

---

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

Sociological Methods & Research  
2020, Vol. 49(?):1–23  
© The Author(s) 2020  
Article reuse guidelines:  
sagepub.com/journals-permissions  
DOI: 10.1177/ToBeAssigned  
journals.sagepub.com/home/smr



H. I. Oberman<sup>1</sup>

## Abstract

Iterative imputation is a popular tool to accommodate missing data. While it is widely accepted that valid inferences can be obtained with this technique, these inferences all rely on algorithmic convergence. There is no consensus on how to evaluate the convergence properties of the method. This paper provides insight into identifying non-convergence of iterative imputation algorithms.

## Keywords

missing data, iterative imputation, non-convergence, mice

## Introduction

Anyone who analyzes person-data may run into a missing data problem. Missing data is not only ubiquitous, but treating it can also be tedious. If a dataset contains just one incomplete observation, statistical inferences are undefined and will not produce any results. To circumvent this, many statistical packages employ list-wise deletion by default (i.e., ignoring incomplete observations). Unfortunately, this *ad hoc* solution may yield wildly invalid results (Van Buuren 2018). An alternative is to *impute* (i.e., fill in) the missing values in the incomplete observations. Subsequently, statistical inferences can be performed on the completed dataset. By repeating this process several times, a distribution of plausible results may be obtained, which reflects the uncertainty in the data due to missingness. This technique is known as ‘multiple imputation’ (MI; Rubin 1976).

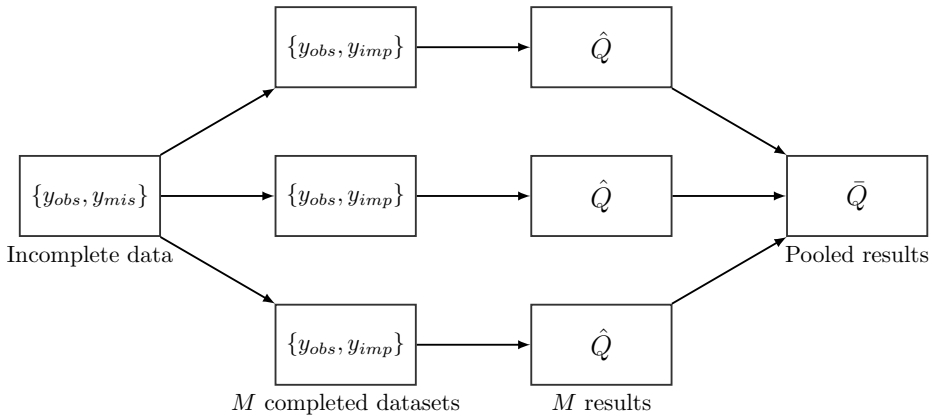
---

<sup>1</sup>Department of Methodology and Statistics, Utrecht University, Utrecht, The Netherlands

### Corresponding author:

Hanne Oberman, Sjoerd Groenman building, Utrecht Science Park, Utrecht, The Netherlands.

Email: [h.i.oberman@uu.nl](mailto:h.i.oberman@uu.nl)



**Figure 1.** Scheme of the main steps in multiple imputation (where  $M = 3$ )—from an incomplete dataset, to  $M$  multiply imputed datasets, to  $M$  estimated quantities of scientific interest  $\hat{Q}$ s, to a single pooled estimate  $\bar{Q}$ .

MI has proven to be a powerful technique to draw valid inferences from incomplete data under many circumstances (Van Buuren 2018).

Figure 1 provides an overview of the steps involved with MI. Missing data  $y_{mis}$  in an incomplete dataset is imputed  $M$  times. The imputed data  $y_{imp}$  is combined with the observed data  $y_{obs}$  to create  $M$  completed datasets. On each completed dataset, an analysis of scientific interest is performed. The quantity of scientific interest (e.g., a regression coefficient) is denoted with  $Q$ . Since  $Q$  is estimated on each completed dataset,  $M$  separate  $\hat{Q}$ -values are obtained. These  $M$  values are combined into a single pooled estimate  $\bar{Q}$ . The premise of multiple imputation is that  $\bar{Q}$  is an unbiased and confidence valid estimate of the true—but missing—data inference.

A popular method to obtain imputations is to use the ‘Multiple Imputation by Chained Equations’ algorithm, shorthand ‘MICE’ (Van Buuren and Groothuis-Oudshoorn 2011). With MICE, imputed values  $y_{imp}$  are drawn from the posterior predictive distribution of the missing values  $y_{mis}$ . The algorithm is named after the iterative algorithmic procedure by which imputed values are generated: a multivariate distribution is obtained by iterating over a sequence of univariate imputations. The iterative nature of algorithms like MICE introduces a potential threat to the validity of the imputations: What if the algorithm has not converged? Are the implications then to be trusted? And can we rely on the inference obtained on the completed data? These are all open questions, because the convergence properties of iterative imputation algorithms have not been systematically studied (Van Buuren 2018). Moreover, there is no scientific consensus on how to evaluate convergence of imputation algorithms (Takahashi 2017). Some default MICE techniques (e.g., ‘predictive mean modeling’) might not yield converged states at all (Murray 2018). Therefore, algorithmic convergence should be monitored carefully.

Currently, the recommended practice for evaluating the convergence of iterative imputation algorithms is through visual inspection. That is, the state space of the algorithm is monitored by plotting a scalar summary  $\theta$  against the iteration number. **[Segway to WHAT TO MONITOR? Also, this paragraph is more or less the same as the third paragraph in the second section. Merge or remove one of the two?]** Typically, the  $\theta$ s under evaluation are chain means and variances (**define!**). These  $\theta$ s may be insufficient because they are summaries of the univariate state space of the algorithm, while MICE is not only concerned with univariate estimates, but the entire multivariate distribution of  $y_{imp}$ . A suggestion by [Van Buuren \(2018\)](#) for a multivariate  $\theta$  to monitor is the quantity of scientific interest  $Q$  (e.g., a regression coefficient). Implementing this, however, might be somewhat advanced for empirical researchers. Moreover, this scalar summary is not model-independent, i.e., it only applies to one model of scientific interest. Yet, one of the advantages of MI is that the missing data problem and scientific problem are solved independently. Therefore, this  $\theta$  is not sufficient either. [Van Buuren \(2018, § 4.5.2\)](#) also proposed multivariable evaluation of the MICE algorithm through eigenvalue decomposition, building on the work of [MacKay and Mac Kay \(2003\)](#). Eigenvalues of a variance-covariance matrix are a measure of the data's total covariance. This would be a model-independent, multivariate scalar summary to monitor, but it has not yet been implemented. **[Segway to HOW TO MONITOR?]** Monitoring convergence visually is insufficient on two counts: 1) it may be challenging to the untrained eye, and 2) only severely pathological cases of non-convergence may be diagnosed ([Van Buuren 2018, § 6.5.2](#)). Therefore, a quantitative, diagnostic evaluation of convergence would be preferred—although not straightforward. Iterative imputation algorithms such as MICE are Markov chain Monte Carlo (MCMC) methods. In MCMC methods, convergence is not from a scalar to a point, but from one distribution to another. The values generated by the algorithm (e.g., imputed values) will vary even after convergence ([Gelman et al. 2013](#)). Since MCMC algorithms do not reach a unique point at which convergence is established, diagnostics may only identify signs of *non-convergence* ([Hoff 2009](#)). Several of such non-convergence diagnostics exist, but it is not known whether these are appropriate for iterative imputation algorithms.

In this paper, we investigate how non-convergence in iterative imputation algorithms may be diagnosed. First, we define which non-convergence diagnostics will be considered, and evaluate how these diagnostics could be appropriate for iterative imputation applications. In the second part of this paper, we will investigate the impact of inducing non-convergence in iterative imputation algorithms by means of model-based simulation in R ([R Core Team 2020](#)). The aim of the simulation is to determine whether unbiased, confidence valid inferences may be obtained under different levels of non-convergence (**and at which point convergence may safely be assumed**). Additionally, we will assess how well different non-convergence diagnostics perform in identifying signs of non-convergence (**and on which scalar summary of the state space of the algorithm these diagnostics should be applied**). The results of this study are guidelines for identifying non-convergence in iterative imputation algorithms, which will aid applied researchers in drawing valid inference from incomplete datasets. **[add more**

**societal relevance: what is the consequence of non-convergence? see example of non-convergence in Figure 2 (e.g., biased estimates, invalid inference, etc.)**] For reasons of brevity, we only focus on the iterative imputation algorithm implemented in the popular *mice* package (Van Buuren and Groothuis-Oudshoorn 2011) in R (R Core Team 2020).

