbl>    <dbl> <dbl> <chr>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)  2876.    269.    10.7   141. <0.0001  1.98   2344.    3408.
## 2 smoke       -255.    114.    -2.24  155. 0.026   1.98   -480.    -30.5
## 3 age          6.07    11.4     0.531  122. 0.6     1.98   -16.6    28.7

```

8.9 Can I pool quantities that are not regression coefficients?

Maybe! See [Marshall et al. \[2009\]](#) and, possibly, the function `pool.scalar()`. For example, it has been suggested that AUC (area under the curve, c-statistic) may be pooled using robust methods [[Marshall et al., 2009](#), Table 1]. Since we will estimate the confidence interval using quantiles of AUC from the imputed datasets, we increase the number of imputations to 200.

```
myroc <- function(x, fm){
  fit <- glm(fm, data = x, family = binomial)
  obs <- x$low
  pred <- predict(fit)
  as.numeric(auc(roc(obs ~ pred, quiet = TRUE)))
}

imp <- mice(dat, m = 200, print = FALSE, seed = 12345)

res <- complete(imp, "long") %>%
  mutate(low = as.numeric(bwt < 2500)) %>%
  group_by(.imp) %>%
  nest() %>%
  mutate(model1 = map_dbl(data, myroc,
                          fm = low ~ smoke),
         model2 = map_dbl(data, myroc,
                          fm = low ~ smoke + age),
         model3 = map_dbl(data, myroc,
                          fm = low ~ smoke + race)) %>%
  select(.imp, model1:model3)

res %>%
  pivot_longer(model1:model3) %>%
  group_by(name) %>%
  summarize(auc = median(value),
           lb = quantile(value, 0.05 / 2),
           ub = quantile(value, 1 - 0.05 / 2))
```

8.10 Can I combine MI with cross-validation or bootstrap sampling?

Generally, it is easiest if MI and resampling methods do not need to be combined, but sometimes it is not possible to avoid it. [Schomaker and Heumann \[2018\]](#) and [Wahl et al. \[2016\]](#) both explored different combinations of bootstrap sampling / cross-validation and MI. They generally suggest that a reasonable approach is to first get bootstrap samples (or for k -fold cross-validation, to split into testing and testing samples), and then to perform MI. Feel free to set up an appointment with me if you need to do this.

8.11 Can I use parallel processing to run mice?

Starting with version 3.15.0 of the `mice` package, a new function `futuremice` takes advantage of parallel process features provided by the `future` and `furrr` packages [[Bengtsson, 2021](#), [Vaughan](#)

and Dancho, 2022]. See also the vignette for `futuremice`, Vink and van Buuren [2019].

For reproducible results, it is recommended to set the option `parallelseed` (not `seed`!). It should be noted that while `future.plan = "multisession"` is the only option available in Windows, `future.plan = "multicore"` could be used with other operating systems (see `?future::plan`). These options are not available while running R interactively in RStudio, so the times given below are from a run of the same code using `Rscript` in batch mode in the RStudio Terminal (`?Rscript`).

```
library(tictoc) # tictoc reports processing time

tic()
imp_seq <- mice(dat, m = 200, seed = 12345, printFlag = FALSE)
toc()
# 8.097 sec elapsed
tic()
imp_parallel <- futuremice(dat, m = 200, n.core = 6,
                          parallelseed = 12345,
                          future.plan = "multisession")
toc()
# 3.292 sec elapsed
parallel::detectCores()
# [1] 14
tic()
imp_parallel <- futuremice(dat, m = 200, n.core = 6,
                          parallelseed = 12345,
                          future.plan = "multicore")
toc()
# 1.427 sec elapsed

