

# Incomplete Data Analysis, 2021/2022

## Multiple imputation of univariate missing data: the **mice** package

School of Mathematics, University of Edinburgh

V. Inácio de Carvalho & M. de Carvalho

Here we will learn how to use the **mice** package in practice. For now, we will only deal with univariate missingness, we will later expand the scope to the case of several variables with missing values. Before proceeding, I leave the reference to the manual of the package

<https://cran.r-project.org/web/packages/mice/index.html>

I will start by simulating some data and then imposing MAR missingness.

```
set.seed(1)
n <- 100
x1 <- runif(n, 0, 5)
x2 <- runif(n, 0, 10)
beta0 <- 5
beta1 <- 3
beta2 <- 1
y <- rnorm(n, beta0 + beta1*x1 + beta2*x2, 1)

x2 <- ifelse(x1 > 4.2, NA, x2)
#checking the percentage of missing values
sum(is.na(x2))/n

## [1] 0.13

#constructing a dataframe with the 3 variables
simdata <- data.frame("y" = y, "x1" = x1, "x2" = x2)
```

The package **mice** has the function `cc` that returns the complete cases. This function is useful when working with real data as it easily allows some exploratory analyses based on the complete cases.

```
require(mice)
cc(simdata)
nrow(cc(simdata))
```

As we have seen back in week 2, **mice** also has a function that allows visualising the missing data patterns.

```
md.pattern(simdata)
```

	y	x1	x2	
87				0
13				1
	0	0	13	13

```
##      y x1 x2
## 87 1  1  1  0
## 13 1  1  0  1
##      0  0 13 13
```

Another function available in `mice` is `md.pairs`, which calculates the number of observations per patterns for all possible pairs of variables. For a pair of variables, there are four possible missing data patterns: both variables are observed (pattern `rr`), the first variable is observed and the second variable is missing (pattern `rm`), the first variable is missing and the second variable is observed (pattern `mr`), and finally the pattern where both variables are missing (pattern `mm`).

```
md.pairs(simdata)
```

```
## $rr
##      y  x1 x2
## y  100 100 87
## x1 100 100 87
## x2  87  87 87
##
## $rm
##      y x1 x2
## y   0  0 13
## x1  0  0 13
## x2  0  0  0
##
## $mr
##      y x1 x2
## y   0  0  0
## x1  0  0  0
## x2 13 13  0
##
## $mm
##      y x1 x2
## y   0  0  0
## x1  0  0  0
## x2  0  0 13
```

Let us now use the package `mice` to impute the values in `x2`. We start with the function `mice()` to perform step 1, i.e., to impute the missing values. We already know that the default in `mice` for continuous variables,

as it is the case of  $x_2$ , is predictive mean matching with  $d = 5$  donors and *Type 1* matching (between the cases with missing values and those with observed values). Also, by default in `mice` we have  $M = 5$ . To know more, type `help(mice)`.

```
imps <- mice(simdata, printFlag = FALSE, seed = 1)
imps
```

```
## Class: mids
## Number of multiple imputations: 5
## Imputation methods:
##      y      x1      x2
##      ""      "" "pmm"
## PredictorMatrix:
##      y x1 x2
## y  0  1  1
## x1 1  0  1
## x2 1  1  0
```

A few comments apply. We set `printFlag = FALSE` which results in silent computation of the missing values and we also use `seed=1` so that our results are reproducible (any other value would obviously work, but fixing the seed outside the function `mice()` will not work). A summary of the imputation results can be obtained by calling the `imps` object. For instance, we see that our saved object `imps` is of class `mids` which stands for *multiply imputed datasets*, which is a special type of object that the `mice` package has set up for storing multiple imputed datasets. We also obtain information about the imputation method used to impute the variables with missing values. In this case only `x2` has missing values and because we have not changed the defaults, unsurprisingly, we have that predictive mean matching was used. Lastly, we have the `predictorMatrix` which, for instance, tell us that `y` and `x1` were used to impute `x2`. It also tells us that in case `y` had missing values, `x1` and `x2` would be used to impute it and similarly for `x1` we would use `y` and `x2`. We can also extract this information from `imps$predictorMatrix`. The default approach in `mice` is to impute one variable based on all other variables.