**[Include sub-questions and make structure of the paper clear!]** How can non-convergence be identified diagnostically? Are common MCMC non-convergence diagnostics appropriate for MICE? And if so, which threshold should be used to diagnose non-convergence? How many iterations are sufficient/needed to be able to diagnose non-convergence? Are the default number of iterations sufficient (i.e., 5 in *mice*, 10 in SPSS and Stata, 30 in *mi*)? How severe is it when the algorithm has not converged? And what are guidelines for practice? Can the parameter of interest  $Q$  be correct when the algorithm is not (yet) converged, and vice versa?

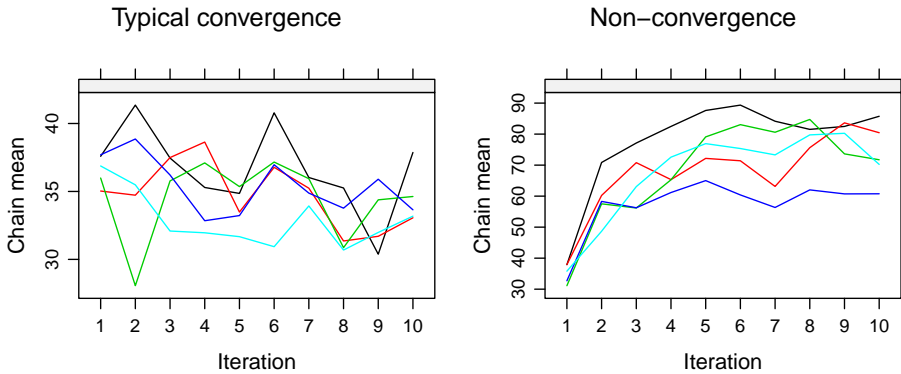
### Some notation

Let  $y$  denote an  $n \times k$  matrix containing the data values on  $k$  variables for all  $n$  units in a sample. The data value of unit  $i$  ( $i = 1, 2, \dots, n$ ) on variable  $j$  ( $j = 1, 2, \dots, k$ ) may be either observed or missing. The number of units  $i$  in dataset  $y$  with at least one missing value, divided by the total number of units  $n$ , is called the missingness proportion  $p_{mis}$ . The collection of observed data values in  $y$  is denoted by  $y_{obs}$ ; the missing part of  $y$  is referred to as  $y_{mis}$ . For each datapoint in  $y_{mis}$ , we sample  $M \times T$  plausible values, where  $M$  is the number of imputations ( $m = 1, 2, \dots, M$ ) and  $T$  is the number of iterations ( $t = 1, 2, \dots, T$ ) in the imputation algorithm. The collection of samples between the initial value (at  $t = 1$ ) and the final imputed value (at  $t = T$ ) will be referred to as an ‘imputation chain’. **[If not defined elsewhere, add:  $\theta$ s are scalar summaries of the algorithm at a certain iteration (e.g., chain means; the average of the imputed values in each imputation chain).]**

### Identifying non-convergence

There are two requirements for convergence of iterative algorithms: mixing and stationarity (Gelman et al. 2013). Without mixing, imputation chains do not intermingle nicely, indicating that the distribution of imputed values differs across imputations. Chains may be ‘stuck’ at a local optimum, instead of sampling imputed values from the entire predictive posterior distribution of the missing values. This may cause underestimation of the variance within chains, which results in spurious, invalid inferences. Without stationarity, there is trending within imputation chains. Trending implies that further iterations would yield a different set of imputations. Iterative imputation algorithms that have not yet reached stationarity, may thus yield biased estimates.

**[Move to INTRO???**] To illustrate what non-convergence looks like in iterative imputation algorithms, we reproduce an example from van Buuren (2018, § 6.5.2). Figure 2 displays a subset of the example: a variable that we will call  $j$ . We see two traceplots of the chain means for  $j$ —i.e., the average of the imputed values for  $j$  in each imputation  $y_{imp,m}$ , plotted against the iteration number. Each line portrays one imputation. The plot on the left-hand side of the figure shows typical convergence



**Figure 2.** Typical convergence versus pathological non-convergence. Please note that this is the same data with a different imputation model, leading to different imputations (e.g., see the units on the y-axis).

of an iterative imputation algorithm. The right-hand side displays pathological non-convergence, induced by purposefully mis-specifying the imputation model. **[Leave this in if it's after mixing and stationarity are defined, otherwise remove?]** In the typical convergence situation, the imputation chains intermingle nicely and there is little to no trending. In the non-convergence plot, there is a lot of trending and some chains do not intermingle. Importantly, the chain means at the last iteration (the imputed value per  $m$ ) are very different between the two plots. The algorithm with the mis-specified model yields imputed values that are on average a factor two larger than those of the typically converged algorithm. That means that non-convergence has an impact on the mean of  $j$  in  $y_{imp,m}$ . This difference (presumably) translates to bias in the  $M$  completed data estimates  $\hat{\mu}_{j,m}$ , and consequently also to the pooled estimate  $\bar{\mu}_j$ . Non-convergence thus leads to biased estimates. This shows the importance of reaching converged states in iterative imputation algorithms.

**[Move to INTRO??]** Convergence has historically been inspected visually, by monitoring the imputation algorithm over iterations. This is typically done through traceplots. In a traceplot, the iteration number is plotted against a certain  $\theta$ .  $\theta$ s are scalar summaries of the state space of the algorithm at a specific iteration. The recommended scalar summaries to evaluate are chain means and chain variances. Additionally, researchers may “monitor some statistic of scientific interest” (Van Buuren 2018, § 6.5.2). As Van Buuren (2018) describes, researchers can specify their quantity of scientific interest  $Q$  (e.g., a regression coefficient) as scalar summary  $\theta$ . Such a user-defined scalar summary, however, may be somewhat advanced for empirical researchers. And moreover, it is not universal to all complete data problems. Focusing on the convergence of outcome parameters may influence the iterative imputation procedure in the sense that the model of evaluation favors the model of interest. The downside to these two approaches is that they either focus on the univariate state space, or primarily track the change over the

iterations of a multivariate outcome conform the scientific model of interest. Ideally, one would like to evaluate a model-independent parameter that summarizes the multivariate nature of the data. Therefore, we propose a  $\theta$  that summarizes the multivariate state space of the algorithm, but is independent from the model of scientific interest. We propose  $\lambda_1$  as such a scalar summary. We define  $\lambda_1$  as the first eigenvalue of the variance-covariance matrix of the completed data. Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_j$  be the eigenvalues of  $\Sigma$  in each of the  $M$  completed datasets  $y_{obs}, y_{imp}$ .  $\lambda_1$  is measure that summarizes the covariances in the completed datasets. **The first eigenvalue has the appealing property that is not dependent on the substantive model of interest. Eigenvalue decomposition is inspired by MacKay and Mac Kay (2003).**

### *Diagnostics under evaluation*

Non-convergence diagnostics for MCMC algorithms typically identify problems in either one of the two requirements for convergence: mixing or stationarity. Non-stationarity within chains may be diagnosed with e.g., autocorrelation ( $AC$ ; Schafer 1997; Gelman et al. 2013), numeric standard error ('MC error'; Geweke 1992), or Raftery and Lewis's (1991) procedure to determine the effect of trending on the precision of estimates. A widely used diagnostic to monitor mixing between chains is the potential scale reduction factor  $\hat{R}$  ('Gelman-Rubin statistic'; Gelman and Rubin 1992). With a recently proposed adaptation,  $\hat{R}$  might also serve to diagnose non-stationarity, but this has not yet been thoroughly investigated (Vehtari et al. 2019). Therefore, we use  $\hat{R}$  and  $AC$  to evaluate mixing and stationarity separately, as recommended by e.g., Cowles and Carlin (1996).

