

# Missing The Point: Non-Convergence in Iterative Imputation Algorithms

Hanne Oberman

19 november 2020

## Contents

<b>1</b>	<b>Abstract</b>	<b>1</b>
<b>2</b>	<b>Intro</b>	<b>2</b>
2.1	One . . . . .	2
2.2	Two . . . . .	2
2.3	Three . . . . .	2
2.4	Original text [remove later] . . . . .	3
<b>3</b>	<b>Identifying non-convergence [remove later]</b>	<b>3</b>
<b>4</b>	<b>Simulation Set-Up</b>	<b>3</b>
4.1	Aims . . . . .	4
4.2	Data generating mechanism . . . . .	4
4.3	Estimands . . . . .	4
4.4	Methods . . . . .	5
4.5	Performance measures . . . . .	5
<b>5</b>	<b>Simulation Results</b>	<b>5</b>
5.1	Diagnostic Methods . . . . .	5
5.2	Performance Measures . . . . .	7
<b>6</b>	<b>Discussion</b>	<b>8</b>
<b>7</b>	<b>Case Study</b>	<b>8</b>
	<b>References</b>	<b>9</b>

*Notation in this manuscript:*

- square brackets are comments for myself/stuff I need to review
- bullet points are things I need to expand

## 1 Abstract

Iterative imputation has become the de facto standard approach to accommodate the ubiquitous problem of missing data. While it is widely accepted that this technique can yield valid inferences, these inferences all rely on algorithmic convergence. Our study provides insight into identifying non-convergence in iterative imputation algorithms. We show that these algorithms can yield correct outcomes even when a converged state has not yet formally been reached. In the cases considered, inferential validity is achieved after five to

ten iterations, much earlier than indicated by diagnostic methods. We conclude that it never hurts to iterate longer, but such calculations hardly bring added value.

## 2 Intro

### 2.1 One

- iterative imputation has become the de facto standard approach to accommodate missing data
- the aim is to draw valid inference: to get unbiased, confidence-valid estimates that incorporate the effect of the missingness
- the estimates are obtained by separating the missing data problem from the scientific problem
- the missing data problem is solved first: by imputing (i.e., filling in) the missing values in the incomplete dataset, using an algorithm such as MICE
- then the analysis of scientific interest can be performed on the completed data, which yields the scientific estimates
- to obtain valid scientific estimates, both the missing data problem and the scientific problem should be solved correctly
- one aspect that is often overlooked is the convergence of the algorithm that is used to generate imputations

### 2.2 Two

- there is no consensus on the convergence properties of iterative imputation algorithms, and there has not been a systematic study on how to evaluate the convergence of these algorithms
- yet, the inferences rely on algorithmic convergence
- the current practice is visual inspection, but this may be inadequate because of reasons 1, 2, and 3 [see original text]
- it is difficult that there is not at single point at which convergence is reached: the algorithm will produce some fluctuations even after convergence [it's convergence in distribution, not to a point]
- therefore, we should look at *non*-convergence
- there are non-convergence identifiers for other iterative algorithms (MCMC), but it is not known whether these are appropriate for imputation algorithms as well

### 2.3 Three

- in this paper, we explore different methods for identifying non-convergence in iterative imputation algorithms
- since the ultimate goal is to have inferential validity, we use this as performance measure in a simulation study: if the estimates are unbiased and confidence-valid, the non-convergence is, apparently, not influential
- we develop guidelines for practice to interpret the non-convergence identifiers, because we typically don't have the ground truth (to calculate bias and coverage)
- we show how all this works by means of a case study

## 2.4 Original text [remove later]

Iterative imputation has become the de facto standard approach to accommodate missing data. With iterative imputation, we ‘impute’ (i.e., fill in) plausible data values in the place of missing values in an incomplete dataset. We obtain these plausible values by sampling from a distribution that we model for each missing entry and update sequentially. Buuren et al. (2006) have shown that with an iterative imputation algorithm such as ‘MICE’, just a few iterations are enough to yield unbiased, confidence-valid statistical inferences.