Now let us look at the imputed values. We can extract them from our `imps` object.

```
imps$imp$x2
```

```
##           1           2           3           4           5
## 4  9.240745  7.410786  9.2861520  9.240745  9.286152
## 6  5.476466  5.980924  4.3147369  5.260277  2.126995
## 7  1.103606  3.287773  3.3548749  2.075451  2.075451
## 18 1.891936  1.891936  0.3554058  2.388687  1.482116
## 21 7.410786  7.410786  9.7617069  9.240745  9.286152
## 29 3.354875  3.354875  1.7344233  2.702601  2.832325
## 52 4.781180  7.293096  3.5672691  5.748722  6.547239
## 61 3.179637  3.179637  5.0044097  3.804939  4.531314
## 70 2.126995  4.525708  3.8049389  3.179637  2.126995
## 76 8.864509  8.770575  8.4061455  7.410786  7.410786
## 77 7.828513  6.304141  6.1464497  8.405070  8.405070
## 80 9.286152  9.286152  8.7705754  9.850952  9.850952
## 94 9.240745  8.770575  9.2861520  7.410786  9.240745
```

The row numbers indicate the record number in the original dataset. We can extract, for instance, the imputed values for the first imputed dataset by simply doing the following:

```
imps$imp$x2[,1]
```

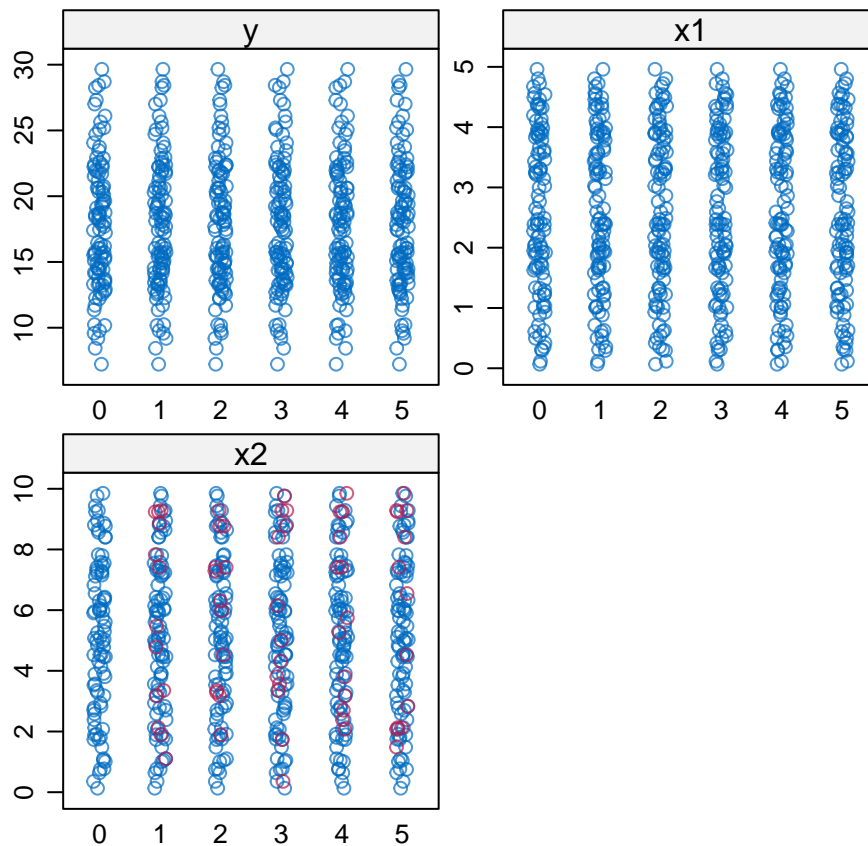
```
## [1] 9.240745 5.476466 1.103606 1.891936 7.410786 3.354875 4.781180 3.179637
## [9] 2.126995 8.864509 7.828513 9.286152 9.240745
```

