

---

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

Sociological Methods & Research  
2020, Vol. 49(?):1–13  
©The Author(s) 2020  
Article reuse guidelines:  
sagepub.com/journals-permissions  
DOI: 10.1177/ToBeAssigned  
journals.sagepub.com/home/smr

SAGE

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

## Abstract

(**Rewrite:**) Multiple imputation by chained equations (mice) is a widely used tool to accommodate missing data. While empirical evidence supports the validity of inferences obtained using mice, there is no consensus on the convergence properties of the method. This paper provides insight into non-convergence of mice algorithms.

## Keywords

MICE; convergence

## Introduction

Missing data problems are ubiquitous and pervasive in the social and behavioral sciences. If a dataset contains just one missing value for a variable of interest, statistical inferences are undefined and will not yield results. Incomplete observations are therefore ignored by default in many statistical packages (i.e., list-wise deletion is employed). Unfortunately, this *ad hoc* solution may yield wildly invalid results (Van Buuren 2018). An alternative is to ‘impute’ (i.e., fill in) every missing value in the incomplete dataset. With imputation techniques, one or several completed datasets are created, on which statistical inferences can be performed. The case with several completed datasets is known as ‘multiple imputation’ (MI; Rubin 1976), and requires an additional step to pool the results across imputations. Because this pooled result will reflect the uncertainty in the data due to missingness, it is in many cases an unbiased and confidence valid estimate of the true (but missing) data inference (Van Buuren 2018). A popular method to obtain imputations

---

<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)

is to use the ‘Multiple Imputation by Chained Equations’ algorithm, shorthand ‘MICE’. MICE is an iterative algorithmic procedure to draw imputations from the posterior predictive distribution of the missing values. This 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 open questions, because the convergence of iterative imputation algorithms has not been systematically studied (Van Buuren 2018). In this paper, we investigate the convergence properties of iterative imputation procedures by means of model based simulation in R `list(title = “R: A Language and Environment for Statistical Computing”, author = list(list(given = “R Core Team”, family = NULL, role = NULL, email = NULL, comment = NULL)), organization = “R Foundation for Statistical Computing”, address = “Vienna, Austria”, year = “2020”, url = “https://www.R-project.org/”).`

Rubin’s rules (1976) may be used to reflect the uncertainty in the inferences due to missingness by pooling the results across ‘imputations’. A popular method to obtain imputations completions, the analysis of scientific interest can be performed. It has been , or to complete these observations by imputing the missing values. Imputation is the method of filling in

by which That is, unless ad hoc solutions are employed (e.g., list-wise deletion), which can yield invalid results. contains missing values, statistical analyses cannot be performed without only be performed a dataset contains missing values, statistical analyses can statistical inferences cannot be performed

Missingness is problematic because statistical inference cannot be performed on the whole data without employing *ad hoc* solutions (e.g., list-wise deletion), which may yield wildly invalid results (Van Buuren 2018). A popular method to obtain statistically valid results despite of missing data is ‘multiple imputation’ (MI; Rubin 1976). The multiple imputation method can be roughly divided into 3 steps: First, each missing value in an incomplete dataset is ‘imputed’ (i.e., filled in) several times by some algorithmic procedure. Next, the analysis of scientific interest is performed on each of the completed datasets. And finally, the results of these analyses are pooled into a single aggregated statistic.

Algorithms are often used to draw imputations from the posterior predictive distribution of the missing values. This 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?

There is no scientific consensus on the convergence properties of MI algorithms (Takahashi 2017). Some default techniques (e.g., ‘predictive mean modeling’) might not yield converged states at all (Murray 2018). Therefore, algorithmic convergence should be monitored carefully. The current practice is visually inspecting imputation chains for signs of non-convergence, but this may be challenging to the untrained eye (Van Buuren 2018, § 6.5.2). Moreover, visually inspecting imputation chains may only diagnose severely pathological cases of non-convergence. Therefore, diagnostic evaluation of convergence would be preferred.

Iterative imputation algorithms such as mice, are MCMC methods, and convergence of MCMC algorithms is not from a scalar to a point but from one distribution to another.