Iterative imputation techniques have proven to be powerful tools to draw valid inference under many missing data circumstances (Rubin 1987; Van Buuren 2018).

With iterative imputation, the validity of the inference depends on the state-space of the algorithm at the final iteration. This introduces a potential threat to the validity of the imputations: What if the algorithm has not converged? Are the imputations then to be trusted? And can we rely on the inference obtained using the imputed data? These remain open questions since the convergence properties of iterative imputation algorithms have not been systematically studied (Van Buuren 2018). [Common convergence diagnostics may not work.] While there is no scientific consensus on how to evaluate the convergence of imputation algorithms (Liu et al. 2014; Zhu and Raghunathan 2015; Takahashi 2017), the current practice is to visually inspect imputations for signs of non-convergence.

Identifying non-convergence through visual inspection may be undesirable for several reasons: 1) it may be challenging to the untrained eye, 2) only severely pathological cases of non-convergence may be diagnosed, and 3) there is not an objective measure that quantifies convergence (Van Buuren, 2018, 6.5.2). Therefore, a quantitative diagnostic method to identify non-convergence would be preferred.

In this paper, we explore some existing and new means of diagnosing non-convergence in iterative algorithms. We evaluate the validity of these methods in the context of missing data imputation using model-based simulation in R (R Core Team 2020). [Brief overview of paper: the diagnostics, the simulation, the results, the discussion, a case study.]

## 3 Identifying non-convergence [remove later]

- Iterative imputation is a sort of MCMC algorithm, so it makes sense to investigate non-convergence in a similar fashion to typical MCMC algorithms.
- With MCMC, the generated values will vary even after convergence, so there is not a unique point at which convergence is established (Gelman et al. 2013). Therefore, diagnostic methods may only identify signs of *non*-convergence (Hoff 2009). Non-convergence occurs when one of two requirements for convergence is not met.
- There are two requirements for convergence of iterative algorithms: mixing and stationarity (Gelman et al. 2013). Mixing implies that generated values intermingle nicely, and stationarity is characterized by the absence of trending between successive draws.
- What are the consequences? → bias, under-coverage [refer to van Buuren (2018) instead of reproducing this]
- We consider autocorrelation and rhat, and monitor four parameters: chain means, chain variances, a scientific estimate, and lambda.
- Explain lambda

## 4 Simulation Set-Up

We investigate non-convergence in iterative imputation through model-based simulation in R (version 4.0.3; R Core Team 2020). We provide a summary of the simulation set-up in Algorithm 1; the complete script and technical details are available from [github.com/hanneoberman/MissingThePoint](https://github.com/hanneoberman/MissingThePoint).

### Algorithm 1: simulation set-up (pseudo-code)

```
for each simulation repetition
  1. simulate complete data
  for each missingness condition
    2. create missingness
    for each iteration
      3. impute missingness
      4. perform analysis of scientific interest
      5. apply non-convergence identifiers
      6. pool results across imputations
      7. compute performance measures
    8. combine outcomes from all iterations
  9. combine outcomes from all missingness conditions
10. aggregate outcomes from all simulation repetitions
```

## 4.1 Aims

With this simulation, we assess the impact of non-convergence on the validity of scientific estimates obtained using the imputation package `{mice}` (Van Buuren and Groothuis-Oudshoorn 2011). Inferential validity is reached when estimates are both unbiased and have nominal coverage across simulation repetitions ( $n_{sim} = 1000$ ). To induce non-convergence, we terminate the algorithm after a varying number of iterations ( $n_{it} = 1, 2, \dots, 50$ ). We differentiate between nine different missingness scenarios that are defined by the data generating mechanism.