**[Add: effective sample size is used if we want independent samples *within* chains. In MI, we only use the final estimate. Also, computing n.eff assumes infinitely many iterations, which is fine in Bayesian analyses where T is often at least one thousand. But the default in iterative imp is often substantially lower.]**

*Potential scale reduction factor* The potential scale reduction factor  $\hat{R}$  was coined by Gelman and Rubin (1992). An updated version of  $\hat{R}$  has been proposed by Vehtari et al. (2019, p. 5). **This version is better suited to detect non-convergence in the tails of distributions, by using three transformations on the scalar summary that it is applied to. Namely, rank-normalization, folding, and localization. The Vehtari et al. version of  $\hat{R}$  may be suitable for iterative imputation. We therefore use their definition of the diagnostic.** Let  $M$  be the total number of chains,  $T$  the number of iterations per chain, and  $\theta$  the scalar summary of interest (e.g., chain means). For each chain ( $m = 1, 2, \dots, M$ ), we estimate the variance of  $\theta$ , and average these to obtain within-chain variance  $W$ .

$$W = \frac{1}{M} \sum_{m=1}^M s_m^2, \text{ where } s_m^2 = \frac{1}{T-1} \sum_{t=1}^T \left( \theta^{(tm)} - \bar{\theta}^{(\cdot m)} \right)^2.$$

We then estimate between-chain variance  $B$  as the variance of the collection of average  $\theta$  per chain.

$$B = \frac{T}{M-1} \sum_{m=1}^M \left( \bar{\theta}^{(\cdot m)} - \bar{\theta}^{(\cdot \cdot)} \right)^2, \text{ where } \bar{\theta}^{(\cdot m)} = \frac{1}{T} \sum_{t=1}^T \theta^{(tm)}, \bar{\theta}^{(\cdot \cdot)} = \frac{1}{M} \sum_{m=1}^M \bar{\theta}^{(\cdot m)}.$$

From the between- and within-chain variances we compute a weighted average,  $\widehat{\text{var}}^+$ , which approximates the total variance of  $\theta$ .  $\widehat{R}$  is then obtained as a ratio between the total variance and the within-chain variance:

$$\widehat{R} = \sqrt{\frac{\widehat{\text{var}}^+(\theta|y)}{W}}, \text{ where } \widehat{\text{var}}^+(\theta|y) = \frac{N-1}{N}W + \frac{1}{N}B.$$

We can interpret  $\widehat{R}$  as potential scale reduction factor since it indicates by how much the variance of  $\theta$  could be shrunk down if an infinite number of iterations per chain would be run (Gelman and Rubin 1992). This interpretation assumes that chains are ‘over-dispersed’ at  $t = 1$ , and reach convergence as  $T \rightarrow \infty$ . Over-dispersion implies that the initial values of the chains are ‘far away’ from the target distribution and each other. When the sampled values in each chain are independent of the chain’s initial value, the mixing component of convergence is satisfied. The variance between chains is then equivalent to the variance within chains, and  $\widehat{R}$ -values will be close to one. High  $\widehat{R}$ -values thus indicate non-convergence. The conventionally acceptable threshold for convergence was  $\widehat{R} < 1.2$  (Gelman and Rubin 1992). More recently, Vehtari et al. (2019) proposed a more stringent threshold of  $\widehat{R} < 1.01$ .

*Autocorrelation* Following the same notation, we define autocorrelation as the correlation between two subsequent  $\theta$ -values within the same chain (Lynch 2007, p. 147). In this study, we only consider *AC* at lag 1, i.e., the correlation between the  $t^{\text{th}}$  and  $t + 1^{\text{th}}$  iteration of the same chain.

$$AC = \left( \frac{T}{T-1} \right) \frac{\sum_{t=1}^{T-1} (\theta_t - \bar{\theta}^{(\cdot m)})(\theta_{t+1} - \bar{\theta}^{(\cdot m)})}{\sum_{t=1}^T (\theta_t - \bar{\theta}^{(\cdot m)})^2}.$$

We can interpret *AC*-values as a measure of stationarity. If *AC*-values are close to zero, there is no dependence between subsequent samples within imputation chains. Negative *AC*-values indicate divergence within imputation chains. **Subsequent sampled values within each imputation chain are less alike.** Positive *AC*-values indicate recurrence. If  $\theta$ -values of subsequent iterations are similar, trending may occur. Negative *AC*-values show no threat to the stationarity component of convergence. On the contrary even—negative *AC*-values indicate that  $\theta$ -values of subsequent iterations diverge from one another, which may increase the variance of  $\theta$  and speed up convergence. As convergence diagnostic, the interest is therefore in positive *AC*-values. The magnitude of *AC*-values may be evaluated statistically, using a Wald test with  $AC = 0$  as null hypothesis (Box et al. 2015).

*Thresholds* Upon convergence,  $\hat{R} = 1$  and  $AC = 0$ , which are unlikely thresholds for MCMC algorithms, because of its convergence to a distribution. Even in the most converged state, the algorithm will show some signs of non-mixing and non-stationarity. Upon *sufficient* convergence, the imputation chains will intermingle such that the only difference between the chains is caused by the randomness induced by the algorithm ( $\hat{R} > 1$ ), and there may be some dependency between subsequent iterations of imputation chains ( $AC > 0$ ). In practice, non-convergence is usually diagnosed when  $\hat{R} > 1.2$  or 1.1 or even 1.01. And a t-test is performed to assess whether  $AC$  is significantly different from zero.

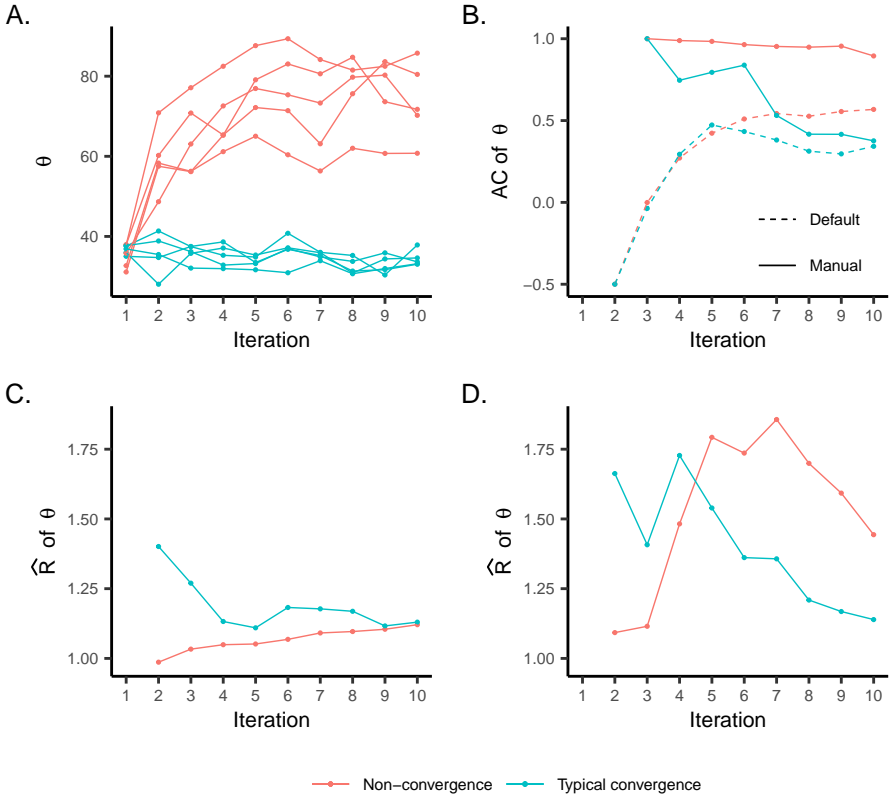
### *In practice*

Before we evaluate the performance of the non-convergence diagnostics  $\hat{R}$  and  $AC$  quantitatively through simulation, we apply them to the example of pathological non-convergence by Van Buuren (2018) that we reproduced above. We want to be able to draw the same conclusions as we did from visually inspecting the two algorithms. But  $\hat{R}$  and  $AC$  should at least be able to distinguish between the two algorithms plotted in Figure 2. From visual inspection, we know that the typically converged algorithm initially portrayed some signs of non-mixing (around  $t = 2$ ), but intermingled nicely overall. Additionally, there was very little trending. The algorithm with pathological non-convergence showed severe non-mixing, although this gradually improved (from  $t = 7$ ). Moreover, there was a lot of trending as well initially (up-to  $t = 6$ ), after which the chains reached a somewhat more stationary state. To assess whether  $\hat{R}$  and  $AC$  may be appropriate non-convergence identifiers for iterative imputation algorithms, they should identify the same trends as the visual inspection. But at least, the following condition must hold: the methods should indicate worse performance for the mis-specified model with pathological non-convergence than the typically converged algorithm (i.e., higher  $\hat{R}$ - and  $AC$ -values).