The (completed) imputed datasets can be extracted by using the `complete` function. As a way of illustrating the usage of this function, I am extracting the first and second completed datasets.

```
com1 <- complete(imps, 1)
com2 <- complete(imps, 2)
com1
com2
```

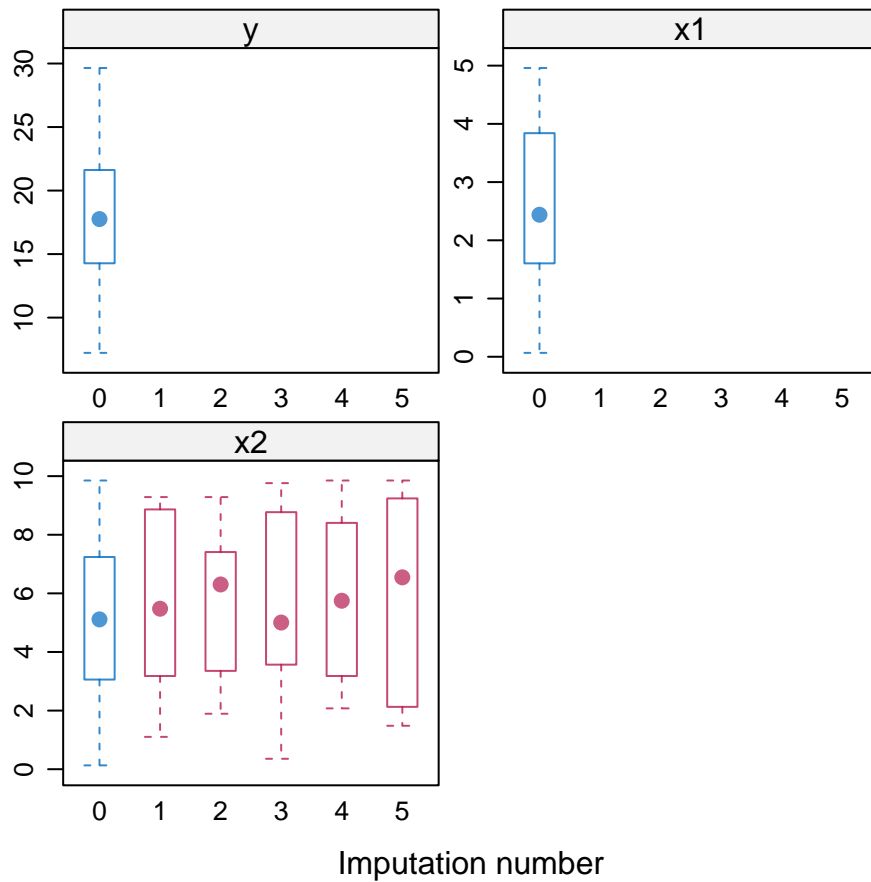
It is also important to visualise the imputation results and the package `mice` provides several plotting tools. This allows us to check whether imputations are plausible. As van Buuren and Groothuis-Oudshoorn say in their paper describing the `mice` package (p. 11): “*Imputations should be values that could have been obtained had they not been missing. Imputations should be close to the data*”. One way to do this is through the `stripplot` function.

```
stripplot(imps)
```