Therefore, generated values will vary even after convergence. The aim of convergence diagnostics for MCMC methods is thus not to establish the point at which convergence is reached, but to diagnose severe non-convergence.

Several MCMC convergence diagnostics exist, but it is not known whether these are appropriate for MICE. The application of convergence diagnostics to MI methods has not been systematically studied (Van Buuren 2018). The open questions are, e.g.: How can non-convergence be diagnosed? How many iterations are sufficient/needed? What are the effects of non-convergence on estimates, predictions and inference? What is an appropriate non-convergence diagnostic? Are the default number of iterations sufficient? The defaults of common MI software: SPSS = 10 (link), Mplus = 100?? (link), Stata = 10 (link), Amelia = NA, because resampling, not convergence (link), en MI = 30 (link).

For reasons of brevity, we only focus on the MI algorithm implemented in the world-leading MI software: the R (R Core Team 2020) package `mice` (Van Buuren and Groothuis-Oudshoorn 2011). The convergence properties of this MI algorithm are investigated through model-based simulation. The results of this simulation study are guidelines for assessing convergence of MI algorithms, which will aid applied researchers in drawing valid inference from incomplete datasets.

## Notation

Let  $y$  denote an  $n \times p$  matrix containing the data values on  $p$  variables for all  $n$  units in a sample. The data value of unit  $i$  ( $i = 1, \dots, n$ ) on variable  $j$  ( $j = 1, \dots, p$ ) may be either observed or missing. The collection of observed data values in  $y$  is denoted by  $y_{obs}$ ; the missing part of  $y$  is referred to as  $y_{mis}$ . (**This is not only notation, but theory too:** For each datapoint in  $y_{mis}$ , we sample  $M \times T$  times plausible values, where  $M$  is the number of imputations ( $m = 1, \dots, M$ ) and  $T$  is the number of iterations ( $t = 1, \dots, T$ ).) 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’.

**Add: missingness percentage = proportion of cases with one or more missing values times 100**

## Convergence diagnostics

**(Start with current guidelines, then continue with possible quant measures of the same things from MCMC lit.)** In iterative algorithmic procedures (Markov chain Monte Carlo methods; e.g., `mice` algorithms or Gibbs samplers) non-convergence may be present as non-stationarity within chains (i.e., trending), or as slow mixing between chains (i.e., no intermingling (source??)). **(add example/guidelines visual inspection here?)**

The stationarity component of convergence may be evaluated with autocorrelation (AC; Schafer 1997; Gelman et al. 2013), numeric standard error (or ‘MC error’; Geweke 1992), and Raftery and Lewis’s (1991) procedure to determine the effect of trending within chains **(or the chain length required to cancel out trending?)**.

The mixing component can be assessed with the potential scale reduction factor  $\hat{R}$  (a.k.a. ‘Gelman-Rubin statistic’; Gelman and Rubin 1992). With an adapted version of

$\hat{R}$ , proposed by [Vehtari et al. \(2019\)](#), the stationarity component of convergence might be evaluated as well. This would make  $\hat{R}$  a general convergence diagnostic. The application of  $\hat{R}$  to assess stationarity has not been thoroughly investigated. Therefore, this study employs both  $\hat{R}$  and autocorrelation to investigate convergence, as recommended by [\(Cowles and Carlin 1996, p. 898\)](#).

Note that all of these methods evaluate the convergence of univariate scalar summaries (e.g., chain means or variances). These convergence diagnostics cannot diagnose convergence of multivariable statistics (i.e., relations between scalar summaries). [Van Buuren \(2018\)](#) proposed to implement multivariable evaluation through eigenvalue decomposition ([MacKay and Mac Kay 2003](#)). This method is outside of the scope of the current study (**not anymore?**).