In Figure 3A, the chain means from Figure 2 are plotted again—now together. Panel B shows two versions of calculating  $AC$ s applied to the  $\theta$ s of panel A: the default calculation with R function `stats::acf()` (R Core Team 2020), and manual calculation as the correlation between  $\theta$  in iteration  $t$  and  $\theta$  in iteration  $t + 1$ . Panel C shows the traditional computation of  $\hat{R}$ , and panel D shows  $\hat{R}$  as computed by implementing Vehtari et al. (2019)'s recommendations. **[Add labels for the two types of Rhat to the figure!]**

When we look at panel B, we conclude something weird. The  $AC$ -values calculated with the default function indicate equal performance for the typical convergence and the pathological non-convergence (up-to  $t = 5$ ), while there is obvious trending in the  $\theta$ s of the latter. Moreover, the best convergence (as indicated by the lowest  $AC$ -value) is observed at  $t = 2$ , but looking at the chain means in panel A, there should be some signs of trending up-to iteration number seven. After consulting the documentation on `stats::acf()` (R Core Team 2020), we conclude that this  $AC$  function is not suited for iterative imputation algorithms. The function is optimized for performance when  $t \geq 50$  (Box et al. 2015), while the default number of iterations in iterative imputation





**Figure 3.** Convergence diagnostics applied on the imputation algorithms of Figure 2.  $\theta$  = chain mean in  $y_{imp,m}$ .

is often much lower. Therefore, we compute  $AC$  manually, see the solid line. These  $AC$  values indicate non-stationarity as expected. We therefore calculate  $AC$  manually.

**Describe panel C here and why it is not optimal]** The  $\hat{R}$ -values in panel D do meet the requirements specified above.  $\hat{R}$  as computed conform [Vehtari et al. \(2019\)](#) does indeed indicate less signs of non-convergence as the number of iterations goes up. (explain the dip in  $R$ hat values at  $t=2$ . Namely, because we can only use 2 of the 3 tricks by [Vehtari et al. 2019](#), if the number of iterations is very low ( $t < 4$ ). That's why the  $\hat{R}$ s are more similar to the traditional GR.)

In conclusion, we will use  $\hat{R}$  conform [Vehtari et al. \(2019\)](#) and  $AC$  computed manually as non-convergence diagnostics in the simulation study. The diagnostics will be applied to four  $\theta$ s of interest: chain means, chain variances, a quantity of scientific interest, and the first eigenvalue of the variance-covariance matrix.

## Simulation study

The aim of the simulation study is to evaluate the impact of inducing non-convergence in the MICE algorithm on several quantities of scientific interest  $Q$ . Subsequently we will evaluate how well  $\hat{R}$  and  $AC$  perform in identifying the effects of non-convergence [**and performance of different thetas**]. And finally, we will formulate an informed advice on the requirements to safely assume *sufficient* convergence in practice. That is, we assume that the algorithm is *sufficiently* converged when each estimate  $\bar{Q}$  is an unbiased and confidence valid estimate of the corresponding estimand  $Q$ . [**and formulate advice on thetas?**]

Non-convergence will be induced by: 1) increasing the missingness proportion  $p_{mis}$  in dataset  $y$  incrementally, and 2) terminating the imputation algorithm at a varying number of iterations  $T$ . The first set of simulation conditions—the missingness proportions—is chosen to reflect the difficulty of the missingness problem. The underlying assumption is that low missingness proportions lead to quick algorithmic convergence, since there is a lot of information in the observed data. Higher missingness proportions should yield slower convergence. Unless, however, the fraction of missing information is so high that the random component in the imputation algorithm outweighs the information in the observed data. Then, convergence to a stable but highly variable state may be reached instantaneously. We assume that the incremental missingness proportions in our study will result in a corresponding increase in signs of non-convergence.

The assumption inherent to the second set of simulation conditions—the number of iterations—is that terminating the imputation algorithm too early causes non-convergence. Generally, the algorithm will not reach convergence if  $T = 1$ , because the imputed values in the first iteration (at  $t = 1$ ) depend on the initial values of the algorithm (which are sampled randomly from the set of observed datapoints). As the number of iterations increases, the imputation chains should become independent of the initial values, until the point at which adding an extra iteration does not lead to a more converged state. We assume that we can induce non-convergence at least in conditions where  $T$  is smaller than the default number of iterations in `mice`,  $T = 5$ .

## Hypotheses

1. We expect that simulation conditions with a high missingness proportion  $p_{mis}$  and a low number of iterations  $T$  will result in biased, invalid estimates of the quantities of scientific interest,  $Q$ s.
2. We hypothesize that  $\hat{R}$  and  $AC$  will correctly identify signs of non-convergence in those simulation conditions where the  $\bar{Q}$ s are *not* unbiased and confidence valid estimates of the  $Q$ s.
3. We hypothesize that the recommended thresholds to diagnose non-convergence with  $\hat{R}$  ( $\hat{R} > 1.2$ ,  $\hat{R} > 1.1$ , and  $\hat{R} > 1.01$ ) may be too stringent for iterative imputation applications. In an empirical study, where  $\hat{R}$  was used to inform the required imputation chain length, it took as many as 50 iterations to overcome the conventional non-convergence threshold  $\hat{R} > 1.2$ . Yet, scientific estimates were insensitive to continued iteration from  $t = 5$  onward (Lacerda et al. 2007). We

therefore suspect that  $\hat{R}$  may over-estimate signs of non-convergence in iterative imputation algorithms, compared to the validity of estimates. In contrast to this, it may occur that signs of non-convergence are under-estimated by  $\hat{R}$ , in exceptional cases where the initial values of the algorithm are not appropriately over-dispersed (Brooks and Gelman 1998, p. 437). In e.g. mice, initial values are chosen randomly from the observed data, hence we cannot be certain of over-dispersion in the initial values. In practice, we do not expect this to cause problems for identifying non-convergence with  $\hat{R}$ .

4. We expect that high  $AC$  values are implausible in iterative imputation algorithms with typical convergence. That is, after only a few iterations, the randomness induced by the algorithm should effectively mitigate the risk of dependency within chains.
5. **[Something about thetas, that multivariate thetas are better at detecting non-convergence than univariate thetas. The non-convergence diagnostics applied on the chain means and chain variances will fail to identify signs of non-convergence when the multivariate Qs are affected.]**

## Set-up

We investigate non-convergence of the MICE algorithm through model-based simulation in R (version 3.6.3; R Core Team 2020). The number of simulation repetitions is 1000. The simulation set-up is summarized in the pseudo-code below. The complete R script of the simulation study is available from [github.com/gerkovink/shinyMice](https://github.com/gerkovink/shinyMice).

```
# pseudo-code of simulation
1. simulate data
for (number of simulation runs from 1 to 1000)
  for (missingness proportions 5%, 25%, 50%, 75% and 95%)
    2. create missingness
    for (number of iterations from 1 to 100)
      3. impute missingness
      4. perform analysis of scientific interest
      5. compute non-convergence diagnostics
      6. pool results across imputations
      7. compute performance measures
    8. combine outcomes of all missingness proportions
  9. aggregate outcomes of all simulation runs
```

*Data-generating mechanism.* In this study, sampling variance is not of interest. Therefore, a single complete dataset may serve as comparative truth in all simulation repetitions (Vink and van Buuren 2014). The data-generating mechanism is a multivariate normal distribution, representing person-data on three predictor variables in a multiple linear regression problem (from an unspecified social scientific field of study). Let the

predictor space be defined as

$$\begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix} \sim \mathcal{N} \left[ \begin{pmatrix} 12 \\ 3 \\ 0.5 \end{pmatrix}, \begin{pmatrix} 4 & 4 & 1.8 & 0 \\ 4 & 16 & 4.8 & 0 \\ 1.8 & 4.8 & 9 & 0 \end{pmatrix} \right].$$

A finite population of  $N = 1000$  is simulated using the `mvtnorm` package (Genz and Bretz 2009). Subsequently, a fourth variable is constructed as outcome  $Y$ . For each unit  $i = 1, 2, \dots, N$ , let

$$Y_i = 1 + 2X_{1i} + .5X_{2i} - X_{3i} + \epsilon_i,$$

where  $\epsilon \sim \mathcal{N}(0, 100)$ . From the complete set, we obtain the true values of several quantities of scientific interest,  $Q$ s. **[MOVE this next part to the start of the Simulation Study section?]  $Q$ s are the targets that will be estimated in each simulation repetition, after first *amputing* the complete data with varying missingness proportions, then imputing the missingness with a varying number of iterations, then estimating the  $Q$ s on the  $M$  completed datasets, and finally pooling the  $M$  estimates  $\hat{Q}$  to obtain a single  $\bar{Q}$  for each  $Q$ . We try to predict the effect of non-convergence (i.e., the difference between  $Q$  and  $\bar{Q}$ ) with the convergence diagnostics  $\hat{R}$  and  $AC$ , by applying these methods to several scalar summaries of the state space of the algorithm,  $\theta$ s.**

*Scientific estimands.* We consider four types of  $Q$ s that are often of interest in empirical research. Namely, two descriptive statistics: the mean  $\mu_j$  and standard deviation  $\sigma_j$  of each variable  $j = Y, X_1, X_2, X_3$ . We also consider the regression coefficients of the predictors:  $\beta_1, \beta_2$ , and  $\beta_3$  in regression equation

$$Y' = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3,$$

where  $Y'$  is the expected value of the outcome. And finally, the proportion of variance explained by the set of predictors: coefficient of determination  $r^2 \times 100$  (**note that lower case  $r$  is used to avoid confusion with non-convergence diagnostic  $\hat{R}$** ). The true values of these scientific estimands are calculated on the complete data with the R functions `base::mean()`, `stats::sd()`, and `stats::lm()` (R Core Team 2020).