## 4.2 Data generating mechanism

Data are generated in each simulation repetition for a complete set of  $n_{obs} = 1000$  cases (i.e., before inducing missingness). We define three multivariately normal random variables, let

$$\begin{pmatrix} Y \\ X_1 \\ X_2 \end{pmatrix} \sim \mathcal{N} \left[ \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 & & \\ 0.5 & 1 & \\ -0.5 & 0.5 & 1 \end{pmatrix} \right].$$

The complete set is amputed according to nine missingness conditions. We use a  $3 \times 3$  factorial design consisting of three missingness mechanisms and three proportions of incomplete cases.

- we use the three missingness mechanisms MCAR, MAR, and MNAR with defaults settings from `mice::ampute()` (e.g., right-tailed MAR)
- we use a multivariate missing data pattern and make 25%, 50%, and 75% of cases incomplete

## 4.3 Estimands

We impute the missing data five times ( $m = 5$ ) using Bayesian linear regression imputation with `mice` (Van Buuren and Groothuis-Oudshoorn 2011). On each imputed dataset, we perform multiple linear regression to predict outcome variable  $Y$  from the other two variables

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2,$$

where  $\hat{Y}$  is the predicted value of the outcome. Our estimands are the regression coefficient  $\beta_1$  and coefficient of determination  $\rho^2$  that we obtain after pooling the regression results across the imputations.

## 4.4 Methods

- non-convergence in iterative algorithms is diagnosed using an identifier and a parameter
- the parameter can be any statistic that we track across iterations, for example chain means (i.e, the average imputed value per imputation)
- the identifier is a calculation of some sort that quantifies non-convergence
- identifiers are historically focused on either non-mixing between chains or non-stationarity within chains
- the popular non-mixing identifier that has recently been updated by Vehtari et al, and should now work for non-stationarity as well
- just to be sure, we also use autocorrelation to quantify trending within chains
- we apply these identifiers to four different univariate and multivariate parameters
- the univariate parameters are those commonly used for visual inspection: chain means and chain variances
- one multivariate parameter that is of immediate interest for our estimand is the estimated value itself: regression coefficient  $\beta_1$
- we also propose a new multivariate parameter that is not dependent on the model of scientific interest: the first eigenvalue of the variance-covariance matrix of the completed data [explain!]

We use eight different methods to diagnose non-convergence: a combination of two non-convergence identifiers—autocorrelation and  $\hat{R}$ —and four parameters—chain means, chain variances,  $\beta_1$ , and  $\lambda$ .

## 4.5 Performance measures

As recommended by Van Buuren (2018), our performance measures are bias, confidence interval width, and coverage rate of the estimands (§ 2.5.2)[or just of the regression estimate?? otherwise add the ciw and cov of  $r^2$  too!!].

# 5 Simulation Results

The following figures display the simulation results for the eight diagnostic methods and four performance measures, contrasted to the number of iterations in the imputation algorithm. Within the figures, we split the results according to the missingness conditions [missingness mechanisms as line types, and proportion of incomplete cases as colors]. Note that these results are averages of the  $n_{sim} = 1000$  simulation repetitions.

## 5.1 Diagnostic Methods

Figure 1 shows the autocorrelations in the chain means (panel A), chain variances (panel B), regression estimates (panel C), and eigenvalues (panel D).

Autocorrelation in the chain means decreases rapidly in the first few iterations (see 1A). The decrease is substantive until  $n_{it} \geq 6$ . This means that there is some initial trending within chains, but the average imputed value quickly reaches stationarity. These results hold irrespective of the missingness condition.

Autocorrelation in the chain variances show us something similar (see 1B). The number of iterations that is required to reach non-improving autocorrelations is somewhat more ambiguous than for chain means, but generally around  $n_{it} \geq 10$ . We do not observe a systematic difference between missingness conditions here either.

There is more autocorrelation in the scientific estimates than in the chain means and chain variances (see 1C). We observe the highest autocorrelations in conditions where 75% of cases are incomplete. Overall, the autocorrelations reach a plateau when  $n_{it} \geq 20$  to 30. There is no clear effect of the missingness mechanisms.