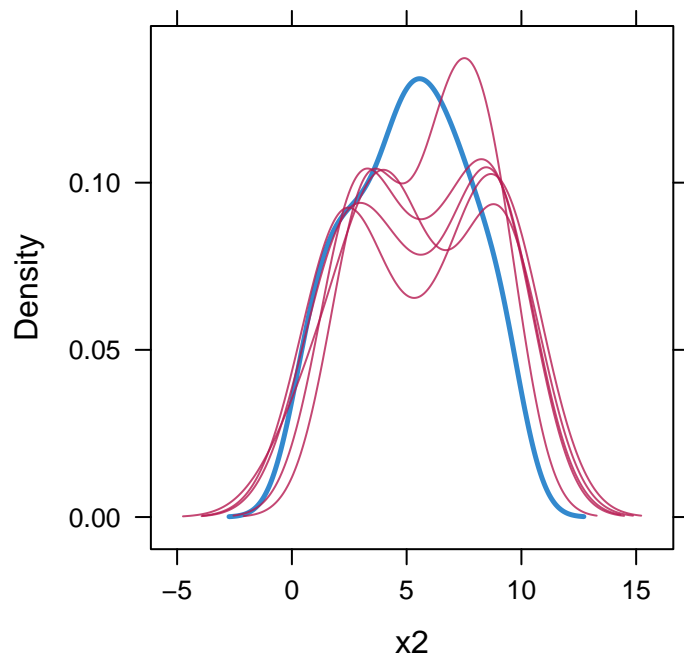
Blue circles denote observed data and red circles imputed data. The panels for `y` and `x1` contain only blue dots because these two variables are fully observed. If there are no large differences between the imputed and observed values then we can conclude that the imputed values are plausible. Here we can see that the red circles follow the blue circles well. If there are discrepancies, interpretation is more difficult, as this may be due to a bad imputation model, due to the missing mechanism not being MAR or due to a combination of both. This plot is most useful when there are not many data points. Alternatively we can use the function `bwplot`, which produces a boxplot of the observed and imputed data.

```
bwplot(imps)
```



There is also the possibility of visualising the kernel density estimates of the observed and imputed data.

```
densityplot(imps)
```



**Aside comment:** note that here the densities assign positive mass to negative  $x_2$  values and we know that this should not be the case ( $x_2$  was simulated from a uniform (0,10) distribution). By looking at the boxplots

in the previous figure, we see that all imputed values are above zero, and so the imputed values are plausible. However, due to the fact that we are using kernel estimates for the densities, which is a nonparametric density estimator, with a reduced sample size, in combination with the fact that the default kernel is the normal one and there values close to zero, leads to mass assigned to negative values.

Adjustments to the defaults used by the predictive mean matching function `mice.impute.pmm` can be made by simply entering the arguments to be altered into the main `mice()` call. They will be automatically passed down to `mice.impute.pmm`. For instance, the number of donors to be sampled from can be set via the `donors` argument. let us now change this argument to three and we will also create  $M = 10$  copies of the dataset (instead of the default  $M = 5$ ).

```
imps_alt <- mice(simdata, m = 10, donors = 3, printFlag = FALSE, seed = 1)
imps_alt
```

```
## Class: mids
## Number of multiple imputations: 10
## Imputation methods:
##      y      x1      x2
##      ""      "" "pmm"
## PredictorMatrix:
##      y x1 x2
## y  0  1  1
## x1 1  0  1
## x2 1  1  0
```

Suppose now that we want to change our method for imputing the missing values. Specifically, suppose that we want to use the method `norm.boot`. There are two possible ways of doing it. The simplest way and feasible only when the number of variables to be imputed is small is to change the method argument directly in the `mice()` call.

```
imps_normb <- mice(simdata, method = "norm.boot", printFlag = FALSE, seed = 1)
imps_normb$imp$x2[,1]
```

```
## [1] 10.194490  4.595104  2.550149  2.953703  9.353378  2.049476  6.779982
## [8]  2.345993  4.452382  9.132930  8.257026  8.847016 10.444725
```

An alternative way is to do a setup run of `mice()` without iterations (`maxit=0`) and to extract and modify the method from there.

```
imps0 <- mice(simdata, maxit = 0)
meth <- imps0$method
meth
```

```
##      y      x1      x2
##      ""      "" "pmm"
```

```
meth["x2"] <- "norm.boot"
imps_norm2 <- mice(simdata, method = meth, printFlag = FALSE,
                  seed = 1)
imps_norm2
```

```
## Class: mids
## Number of multiple imputations: 5
## Imputation methods:
##      y      x1      x2
##      ""      "" "norm.boot"
## PredictorMatrix:
##      y x1 x2
```

```
## y 0 1 1
## x1 1 0 1
## x2 1 1 0
```

```
imps_norm2$imp$x2[,1]
```

```
## [1] 10.194490 4.595104 2.550149 2.953703 9.353378 2.049476 6.779982
## [8] 2.345993 4.452382 9.132930 8.257026 8.847016 10.444725
```