*Methods.* In each simulation repetition, we *ampute* the complete dataset five times to obtain five different incomplete datasets,  $y$ s. The ratio between  $y_{obs}$  and  $y_{mis}$  in each  $y$  depends on the missingness proportion:  $p_{mis} = .05, .25, .5, .75, .95$ . We ampute the complete data with the `mice` package (function `mice::ampute()`; Van Buuren and Groothuis-Oudshoorn 2011). We consider all possible univariate and multivariate patterns of missingness, conform a ‘missing completely at random’ missingness mechanism (Rubin 1987). I.e., the probability to be missing is the same for all  $N \times k$  cells in  $y$ , conditional on the missingness proportion  $p_{mis}$ .

Missing datapoints in each incomplete dataset  $y$  are imputed with the `mice()` function (Van Buuren and Groothuis-Oudshoorn 2011). All imputation procedures are performed with Bayesian linear regression imputation, and five imputation chains ( $M =$

5). Because  $M = 5$ , each run of the `mice()` function results in five sets of imputations:  $y_{imp,m}$ , where  $m = 1, 2, \dots, 5$ . The number of iterations in each run of the algorithm varies between simulation conditions ( $T = 1, 2, \dots, 100$ ).

From the 5 imputations  $y_{imp,m}$  we obtain the chain means and chain variances for each variable  $j$ . These will serve as scalar summary  $\theta$  to apply non-convergence diagnostics  $\hat{R}$  and  $AC$  on. We then use `mice::complete()` to combine the imputed data  $y_{imp,m}$  with the observed data  $y_{obs}$ , resulting in five completed datasets  $\{y_{obs}, y_{imp,m}\}$ .

On the 5 completed datasets  $\{y_{obs}, y_{imp,m}\}$ , we:

- Estimate  $\hat{Q}$ s for  $Q = \mu$  and  $Q = \sigma$  with `base::mean()`, and `stats::sd()` (R Core Team 2020). We pool the  $\hat{Q}$ s with `base::mean()` to get  $\bar{Q}$ s.
- Compute  $\lambda_1$ , to use as  $\theta$  in the two convergence diagnostics. By definition, the first eigenvalue of the variance-covariance matrix,  $\lambda_1$ , is equal to the variance of the principal component solution of each  $\{y_{obs}, y_{imp,m}\}$ , which we calculate with `stats::princomp()` (R Core Team 2020).
- Perform multiple linear regression with `stats::lm()` (R Core Team 2020) to get the  $\hat{Q}$ s for  $Q = \beta$  and  $Q = r^2$ . We pool the  $\hat{Q}$ s conform Vink and van Buuren (2014) to get  $\bar{Q}$ s for the  $\beta$ s and  $r^2$  [check finite pooling of  $r^2$ , instead of `mice::pool.r.squared()`]. Additionally, we use the  $\hat{Q}$ s of  $Q = \beta$  as  $\theta$  to apply the convergence diagnostics to.

We apply non-convergence diagnostics  $\hat{R}$  and  $AC$  to each scalar summary  $\theta$ . We calculate  $\hat{R}$  by implementing Vehtari et al.'s (2019) recommendations, and  $AC$  as the correlation between the  $t^{th}$  and the  $(t + 1)^{th}$  iteration.

*Performance measures.* We assess the impact of inducing non-convergence in the iterative imputation algorithm by comparing  $\bar{Q}$ s with  $Q$ s. For each  $Q$ , we calculate bias as  $\bar{Q} - Q$ . For the regression coefficients  $Q = \beta_{1,2,3}$ , we also compute the empirical coverage rate (CR). CR is defined as the percentage of simulation repetitions in which the 95% confidence interval (CI) around  $\bar{Q}$  covers the true estimand  $Q$ . Let

$$CI = \bar{Q} \pm t_{(M-1)} \times SE_{\bar{Q}},$$

where  $t_{(M-1)}$  is the quantile of a  $t$ -distribution with  $M - 1$  degrees of freedom, and  $SE_{\bar{Q}}$  is the square root of the pooled variance estimate. If we obtain nominal coverage (CR = 95%), we can conclude that the  $\bar{Q}$ s are confidence valid estimates of the  $Q$ s. Finally, we inspect CI width (CIW): the difference between the lower and upper bound of the 95% confidence interval around  $\bar{Q}$ . CIW is of interest because too short CIs indicate underestimation of the variance in  $\bar{Q}$ , which may yield spurious inferences.

With these performance measures, we also evaluate how well the non-convergence diagnostics  $\hat{R}$  and  $AC$  can identify simulation conditions in which the  $\bar{Q}$ s are *not* unbiased, confidence valid estimates of the  $Q$ s. **[Rephrase this:] We measure the performance of  $\hat{R}$  and  $AC$  by the four scalar summaries  $\theta$ : chain means (i.e., the mean in each  $y_{imp,m,j}$ ) versus  $Q = \mu$ , chain variances (i.e., the variance in each  $y_{imp,m,j}$ ) versus  $Q = \sigma$ , the estimated regression coefficients  $\beta$  in each**

$\{y_{obs}, y_{imp,m}\}$  (i.e., a user-defined statistic of scientific interest  $\hat{Q}$ ) against  $Q = \beta$ , and the first eigenvalue of the variance-covariance matrix in each  $\{y_{obs}, y_{imp,m}\}$ ,  $\lambda_1$  versus  $Q = r^2$ .

[Remove this?] Primarily, we want to know what the effect is of inducing non-convergence on the estimated  $Q$ s. More specifically, we are interested in the bias of  $\bar{Q}$  for all  $Q$ s and the confidence validity of  $\bar{Q}$  for  $Q = \beta$ . Therefore, we evaluate these performance measures against the simulation conditions (the number of iterations  $T$  and the missingness proportion  $p_{mis}$ ). We look at averages across simulation repetitions. Secondly, we want to know if conditions in which the  $Q$ s are affected by non-convergence are also identified by the non-convergence diagnostics  $\hat{R}$  and  $AC$ . We do this, again, by looking at the performance measures for the  $Q$ s, but now we evaluate these against what  $\hat{R}$  and  $AC$  indicate about signs of non-convergence in the  $\theta$ s.