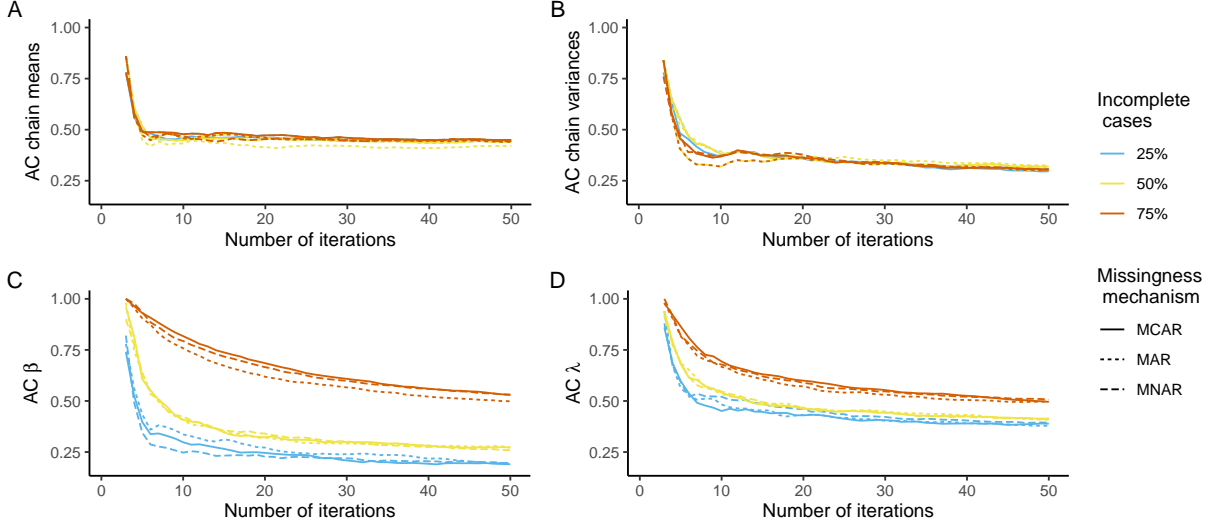


Figure 1: Autocorrelation (AC) with different parameters.

The autocorrelation in the eigenvalues exhibits a similar trend to the autocorrelation in the scientific estimates (see 1D). Trending in this parameter diminishes when  $n_{it} \geq 20$ .

Figure 2 shows the potential scale reduction factor in the chain means (panel A), chain variances (panel B), regression estimates (panel C), and eigenvalues (panel D).