The setup run is also useful to customize our imputation model. Variables in the columns of the `predictorMatrix` can be switched on or off by using a 1 or a 0 to include or exclude them from the imputation model, respectively. In this way the imputation models for each variable with missing data can be customized (remember that the default is to use all variables in the dataset to impute the variable(s) with missing data). In the hypothetical case that we only want to impute `x2` using `y`, and not both `y` and `x1` (note that this is only to exemplify how to customize the imputation model, I am not saying this is necessarily the way to go in this case).

```
pred <- imps0$predictorMatrix
pred[3,2] <- 0
imps_norm_pred <- mice(simdata, method = meth, predictorMatrix = pred, printFlag = FALSE,
                      seed = 1)
imps_norm_pred
```

```
## Class: mids
## Number of multiple imputations: 5
## Imputation methods:
##           y           x1           x2
##           ""           "" "norm.boot"
## PredictorMatrix:
##    y x1 x2
## y  0  1  1
## x1 1  0  1
## x2 1  0  0
```

We will now proceed to step 2, and we will use the function `with()`. Suppose that our substantive model, i.e., our model of interest, is the model we have used to generate the data, that is

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon, \quad \varepsilon \sim N(0, 1).$$

Just for the sake of illustration, I will be using the completed datasets stored in the object `imps`, using `mice`'s defaults.

```
fits <- with(imps, lm(y ~ x1 + x2))
class(fits)
```

```
## [1] "mira" "matrix"
```

The object `fits` contains the results of fitting  $M = 5$  complete data linear models based on the imputed datasets. The class of `fits` is `mira`, which stands for *multiply imputed repeated analysis*. We can extract the results and corresponding summary of the, say, first and second imputed datasets by doing

```
fits$analyses[[1]]
```

```
##
## Call:
## lm(formula = y ~ x1 + x2)
##
## Coefficients:
## (Intercept)          x1          x2
```

```
##          4.938          2.968          1.018
summary(fits$analyses[[1]])

##
## Call:
## lm(formula = y ~ x1 + x2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8393 -0.6519 -0.1273  0.6897  2.3144
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.93789    0.29489   16.75  <2e-16 ***
## x1          2.96771    0.07628   38.90  <2e-16 ***
## x2          1.01751    0.03838   26.51  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.015 on 97 degrees of freedom
## Multiple R-squared:  0.9591, Adjusted R-squared:  0.9582
## F-statistic: 1136 on 2 and 97 DF,  p-value: < 2.2e-16
```

```
fits$analyses[[2]]
```

```
##
## Call:
## lm(formula = y ~ x1 + x2)
##
## Coefficients:
## (Intercept)          x1          x2
##         4.888         2.925         1.041
summary(fits$analyses[[2]])

##
## Call:
## lm(formula = y ~ x1 + x2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8787 -0.7282 -0.1435  0.5983  2.8223
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.88841    0.31110   15.71  <2e-16 ***
## x1          2.92544    0.07997   36.58  <2e-16 ***
## x2          1.04074    0.04140   25.14  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.063 on 97 degrees of freedom
## Multiple R-squared:  0.9551, Adjusted R-squared:  0.9541
## F-statistic: 1031 on 2 and 97 DF,  p-value: < 2.2e-16
```

The final step is to combine (pool) the analyses to the final estimates using the pool function.



```
ests <- pool(fits)
#ests
summary(ests, conf.int = TRUE)
```

```
##          term estimate  std.error statistic    df p.value    2.5 %    97.5 %
## 1 (Intercept) 4.948381 0.30048570  16.46794 90.28135      0 4.3514390 5.545323
## 2          x1 2.947302 0.07947479   37.08474 79.49713      0 2.7891265 3.105477
## 3          x2 1.022477 0.04129880   24.75803 64.66538      0 0.9399896 1.104965
```

The object `ests` is of class `mipo`, meaning *multiply imputed pooled outcomes*. Its printed output resembles the output of an `lm` object, but note that its content is different: `pool` gathers the data in `mipo` in a *mira* way that makes summarising the statistics using `summary` easier. One cannot therefore use `residuals` or `predict` to obtain residuals or predictions from the final estimated model.