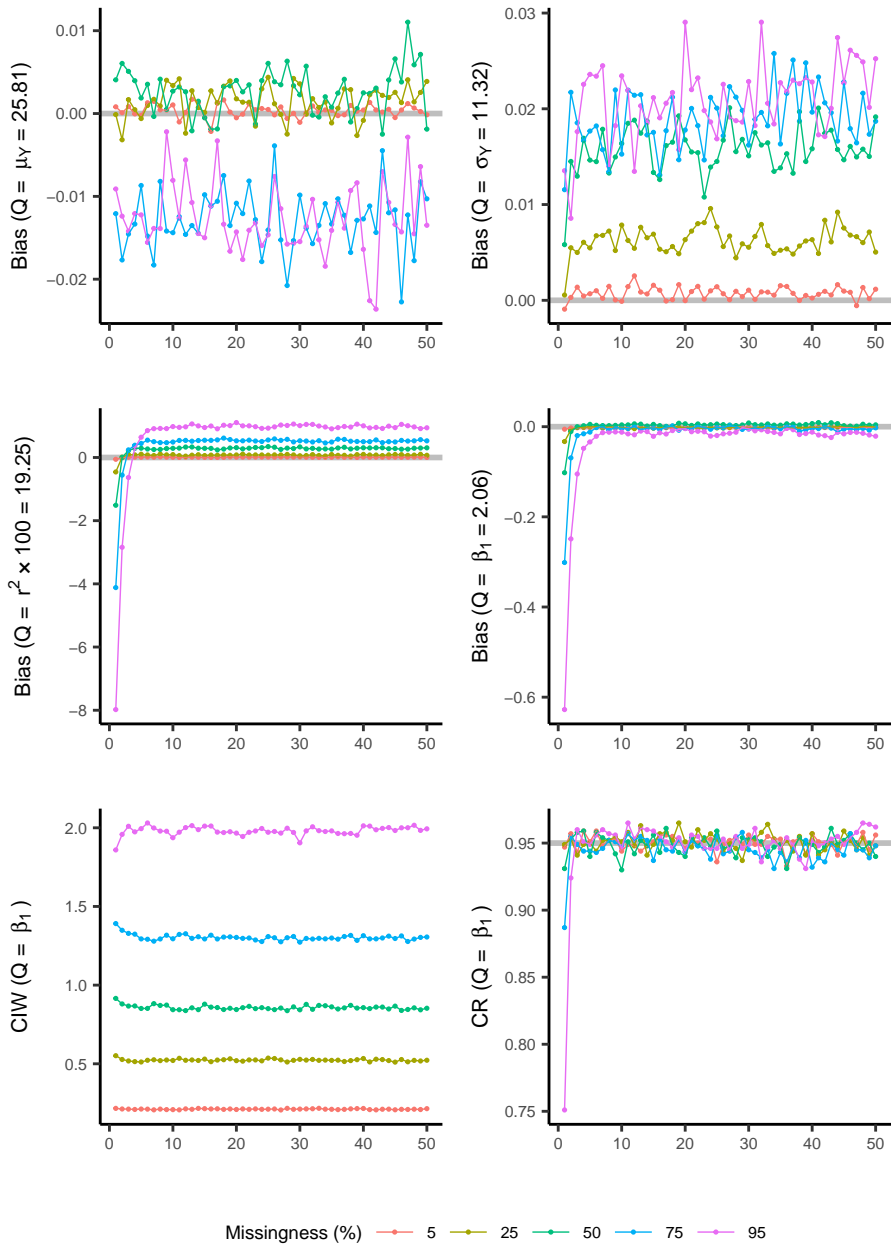
## Results

[Add to this section: 1) emphasize that the iterations in the plots are averages across repetitions, not within one run of MICE; 2) Add more info about figure legends and axes]

Our results show the impact of inducing non-convergence in an iterative imputation algorithm on several scientific estimates, and whether the non-convergence diagnostics  $\hat{R}$  and  $AC$  identify the signs of non-convergence in the algorithm. For reasons of brevity, we only discuss the non-convergence diagnostics for estimates with the worst performance in terms of bias. Bias in the descriptive statistics ( $Q = \mu$  and  $Q = \sigma$ ) is most pronounced for the outcome variable  $Y$ . Of the regression coefficients ( $Q = \beta$ ) we observe the largest bias for  $\beta_1$  (the effect of  $X_1$  on  $Y$ ). There is only one estimate for  $Q = r^2$  to consider. In the figures, we present simulation conditions where the number of iterations is  $1 \leq T \leq 50$ , since the results are more or less stable for conditions where  $T \geq 30$ . Full results are available from [github.com/gerkovink/shinyMice](https://github.com/gerkovink/shinyMice).

### Quantities of scientific interest

Figure 4 displays the influence of inducing non-convergence on the estimated quantities of scientific interest,  $Q$ s. [Add panel labels, describe panels] As expected, conditions with a higher proportion of missingness and/or a lower number of iterations are impacted more severely by non-convergence. I.e., on average these conditions portray more extreme bias in the estimated  $Q$ s, and non-nominal coverages. Roughly speaking, this means that the algorithm is indeed not converged at  $t = 1$ , and converges gradually as  $t$  increases to  $T$ . The point at which an additional iteration does not lead to an improvement in the estimates depends on quantity of scientific interest  $Q$  and the missingness proportion  $p_{mis}$ . The highest number of iterations that is necessary to reach this point across all simulation conditions is  $T = 7$ . This implies that, under the current specifications, no more than seven iterations are necessary to obtain *sufficient* performance in terms of bias and confidence validity.



**Figure 4.** Impact of non-convergence on statistical inferences.

$Q = \mu$ . The estimated univariate means seem unaffected by the number of iterations in the algorithm. The magnitude of the bias in  $\bar{Q}$  depends exclusively on the missingness proportion  $p_{mis}$ , with higher missingness proportions causing more extreme bias. Note however, that the simulation condition with the largest overall bias resulted in a  $\bar{Q}$  under-estimated the true value of  $Q$  by only 0.02 units (where  $Q = \mu_Y = 25.81$ ;  $\sigma_Y = 11.32$ ). We conclude that—conditional on the missingness proportion—all simulation conditions ( $T \geq 1$ ) have equal performance. According to this scientific estimand  $Q$ , non-convergence does not impact the validity of the inferences. **[Add: up-to 50% missingness everything is relatively unbiased, while 75% and 95% missingness show very different results. Why is that? Add sd of Y in the figure to show that the bias is relative: a bias of 0.02 for a variable with a mean of 25.81, and sd of 11.32 is nothing.]**

$Q = \sigma$ . The estimated standard deviations show impact of non-convergence similar to the means. The magnitude of the bias  $\bar{Q}$  follows the incremental order of the missingness proportions. The number of iterations in the algorithm has little effect on the bias. Only the estimates in simulation conditions with just one iteration ( $T = 1$ ) differ from other conditions that have the same missingness proportion but a higher number of iterations. Interestingly, the bias in conditions where  $T = 1$  is *less* severe than after continued iteration. For the least possible bias, we may use any number of iterations  $T \geq 1$ , but only  $p_{mis} = 0.05$  results in completely unbiased  $\bar{Q}$ s.

$Q = \beta$ . For the estimated regression coefficients, we consider bias, confidence interval width (CIW), and coverage rate (CR) as performance measures. The magnitude of the bias, again, depends on the missingness proportion. But for this  $Q$  we also observe a clear effect of the number of iterations. Conditions with a lower number of iterations yield more bias in  $\bar{Q}$ . This is especially pronounced in conditions with high missingness proportions: with 95% missingness we need at least seven iterations before an additional iteration does not lead to improved estimates. Conditions with a lower missingness proportion (e.g.,  $p_{mis} = 0.5$ ), produce *optimal* estimates with as many as two iterations. Approximately unbiased estimates are obtained for all conditions where  $p_{mis} \leq 0.75$  and  $T \leq 6$ . **[Add CIW and CR]**

$Q = r^2$ . The bias in the estimated explained variance perfectly follows both sets of simulation conditions: The worst performance is observed for the highest missingness proportion and the lowest number of iterations, while the best performance occurs with less missingness and more iterations. Only conditions with the lowest missingness proportion ( $p_{mis} = 0.05$ ) results in unbiased estimates for every number of iterations. In other conditions there is at least some bias irrespective of  $T$ . The inferences do not improve in any condition where  $T \geq 6$ . This suggests that the algorithm is as converged as necessary to obtain valid inferences.

### Non-convergence diagnostics

In Figure 5 we see the results of the non-convergence diagnostics  $\hat{R}$  and  $AC$ . **[Add panel labels and explain how to read the plots]** As expected, the diagnostics identify more



signs of non-convergence in simulation conditions with a lower number of iterations. The missingness proportions, however, only seem to impact multivariate convergence. When  $\hat{R}$  and  $AC$  are applied to univariate  $\theta$ s, they identify the same amount of non-convergence, irrespective of the missingness proportion. When the diagnostics are applied to multivariate  $\theta$ s, they generally identify more signs of non-convergence for higher missingness proportions.

None of the simulation conditions that we consider resulted in complete convergence of the algorithm (defined as  $\hat{R} = 1$  and  $AC = 0$ ). The thresholds for diagnosing non-convergence that are used in practice seem too conservative compared to the performance measures:  $\hat{R}$  and  $AC$  identify signs of non-convergence in conditions where the  $\bar{Q}$ s are unbiased and confidence valid estimates of the  $Q$ s. The most recent recommended threshold for diagnosing non-convergence with  $\hat{R}$  is never overcome in this study (all  $\hat{R} > 1.01$ ). Apparently, this threshold is too stringent for iterative imputation algorithms, and will not be discussed further.

$\theta = \text{chain mean}$ . When  $\hat{R}$  and  $AC$  are applied to the chain means, we notice no impact of the missingness proportions. The results only depend on the number of iterations per simulation condition.

**[MAYBE remove all this because it has no relation to the bias in Q] The highest  $\hat{R}$ -values are obtained for the condition  $T = 2$ . We observe the same ‘dip’ in  $\hat{R}$  as we did in Figure 3 for conditions with three iterations.  $\hat{R}$  gradually decreases as the number of iterations increases across conditions. By the traditional threshold  $\hat{R} > 1.2$  we should conclude that there is non-convergence in conditions where  $T < 8$ . The common threshold  $\hat{R} > 1.1$  is overcome in conditions where  $T > 14$ .**

**Autocorrelations indicate some signs of non-stationarity in all simulation conditions, but primarily in conditions where  $T < 6$ . The threshold for diagnosing non-convergence is not reached in any condition. We would thus conclude that the algorithm is sufficiently converged regardless of the number of iterations or missingness proportion.**

We evaluate the diagnostics with chain means as  $\theta$  against the bias in  $Q = \mu$ . According to the bias in the univariate means, we should conclude that two iterations is sufficient, and we should see that  $p_{mis} = 0.75$  and  $p_{mis} = 0.95$  have worse convergence than the other missingness proportions. This is not the case.