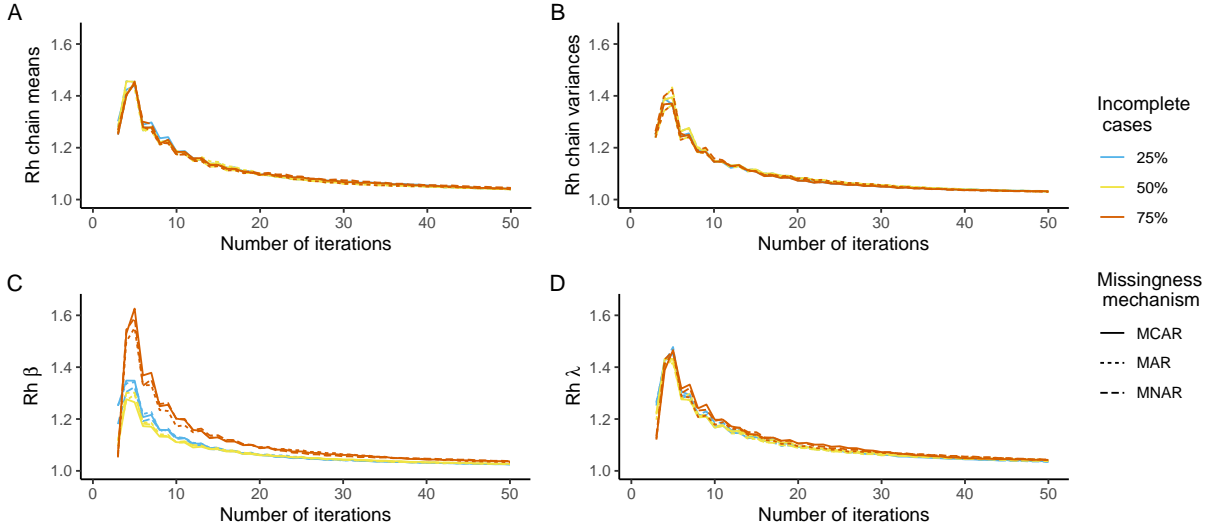


Figure 2: Potential scale reduction factor (Rh) with different parameters.

We observe that  $\hat{R}$ -values of the chain means generally decreases as a function of the number of iterations (see 2A). An exception to this observation is the initial increase when  $3 \leq n_{it} \leq 5$  [interpret?? due to initialization or is there really more mixing initially??]. After the first couple of iterations, the mixing between chain means generally improves until  $n_{it} \geq 30$  to 40. There is no apparent differentiation between the missingness conditions.

The mixing between chain variances mimics the mixing between chain means almost perfectly (see 2B).

Irrespective of the missingness condition, the  $\hat{R}$ -values taper off around  $n_{it} \geq 30$ .

With the regression estimate as parameter we observe very similar  $\hat{R}$ -values than with chain means and chain variances (see 2C). We do, however, see some differences between missingness conditions. Conditions where 75% of the cases are incomplete show more extreme non-mixing. The overall trend remains the same: about 30 iterations are required before mixing stops improving substantially.

$\hat{R}$ -values of the eigenvalues show a trend similar to the chain means and chain variances (see 2D). [add something about conditions??]

[These rhat plots all show some initialization before the fifth iteration: is rhat useful before that??]

## 5.2 Performance Measures

In Figure 3 we show the performance measures: bias in the regression estimate (panel A), bias in the coefficient of determination (panel B), the empirical coverage rate of the regression estimate (panel C) and the average confidence interval width of this estimate (panel D).