fit <- with(imp_parallel, lm(bwt ~ age * smoke + ftv))
summary(pool(fit), conf.int = TRUE) %>%
  mutate(p.value = format.pval(p.value, eps = 0.0001))
#           term estimate std.error statistic  df p.value conf.low conf.high
# 1 (Intercept)   2.441    0.311     7.85 162 <1e-04  1.82663  3.0553
# 2           age   0.027    0.013     2.01 160   0.05  0.00048  0.0531
# 3          smoke  0.575    0.551     1.04 130   0.30 -0.51498  1.6647
# 4           ftv   0.038    0.061     0.63 137   0.53 -0.08208  0.1580
# 5    age:smoke  -0.041    0.024    -1.73 122   0.09 -0.08869  0.0059
```

9 References

V Audigier, I R White, S Jolani, T P A Debray, M Quartagno, J Carpenter, S van Buuren, and M Resche-Rigon. Multiple imputation for multilevel data with continuous and binary variables. 33(2):160–183. doi:[10.1214/18-STS646](https://doi.org/10.1214/18-STS646).

- Henrik Bengtsson. A unifying framework for parallel and distributed processing in r using futures. *The R Journal*, 13(2):208–227, 2021. doi:[10.32614/RJ-2021-048](https://doi.org/10.32614/RJ-2021-048). URL <https://doi.org/10.32614/RJ-2021-048>.
- O Harel, EM Mitchell, NJ Perkins, SR Cole, EJ Tchetgen Tchetgen, BL Sun, and EF Schisterman. Multiple Imputation for Incomplete Data in Epidemiologic Studies. *American Journal of Epidemiology*, 187(3):576–584, 11 2017. ISSN 0002-9262. doi:[10.1093/aje/kwx349](https://doi.org/10.1093/aje/kwx349).
- P Madley-Dowd, R Hughes, K Tilling, and J Heron. The proportion of missing data should not be used to guide decisions on multiple imputation. *Journal of Clinical Epidemiology*, 110:63 – 73, 2019. ISSN 0895-4356. doi:[10.1016/j.jclinepi.2019.02.016](https://doi.org/10.1016/j.jclinepi.2019.02.016).
- A Marshall, DG Altman, RL Holder, and Pk Royston. Combining estimates of interest in prognostic modelling studies after multiple imputation: current practice and guidelines. *BMC Medical Research Methodology*, 9(1):57, Jul 2009. ISSN 1471-2288. doi:[10.1186/1471-2288-9-57](https://doi.org/10.1186/1471-2288-9-57).
- A Miles. Obtaining predictions from models fit to multiply imputed data. *Sociological Methods & Research*, 45(1):175–185, 2016. doi:[10.1177/0049124115610345](https://doi.org/10.1177/0049124115610345).
- AB Pedersen, EM Mikkelsen, D Cronin-Fenton, NR Kristensen, TM Pham, L Pedersen, and I Petersen. Missing data and multiple imputation in clinical epidemiological research. *Clinical Epidemiology*, 2017(9):157–166, 2017. doi:[10.2147/CLEP.S129785](https://doi.org/10.2147/CLEP.S129785).
- NJ Perkins, SR Cole, O Harel, EJ Tchetgen Tchetgen, BL Sun, EM Mitchell, and EF Schisterman. Principled Approaches to Missing Data in Epidemiologic Studies. *American Journal of Epidemiology*, 187(3):568–575, 11 2017. ISSN 0002-9262. doi:[10.1093/aje/kwx348](https://doi.org/10.1093/aje/kwx348).
- A Robitzsch, S Grund, and T Henke. *miceadds: Some additional multiple imputation functions, especially for mice*, 2019. URL <https://CRAN.R-project.org/package=miceadds>. R package version 3.3-33.
- M Schomaker and C Heumann. Bootstrap inference when using multiple imputation. *Statistics in Medicine*, 37(14):2252–2266, 2018. doi:[10.1002/sim.7654](https://doi.org/10.1002/sim.7654).
- JAC Sterne, IR White, JB Carlin, M Spratt, P Royston, MG Kenward, AM Wood, and JR Carpenter. Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls. *BMJ*, 338, 2009. ISSN 0959-8138. doi:[10.1136/bmj.b2393](https://doi.org/10.1136/bmj.b2393).
- S van Buuren. *Flexible Imputation of Missing Data*. 2nd edition, 2018. URL <https://stefvanbuuren.name/fimd>.
- S van Buuren and K Groothuis-Oudshoorn. mice: Multivariate imputation by chained equations in R. *Journal of Statistical Software*, 45, 2011. doi:[10.18637/jss.v045.i03](https://doi.org/10.18637/jss.v045.i03).
- Davis Vaughan and Matt Dancho. *furrr: Apply Mapping Functions in Parallel using Futures*, 2022. URL <https://CRAN.R-project.org/package=furrr>. R package version 0.3.1.
- G Vink and S van Buuren. *micevignettes*, 2019. URL <https://www.gerkovink.com/miceVignettes/>.
- Paul von Hippel. How many imputations do you need? a two-stage calculation using a quadratic rule. *Sociological Methods & Research*, 2018. doi:[10.1177/0049124117747303](https://doi.org/10.1177/0049124117747303). URL <https://statisticalhorizons.com/how-many-imputations>.
- S Wahl, A-L Boulesteix, A Zierer, B Thorand, and MA van de Wiel. Assessment of predictive performance in incomplete data by combining internal validation and multiple imputation. *BMC Medical Research Methodology*, 16(1):144, Oct 2016. ISSN 1471-2288. doi:[10.1186/s12874-016-0239-7](https://doi.org/10.1186/s12874-016-0239-7).

IR White and P Royston. Imputing missing covariate values for the Cox model. *Statistics in Medicine*, 28(15):1982–1998, 2009. ISSN 1097-0258. doi:[10.1002/sim.3618](https://doi.org/10.1002/sim.3618).

IR White, P Royston, and AM Wood. Multiple imputation using chained equations: Issues and guidance for practice. *Statistics in Medicine*, 30(4):377–399, 2011. ISSN 1097-0258. doi:[10.1002/sim.4067](https://doi.org/10.1002/sim.4067).

A Recent Changes

Version 1.0 Reworked the original "Multiple Imputation using Stata" document for R.

Version 1.1 Added the [von Hippel \[2018\]](#) method of determining an appropriate number of imputations. Added some comments on reporting analysis that uses MI.

Version 1.2 Removed the hierarchical data example using `hmi`, which is no longer on CRAN. Changed the code slightly to reflect the use of `pivot_longer` rather than `reshape`.

Version 1.3 Added note about `micemd` package for multiple imputation with multilevel data.

Version 1.4 Added `predict_mi` function, and use of parallel processing with `futuremice`. Added example of `pool.scalar` to the FAQs. Rewrote code to generally use `tidverse` (`dplyr`, `purrr`, etc) syntax. We now use slightly modified versions of `summary.mipo` and `tidy.lm` to clean up the output:

```
# to remove when this issue is fixed: https://github.com/amices/mice/issues/719
summary <- function(x, ...){
  mice::summary.mipo(x, ...) %>%
  dplyr::select(-'2.5 %', -'97.5 %') %>%
  dplyr::mutate(p.value = format.pval(p.value, eps = 0.0001,
                                     sci = FALSE, digits = 2)) %>%
  as_tibble()
}

tidy <- function(x, ...){
  out <- broom::tidy(x, ...)
  if("p.value" %in% names(out)){
    out <- out %>%
    dplyr::mutate(p.value = format.pval(p.value, eps = 0.0001,
                                       sci = FALSE, digits = 2))
  }
  out
}
```

B Code to make datasets used here

```

library(mice)
set.seed(987654)
dat0 <- MASS::birthwt[, c("bwt", "age", "lwt", "race",
                          "smoke", "ptl", "ht", "ui", "ftv")]
patt <- ampute(dat0)$pattern
patt <- rbind(patt,
             c(1, 1, 0, 1, 1, 0, 0, 0, 0),
             c(0, 1, 0, 1, 1, 0, 0, 0, 0),
             c(0, 1, 0, 0, 0, 1, 1, 1, 1))
dat <- ampute(dat0, pattern = patt)$amp
dat$race <- factor(dat$race, 1:3, c("white", "black", "other"))
phr <- potthoffroy
idmis <- c(3,6,9,10,13,16,23,24,27)
phr[idmis, 4] <- NA

```

```

sessionInfo()

## R version 4.5.2 (2025-10-31)
## Platform: aarch64-apple-darwin20
## Running under: macOS Tahoe 26.3.1
##
## Matrix products: default
## BLAS: /System/Library/Frameworks/Accelerate.framework/Versions/A/Frameworks/vecLib.framework
## LAPACK: /Library/Frameworks/R.framework/Versions/4.5-arm64/Resources/lib/libRlapack.dylib; I
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## time zone: Europe/Zurich
## tzcode source: internal
##
## attached base packages:
## [1] stats graphics grDevices utils datasets methods base
##
## other attached packages:
## [1] howManyImputations_0.2.5 gtsummary_2.5.0 broom.mixed_0.2.9.6
## [4] lme4_2.0-1 Matrix_1.7-4 survival_3.8-3
## [7] broom_1.0.12 lubridate_1.9.4 forcats_1.0.1
## [10] stringr_1.6.0 dplyr_1.2.1 purrr_1.2.1
## [13] readr_2.1.6 tidyr_1.3.2 tibble_3.3.1
## [16] ggplot2_4.0.2 tidyverse_2.0.0 mice_3.19.0
## [19] knitr_1.50
##
## loaded via a namespace (and not attached):
## [1] tidymodels_1.2.1 farver_2.1.2 S7_0.2.1 fastmap_1.2.0
## [5] digest_0.6.39 rpart_4.1.24 timechange_0.3.0 lifecycle_1.0.5
## [9] magrittr_2.0.4 compiler_4.5.2 rlang_1.2.0 tools_4.5.2
## [13] utf8_1.2.6 gt_1.1.0 labeling_0.4.3 xml2_1.4.1

```

```
## [17] RColorBrewer_1.1-3 withr_3.0.2 nnet_7.3-20 grid_4.5.2
## [21] jomo_2.7-6 future_1.70.0 globals_0.19.1 scales_1.4.0
## [25] iterators_1.0.14 MASS_7.3-65 cli_3.6.6 rmarkdown_2.30
## [29] reformulas_0.4.4 generics_0.1.4 rstudioapi_0.17.1 tzdb_0.5.0
## [33] commonmark_2.0.0 minqa_1.2.8 DBI_1.2.3 splines_4.5.2
## [37] parallel_4.5.2 mitools_2.4 vctrs_0.7.3 boot_1.3-32
## [41] glmnet_4.1-10 litedown_0.8 hms_1.1.4 mitml_0.4-5
## [45] listenv_0.10.1 foreach_1.5.2 glue_1.8.0 parallelly_1.46.1
## [49] nloptr_2.2.1 pan_1.9 codetools_0.2-20 stringi_1.8.7
## [53] shape_1.4.6.1 gtable_0.3.6 furr_0.3.1 pillar_1.11.1
## [57] miceadds_3.19-16 htmltools_0.5.8.1 R6_2.6.1 Rdpack_2.6.6
## [61] evaluate_1.0.5 lattice_0.22-7 markdown_2.0 highr_0.11
## [65] rbibutils_2.4.1 backports_1.5.0 cards_0.7.1 Rcpp_1.1.1
## [69] nlme_3.1-168 xfun_0.54 fs_1.6.6 pkgconfig_2.0.3
```