$\theta = \text{chain variance}$ . The non-convergence diagnostics applied to the chain variances as  $\theta$  yield highly similar results to the chain means. We therefore evaluate these results immediately against the corresponding  $Q$ , without describing the  $\hat{R}$  and  $AC$  values. While the bias in the estimated standard deviations mainly depends on the missingness proportion of each condition,  $\hat{R}$  and  $AC$  are influenced mostly by the number of iterations. We thus conclude that there is no added value of evaluating non-convergence with chain variances as  $\theta$  in the diagnostics.

$\theta = \hat{Q}$ . The non-convergence diagnostics are applied to the estimated regression coefficients in each imputation, and evaluated against the bias in the corresponding pooled estimate.

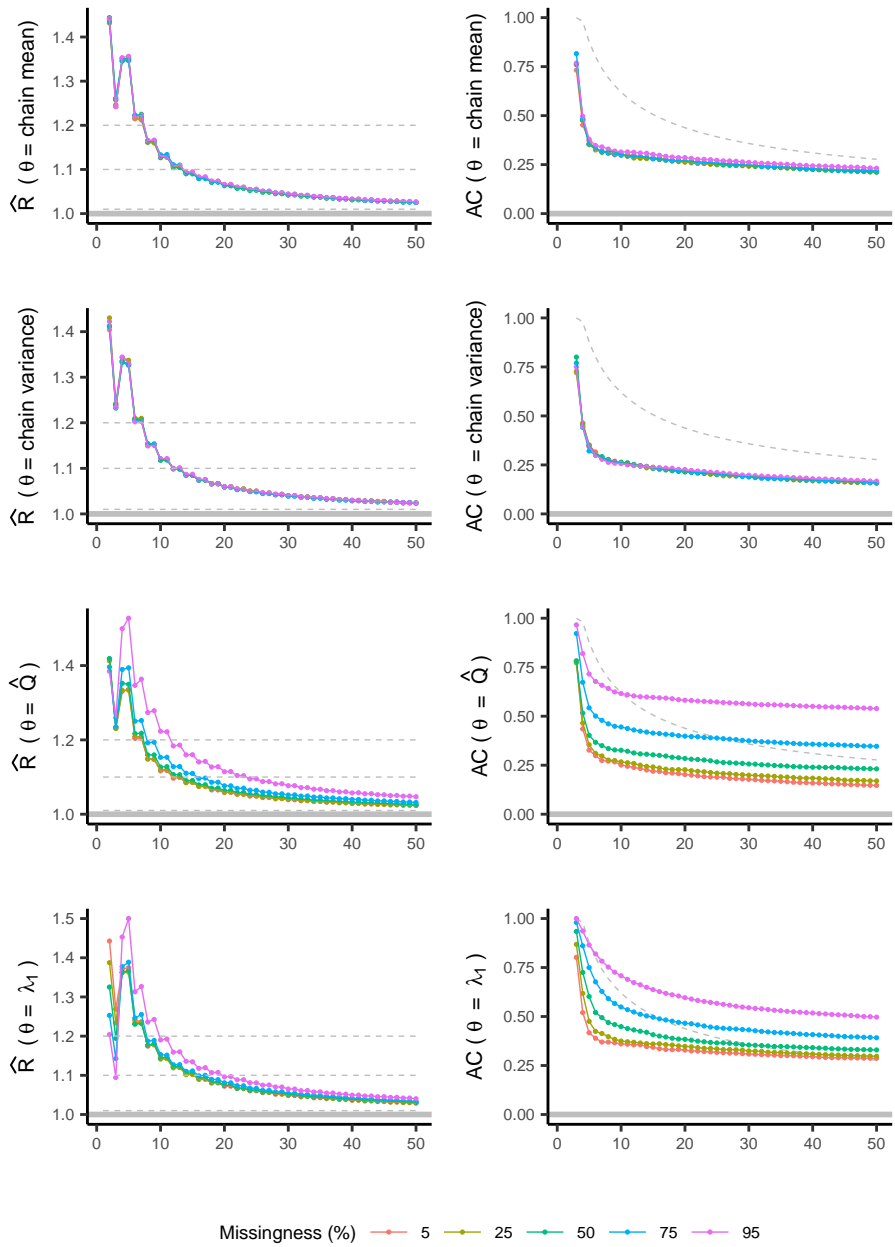


Figure 5. Non-convergence diagnostics. [Split up into 2 or 4 figures??]

$\hat{R}$  indicates signs of non-convergence across all missingness proportions and each number of iterations, in the expected order: simulation conditions with higher missingness proportions and less iterations are identified as the worst converged. The thresholds to diagnose non-convergence are overcome in different conditions, depending on the number of iterations and the missingness proportion. The most lenient threshold ( $\hat{R} > 1.2$ ) is overcome at  $T \geq 9$  for  $p_{mis} = 0.05$ , and  $T \geq 12$  for  $p_{mis} = 0.95$ . The conventional  $\hat{R} > 1.1$  is surpassed at  $T \geq 12$  and  $T \geq 24$ , respectively.

For the autocorrelations applied to the  $\hat{Q}$ s, we see an even stronger effect of the missingness proportions. Especially the conditions where  $p_{mis} = 0.75$  and  $p_{mis} = 0.95$  are identified as non-stationary. The threshold used to evaluate whether  $AC = 0$ , however, is only reached in conditions with at least ten iterations. This suggests that the default number of iterations may be too low to detect non-stationarity with  $AC$ . However, the critical value for diagnosing non-convergence with  $AC$  decreases as a function of the number of iterations. The conclusion to identify non-stationarity in conditions where  $T \geq 10$  and  $p_{mis} = 0.95$  might therefore be spurious. **[Also, the  $AC$  across missingness proportions shows the opposite of what you would expect: shouldn't there be higher autocorrelations in conditions with comparatively more observed data in the completed datasets? I don't understand this.]**

Compared to the bias in  $\bar{Q}$ , the diagnostics applied to the  $\hat{Q}$ s over-estimate the severity of the non-convergence. This is odd, since the  $\theta$  and  $Q$  are directly related.

$\theta = \lambda_1$ . Finally, we evaluate the non-convergence diagnostics applied to the first eigenvalue of the variance-covariance matrices of the completed data. Similar to the results the regression coefficients as  $\theta$ s, we see effects of both the missingness proportion and number of iterations in each simulation condition. The most signs of non-convergence are identified in conditions where  $p_{mis}$  is large and  $T$  is small.

Notwithstanding the similarities with  $\theta = \hat{Q}$ , there is a big difference in  $\hat{R}$ -values for conditions where  $T < 4$ . In these conditions,  $\hat{R}$  behaves oppositely of our expectation: conditions with the lowest missingness proportions yield the highest  $\hat{R}$ -values. In conditions with a higher number of iterations, this artifact disappears. This suggests that the combination of  $\hat{R}$  and  $\lambda_1$  is only appropriate for imputation algorithms where the number of iterations is at least four. Interpreting the  $\hat{R}$  values of conditions where  $T < 3$  may lead to the incorrect conclusion that only conditions with *low* missingness proportions show signs of non-convergence. For applications as a model-independent multivariate  $\theta$ , we should therefore only interpret the  $\hat{R}$ -values when  $T \geq 4$ . Using this cut-off we conveniently ignore the initial 'dip' at  $T = 3$  as well. The traditional threshold  $\hat{R} > 1.2$  is overcome in conditions where  $p_{mis} = .95$  and  $T \geq 10$ , and for all other missingness proportions in conditions where  $T \geq 8$ . With the  $\hat{R} > 1.1$  threshold, the number of iterations are  $T \geq 20$  and  $T \geq 16$ , respectively.

The  $AC$ -values largely depend on the missingness proportion in each condition. Conditions where  $p_{mis} = .95$  show statistically significant autocorrelations for any number of iterations. The other extreme, conditions where  $p_{mis} = .05$ , does not yield  $AC$  values above the threshold. The  $AC$  values with this missingness proportion are equivalent for all conditions where  $T \geq 6$ , which suggests that the algorithm is as

stationary as it can be. A similar point is only observed after twenty iterations in conditions with higher missingness proportions.