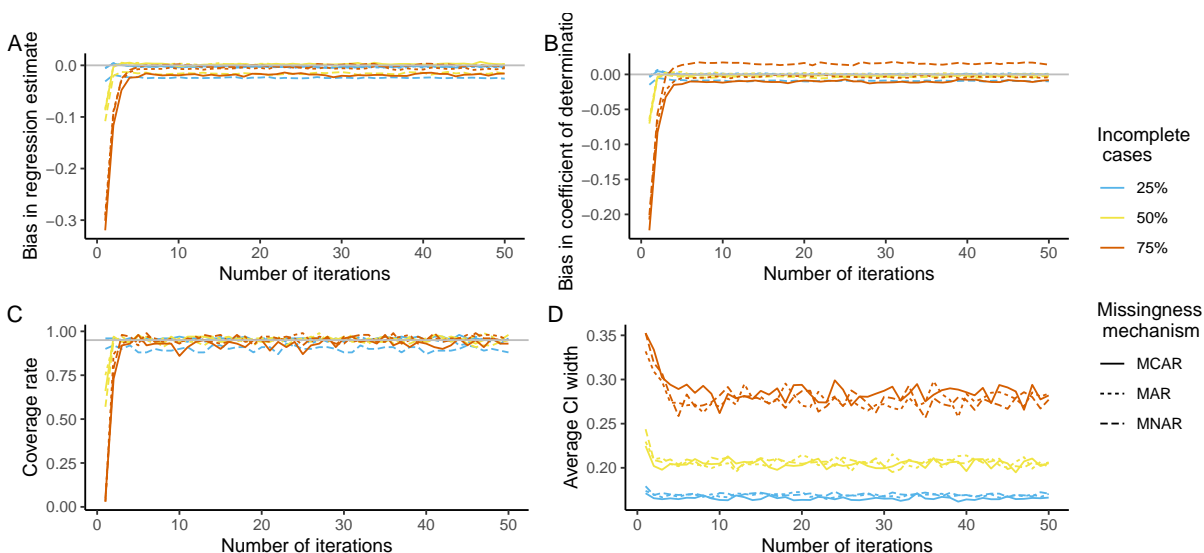


Figure 3: Performance measures.

We see that within a few iterations the bias in the regression estimate approaches zero (see 3A). When  $n_{it} \geq 6$ , even the worst-performing conditions (e.g., with a proportion of incomplete cases of 75%) produce stable, non-improving estimates [regression coefficient is underestimated because there is less info to estimate the relation??].

Equivalent to the bias in the regression estimate, the bias in the coefficient of determination tapers off within a couple of iterations (see 3B). We observe stable estimates in all conditions when  $n_{it} \geq 6$  [interpret the over-estimation in MNAR+75% condition?].

Nominal coverage is quickly reached (see 3C). After just three iterations, the coverage rates are non-improving in every missingness condition [but MNAR with 5% incomplete cases does not reach nominal coverage  $\rightarrow$  due to bias in the estimate in combination with very narrow CI (see CIW!)].

The average confidence interval width decreases quickly with every added iteration until a stable plateau is reached (see 3D). Depending on the proportion of incomplete cases this takes up-to  $n_{it} \geq 9$ .

## 6 Discussion

Our study found that—in the cases considered—inferential validity was achieved after five to ten iterations, much earlier than indicated by the non-convergence identifiers. Of course, it never hurts to iterate longer, but such calculations hardly bring added value.

- Convergence diagnostics keep improving substantially until  $n\_it = 20-30$
- Performance measures do not improve after  $n\_it = 9$
- [methodological explanation is that  $rhat$  and  $ac$  have a lag (few it to inform your statistic)  $\rightarrow$  will always indicate convergence slower than inferential validity is reached]
- Univariate thetas may under-estimate non-convergence.
- Determining non-stationarity with  $\lambda$  is more difficult than with  $qhat$  :(

In short, the validity of iterative imputation stands or falls with algorithmic convergence—or so it's thought. We have shown that iterative imputation algorithms can yield correct outcomes even when a converged state has not yet formally been reached. Any further iterations just burn computational resources without improving the statistical inferences.

## 7 Case Study

- We use real data: the `boys` dataset from the `mice` package
- We are interested in predicting age from the other variables, in particular in the regression coefficient of `hgt`
- We compare non-convergence identified using visual inspection versus  $rhat$  in the chain variances, scientific estimate and  $\lambda$ .
- The figures show results of a `mice` run with 20 iterations but otherwise default settings.