The column `estimate` correspond to the pooled regression coefficients and their corresponding standard error is available in `std.error`. By further inspecting the output we have columns corresponding to `ubar`, which is the within-imputation variance  $\bar{U}$ , `b` corresponds to the between-imputation variance, `rvi`, which stands for *relative increase in variance* due to the missing values and as we have learned last week, its expression is given by  $\frac{B + \frac{B}{M}}{\bar{U}}$ , the column corresponding to `lambda`, which is the proportion of variance in the parameter of interest due to the missing values and which is given by  $\frac{B + \frac{B}{M}}{V_T}$ . Finally, `fmi` contains the *fraction of missing information* as defined in Rubin (1987), and it depends on `rvi` but we will not study it further.

We can also only select the columns we are interested in from the summary, as illustrated below. Further note that the argument `conf.int = TRUE` computes a 95% confidence interval for the (pooled) coefficient estimates.

```
summary(ests, conf.int = TRUE)[, c(2, 3, 7, 8)]
```

```
##  estimate  std.error    2.5 %    97.5 %
## 1 4.948381 0.30048570 4.3514390 5.545323
## 2 2.947302 0.07947479 2.7891265 3.105477
## 3 1.022477 0.04129880 0.9399896 1.104965
```

For **linear** regression models, the pooled  $R^2$  can be calculated using the function `pool.r.squared()`.

```
pool.r.squared(fits, adjusted = TRUE)
```

```
##          est    lo 95    hi 95    fmi
## adj R^2 0.9579276 0.937119 0.9719522 0.07208382
```

The arguments `adjusted` specifies whether the adjusted  $R^2$  or the standard  $R^2$  is returned.