In comparison to the  $Q$ s, the diagnostics again over-estimate the signs of non-convergence. Conditions with low missingness proportions (e.g., 5% or 25% of the complete data) yield unbiased, confidence valid estimates with as little as two iterations. The convergence diagnostics, however, indicate improved convergence up-to eight iterations according to  $AC$  and thirty to forty iterations according to  $\hat{R}$ . The model-independent  $\theta \lambda_1$  seems as well-equipped to diagnose signs of non-convergence as  $\theta = \hat{Q}$ . The univariate  $\theta$ s do not differentiate between bias induced by the missingness proportions.

### *In short [remove later!]*

- Is there more bias in conditions where non-convergence was induced? (NOPE for univariates, YES for multivar.).
- Is the point of *sufficient* convergence (at which the  $\bar{Q}$ s are unbiased and confidence valid estimates of the  $Q$ s) the same point as identified by the non-convergence diagnostics? (NOPE, they are too conservative).
- Is one of the thetas better? (YES, while it doesn't really matter as long as it's multivariate,  $\lambda_1$  is model-independent and therefore preferred).
- What should the thresholds be? (much HIGHER than recommended for  $\hat{R}$ , difficult to say for  $AC$ ).
- And the default nr of iterations? (MAYBE increase to 10 to detect significant  $AC$  and be sure not to miss the dip in  $\hat{R}$ ?).

## Discussion

In short, iterative imputation algorithms may yield approximately unbiased, confidence valid estimates under different levels of induced non-convergence. Univariate estimates seem robust against terminating the algorithm early: There is no clear effect of the number of iterations on the bias in estimated means and standard deviations. The bias in these  $Q$ s was only impacted by the missingness proportion  $p_{mis}$ , with higher  $p_{mis}$  resulting in more bias. Bias in multivariate estimates depended on both  $p_{mis}$  and the number of iterations. However, the only effect of the number of iterations was observed when the algorithm was terminated at seven iterations or earlier. This suggests that continued iteration beyond seven iterations does not yield better estimates. Yet, the non-convergence diagnostics  $\hat{R}$  and  $AC$  indicated that the iterative imputation algorithm did not reach a stable state before twenty to thirty iterations. **[MERGE or remove]  $\hat{R}$  and autocorrelation indicate that there are signs of non-convergence in the algorithm up-to at least twenty iterations, while unbiased, confidence valid estimates may be obtained with as little as one iteration. This observation is in agreement with one of [still add the others!]\*\* the simulation hypotheses:  $\hat{R}$  over-estimates the severity of non-convergence when applied to MI procedures.\*\***

Over-estimation of the signs of non-convergence may be due to the methods (and their thresholds) or due to the  $Q$ s (descriptive and multivariate linear regression, not

higher dimensional/more complex  $Q$ s). More complicated  $Q$ s (e.g., higher-order effects or variance components) might show bias, under- or over-coverage at higher  $T$ . Which might be why  $\hat{R}$  and  $AC$  identify non-convergence.

### *Recommendations for empirical researchers*

For empirical researchers:

1. Check traceplots for signs of pathological non-convergence. Adjust imputation model if necessary.
2. Monitor non-convergence diagnostics  $\hat{R}$  and  $AC$ . Use the thresholds for  $\hat{R} > 1.2$  and autocorrelation ????. Keep iterating until these thresholds are reached.
3. Do not use the traditional calculation for  $\hat{R}$ , or the R function `stats::acf()` to compute  $AC$ . Instead, calculate  $\hat{R}$  conform [Vehtari et al. \(2019\)](#) and compute autocorrelations manually (see e.g., [github.com/gerkovink/shinyMice](https://github.com/gerkovink/shinyMice)).
4. Track your own scalar summary of interest (see e.g., [Van Buuren \(2018\)](#)). Compute  $\hat{R}$  and autocorrelation values for this scalar summary.
5. **Something** about the novel  $\theta$  that is ‘substantive model-independent’.

### *Recommendations for future research*

Further research is needed to investigate the performance of iterative imputation algorithms under clear violation of convergence, e.g. dependency between predictors (predictors with very high correlations). For example:

- Induce non-convergence with mis-specified imputation model.
- Introduce more difficult missingness problems, i.e., M(N)AR instead of MCAR. That is, proper performance under a ‘missing completely at random’ missingness mechanism is necessary to demonstrate the appropriateness of  $\hat{R}$  and  $AC$  as non-convergence diagnostics. However, results may not be extrapolated to other missingness mechanisms.
- Imputation technique with questionable convergence properties, e.g., PMM.
- Apply recommendations to empirical data and write a vignette for applied researchers.
- Implement recommendations in ShinyMICE.
- Develop an additional convergence diagnostic unique for iterative imputation:  $AC$  across imputations, instead of iterations.

In short, we have shown that an iterative imputation algorithm can yield correct outcomes, even when a converged state is not reached according to two common non-convergence diagnostics. Continued iteration just burns a lot of computing time, without improving the statistical inferences. To conclude, in the case of drawing inference from incomplete data, convergence of the iterative imputation algorithm is often a convenience but not a necessity.

## References

- Box GEP, Jenkins GM, Reinsel GC and Ljung GM (2015) *Time Series Analysis: Forecasting and Control*. John Wiley & Sons. ISBN 978-1-118-67492-5. Google-Books-ID: rNt5CgAAQBAJ.
- Brooks SP and Gelman A (1998) General Methods for Monitoring Convergence of Iterative Simulations. *Journal of Computational and Graphical Statistics* 7(4): 434–455. DOI: 10.1080/10618600.1998.10474787.
- Cowles MK and Carlin BP (1996) Markov chain Monte Carlo convergence diagnostics: a comparative review. *Journal of the American Statistical Association* 91(434): 883–904.
- Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A and Rubin DB (2013) *Bayesian Data Analysis*. Philadelphia, PA, United States: CRC Press LLC.
- Gelman A and Rubin DB (1992) Inference from Iterative Simulation Using Multiple Sequences. *Statistical Science* 7(4): 457–472. DOI:10.1214/ss/1177011136.
- Genz A and Bretz F (2009) *Computation of multivariate normal and t probabilities*. Lecture notes in statistics. Heidelberg: Springer-Verlag. ISBN 978-3-642-01688-2.
- Geweke J (1992) Evaluating the accuracy of sampling-based approaches to the calculations of posterior moments. *Bayesian statistics* 4: 641–649.
- Hoff PD (2009) *A First Course in Bayesian Statistical Methods*. Springer Texts in Statistics. New York, NY: Springer New York. DOI:10.1007/978-0-387-92407-6.
- Lacerda M, Ardington C and Leibbrandt M (2007) Sequential regression multiple imputation for incomplete multivariate data using Markov chain Monte Carlo. Technical report, University of Cape Town, South Africa.
- Lynch SM (2007) *Introduction to applied Bayesian statistics and estimation for social scientists*. Springer Science & Business Media.
- MacKay DJ and Mac Kay DJ (2003) *Information theory, inference and learning algorithms*. Cambridge university press.
- Murray JS (2018) Multiple Imputation: A Review of Practical and Theoretical Findings. *Statistical Science* 33(2): 142–159. DOI:10.1214/18-STS644.
- R Core Team (2020) *R: A language and environment for statistical computing*. Vienna, Austria.
- Raftery AE and Lewis S (1991) How many iterations in the Gibbs sampler? Technical report, Washington University Seattle, Department of Statistics, United States.
- Rubin DB (1976) Inference and Missing Data. *Biometrika* 63(3): 581–592. DOI:10.2307/2335739.
- Rubin DB (1987) *Multiple Imputation for nonresponse in surveys*. Wiley series in probability and mathematical statistics Applied probability and statistics. New York, NY: Wiley.
- Schafer JL (1997) *Analysis of incomplete multivariate data*. Chapman and Hall/CRC.
- Takahashi M (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: 37. DOI:10.5334/dsj-2017-037.
- Van Buuren S (2018) *Flexible imputation of missing data*. Chapman and Hall/CRC.
- Van Buuren S and Groothuis-Oudshoorn K (2011) mice: Multivariate Imputation by Chained Equations in R. *Journal of Statistical Software* 45(1): 1–67. DOI:10.18637/jss.v045.i03.

Vehtari A, Gelman A, Simpson D, Carpenter B and Bürkner PC (2019) Rank-normalization, folding, and localization: An improved  $\widehat{R}$  for assessing convergence of MCMC URL <http://arxiv.org/abs/1903.08008>.

Vink G and van Buuren S (2014) Pooling multiple imputations when the sample happens to be the population URL <http://arxiv.org/abs/1409.8542>.