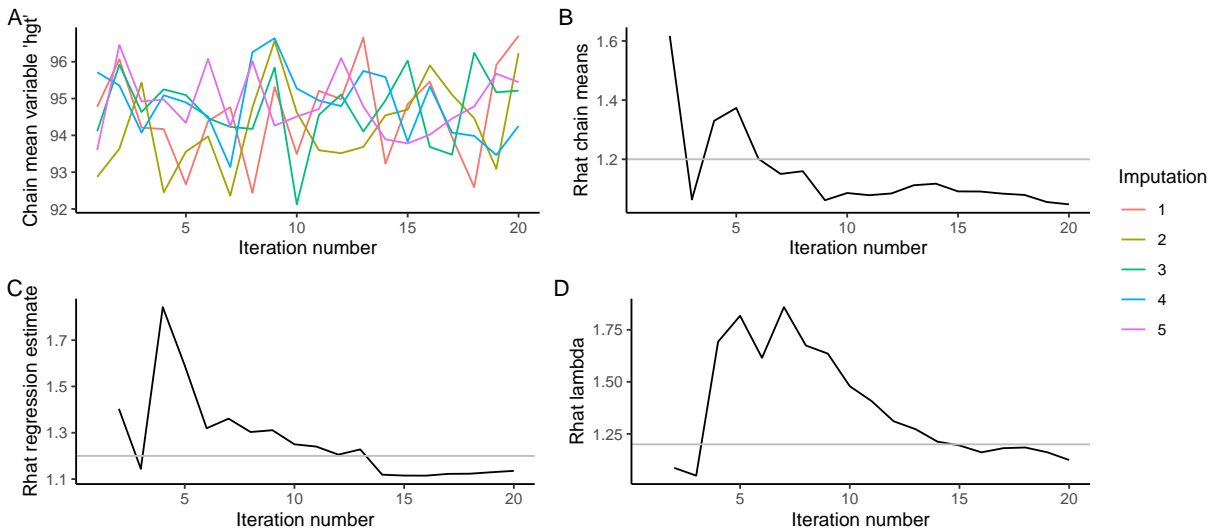


Figure 4: Case study.

From the traceplot of the chain means (see 4A) it seems that mixing improves up-to 10 iterations, while trending is only apparent in the first three iterations.



This figure (4B) shows that 7 iterations are required before the  $\hat{R}$ -values of the chain means drop below the threshold for non-convergence.

The  $\hat{R}$ -values for the scientific estimate reaches the threshold much sooner, when  $n_{it} = 14$  (see 4C).

According to the  $\hat{R}$ -values with  $\lambda$  as parameter, at least 15 iterations are required (see 4D).

## References

- Buuren, S. Van, J. P. L. Brand, C. G. M. Groothuis-Oudshoorn, and D. B. Rubin. 2006. “Fully Conditional Specification in Multivariate Imputation.” *Journal of Statistical Computation and Simulation* 76 (12): 1049–64. <https://doi.org/10.1080/10629360600810434>.
- Gelman, Andrew, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin. 2013. *Bayesian Data Analysis*. Philadelphia, PA, United States: CRC Press LLC.
- Hoff, Peter D. 2009. *A First Course in Bayesian Statistical Methods*. Springer Texts in Statistics. New York, NY: Springer New York. <https://doi.org/10.1007/978-0-387-92407-6>.
- Liu, J., A. Gelman, J. Hill, Y.-S. Su, and J. Kropko. 2014. “On the Stationary Distribution of Iterative Imputations.” *Biometrika* 101 (1): 155–73. <https://doi.org/10.1093/biomet/ast044>.
- R Core Team. 2020. *R: A Language and Environment for Statistical Computing*. Vienna, Austria.
- Rubin, Donald B. 1987. *Multiple Imputation for Nonresponse in Surveys*. Wiley Series in Probability and Mathematical Statistics Applied Probability and Statistics. New York, NY: Wiley.
- Takahashi, Masayoshi. 2017. “Statistical Inference in Missing Data by MCMC and Non-MCMC Multiple Imputation Algorithms: Assessing the Effects of Between-Imputation Iterations.” *Data Science Journal* 16 (July): 37. <https://doi.org/10.5334/dsj-2017-037>.
- Van Buuren, Stef. 2018. *Flexible Imputation of Missing Data*. Chapman and Hall/CRC.
- Van Buuren, Stef, and Karin Groothuis-Oudshoorn. 2011. “Mice: Multivariate Imputation by Chained Equations in R.” *Journal of Statistical Software* 45 (1): 1–67. <https://doi.org/10.18637/jss.v045.i03>.
- Zhu, Jian, and Trivellore E. Raghunathan. 2015. “Convergence Properties of a Sequential Regression Multiple Imputation Algorithm.” *Journal of the American Statistical Association* 110 (511): 1112–24. <https://doi.org/10.1080/01621459.2014.948117>.