To conclude, let us check the effect of the choice of  $M$  on the results which, of course, in practice, depends on the particular analysis we are doing.

```
#using the default M=5 but changing the seed
ests_seed2 <- pool(with(mice(simdata, printFlag = FALSE, seed = 11), lm(y ~ x1 + x2)))
ests_seed3 <- pool(with(mice(simdata, printFlag = FALSE, seed = 111), lm(y ~ x1 + x2)))
```

```
summary(ests, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]
```

```
##  estimate  std.error p.value    2.5 %    97.5 %
## 1 4.948381 0.30048570      0 4.3514390 5.545323
## 2 2.947302 0.07947479      0 2.7891265 3.105477
## 3 1.022477 0.04129880      0 0.9399896 1.104965
```

```
summary(ests_seed2, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]
```

```
##  estimate  std.error p.value    2.5 %    97.5 %
```

```
## 1 4.975028 0.30434507      0 4.3696852 5.580371
## 2 2.926547 0.08158902      0 2.7635916 3.089502
## 3 1.023931 0.04061770      0 0.9430113 1.104850
```

```
summary(ests_seed3, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]
```

```
##      estimate  std.error    p.value      2.5 %   97.5 %
## 1 4.951749 0.32017344 0.000000e+00 4.3076009 5.595897
## 2 2.944345 0.09815355 1.110223e-15 2.7365704 3.152119
## 3 1.023076 0.03941279 0.000000e+00 0.9446045 1.101547
```

```
#using the M=20 and changing the seed
```

```
ests_seed1_20 <- pool(with(mice(simdata, printFlag = FALSE, seed = 1, m = 20), lm(y ~ x1 + x2)))
ests_seed2_20 <- pool(with(mice(simdata, printFlag = FALSE, seed = 11, m = 20), lm(y ~ x1 + x2)))
ests_seed3_20 <- pool(with(mice(simdata, printFlag = FALSE, seed = 111, m = 20), lm(y ~ x1 + x2)))
```

```
summary(ests_seed1_20, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]
```

```
##      estimate  std.error p.value      2.5 %   97.5 %
## 1 4.926221 0.30482933      0 4.3202325 5.532209
## 2 2.959692 0.08561980      0 2.7885989 3.130786
## 3 1.022750 0.04034451      0 0.9424956 1.103004
```

```
summary(ests_seed2_20, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]
```

```
##      estimate  std.error p.value      2.5 %   97.5 %
## 1 4.958667 0.30200150      0 4.3586337 5.558700
## 2 2.955109 0.08509408      0 2.7852253 3.124992
## 3 1.017642 0.03943637      0 0.9392827 1.096000
```

```
summary(ests_seed3_20, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]
```

```
##      estimate  std.error p.value      2.5 %   97.5 %
## 1 4.938306 0.30942374      0 4.3228618 5.553751
## 2 2.947024 0.08790346      0 2.7710699 3.122978
## 3 1.024396 0.04017278      0 0.9445281 1.104264
```

```
#using the M=50 and changing the seed
```

```
ests_seed1_50 <- pool(with(mice(simdata, printFlag = FALSE, seed = 1, m = 50), lm(y ~ x1 + x2)))
ests_seed2_50 <- pool(with(mice(simdata, printFlag = FALSE, seed = 11, m = 50), lm(y ~ x1 + x2)))
ests_seed3_50 <- pool(with(mice(simdata, printFlag = FALSE, seed = 111, m = 50), lm(y ~ x1 + x2)))
```

```
summary(ests_seed1_50, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]
```

```
##      estimate  std.error p.value      2.5 %   97.5 %
## 1 4.949031 0.30095407      0 4.3511043 5.546957
## 2 2.949120 0.08577777      0 2.7780191 3.120221
## 3 1.021698 0.04022434      0 0.9417208 1.101674
```

```
summary(ests_seed2_50, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]
```

```
##      estimate  std.error p.value      2.5 %   97.5 %
## 1 4.951577 0.30207745      0 4.3513394 5.551814
## 2 2.951857 0.08441553      0 2.7836038 3.120110
## 3 1.020227 0.03971142      0 0.9413022 1.099152
```

```
summary(ests_seed3_50, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]
```

```
##      estimate  std.error p.value      2.5 %   97.5 %
```

```
## 1 4.957735 0.30283071      0 4.3559624 5.559507
## 2 2.947614 0.08337180      0 2.7815445 3.113683
## 3 1.020432 0.03978913      0 0.9413519 1.099513

#using the M=100 and changing the seed
ests_seed1_100 <- pool(with(mice(simdata, printFlag = FALSE, seed = 1, m = 100), lm(y ~ x1 + x2)))
ests_seed2_100 <- pool(with(mice(simdata, printFlag = FALSE, seed = 11, m = 100), lm(y ~ x1 + x2)))
ests_seed3_100 <- pool(with(mice(simdata, printFlag = FALSE, seed = 111, m = 100), lm(y ~ x1 + x2)))

summary(ests_seed1_100, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]

##   estimate  std.error p.value    2.5 %   97.5 %
## 1 4.943259 0.30348831      0 4.3401883 5.546330
## 2 2.944795 0.08379677      0 2.7779059 3.111683
## 3 1.024108 0.03996190      0 0.9446866 1.103530

summary(ests_seed2_100, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]

##   estimate  std.error p.value    2.5 %   97.5 %
## 1 4.950286 0.30149773      0 4.3512125 5.549359
## 2 2.945878 0.08419970      0 2.7781334 3.113623
## 3 1.022392 0.03931929      0 0.9442725 1.100511

summary(ests_seed3_100, conf.int = TRUE)[, c(2, 3, 6, 7, 8)]

##   estimate  std.error p.value    2.5 %   97.5 %
## 1 4.957679 0.30407645      0 4.3533653 5.561993
## 2 2.941578 0.08564185      0 2.7708611 3.112294
## 3 1.022356 0.03970532      0 0.9434533 1.101258
```

The (pooled) estimates, standard errors, and the bounds of the intervals get more stable as  $M$  increases and we can be more confident in any one specific run. Note that whatever value of  $M$  we choose, there will always be some variation in results between repeat runs. The point is that with a sufficiently large  $M$ , the results will with high probability only differ by a small amount.