### Potential scale reduction factor

To define  $\hat{R}$ , we follow notation by [\(Vehtari et al. 2019, p. 5\)](#). 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 mean or chain variance). 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 over-estimates the total variance of  $\theta$ .  $\hat{R}$  is then obtained as a ratio between the over-estimated total variance and the within-chain variance:

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

We can interpret  $\hat{R}$  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 all chains sample independent of their initial values, the mixing component of convergence is satisfied, and  $\hat{R}$ -values will be close to one. High  $\hat{R}$ -values thus indicate

non-convergence. The conventionally acceptable threshold for convergence was  $\hat{R} < 1.2$  (Gelman and Rubin 1992). More recently, Vehtari et al. (2019) proposed a more stringent threshold of  $\hat{R} < 1.01$ .

### 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  $t^{th}$  and  $t + 1^{th}$  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. Positive  $AC$ -values, however, 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. Moreover, the magnitude of  $AC$ -values may be evaluated statistically, but that is outside of this note’s scope.

Negative  $AC$ -values indicate divergence within imputation chains. Subsequent sampled values within each imputation chain are less alike. It may not be wise (**why? and how would we know? e.g., wider distributions? didn’t seem like it...**) to terminating the algorithm when  $AC$  is negative. **High  $AC$ -values are implausible in MI procedures. That is, the randomness induced by the MI algorithm effectively mitigates the risk of dependency within chains.**

In short, convergence is reached when there is no dependency between subsequent iterations of imputation chains ( $AC = 0$ ), and chains intermingle such that the only difference between the chains is caused by the randomness induced by the algorithm ( $\hat{R} = 1$ ).

### Simulation Hypothesis

This study evaluates whether  $\hat{R}$  and  $AC$  could diagnose convergence of multiple imputation algorithms. We assess the performance of the two convergence diagnostics against the recommended evaluation criteria for MI methods (i.e., average bias, average confidence interval width, and empirical coverage rate across simulations; Van Buuren 2018, § 2.5.2). That is, there is no baseline measure available to evaluate performance against.

Based on an empirical finding (Lacerda et al. 2007), we hypothesize that  $\hat{R}$  will over-estimate non-convergence of MI algorithms. The threshold of  $\hat{R} < 1.01$  will then be too stringent for diagnosing convergence. This over-estimation may, however, be diminished

because  $\hat{R}$  can falsely diagnose convergence if initial values of the algorithm are not appropriately over-dispersed (Brooks and Gelman 1998, p. 437). In `mice`, initial values are chosen randomly from the observed data. Therefore, we cannot be certain that the initial values are over-dispersed. We expect this to have little effect on the hypothesized performance of  $\hat{R}$ . **(Add actual hypothesis.)** No hypothesis was formulated about the performance of  $AC$  as convergence diagnostic.

## Simulation study

The convergence of `mice` is investigated through model-based simulation in R (version 3.6.3; R Core Team 2020). 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).

A finite population of  $N = 1000$  is simulated to solve a multiple linear regression problem, where dependent variable  $Y$  is regressed on independent variables  $X_1$ ,  $X_2$  and  $X_3$  **(check notation with betas here, and add what the quantity/-ies of scientific interest is/are!)**:

$$Y \sim \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3.$$

The data generating model is a multivariate normal distribution

$$\begin{pmatrix} X_1 \\ X_2 \\ X_3 \\ \epsilon \end{pmatrix} \sim N \left[ \begin{pmatrix} 12 \\ 3 \\ 0.5 \\ 0 \end{pmatrix}, \begin{pmatrix} 4 & 4 & 1.8 & 0 \\ 4 & 16 & 4.8 & 0 \\ 1.8 & 4.8 & 9 & 0 \\ 0 & 0 & 0 & 100 \end{pmatrix} \right]$$

Outcome variable  $Y$  is subsequently calculated as

$$Y = 1 + 2X_1 + .5X_2 - X_3 + \epsilon.$$

**Add the function `mvtnorm::rmvnorm()` back in** use citation(“mvtnorm”)

The complete data is ‘amputed’ once for each simulation repetition with function `mice::ampute()`. **(change this:** The missingness is univariate, and the probability to be missing is the same for all four variables, namely 20% (`prop = 0.8`, `mech = "MCAR"`). This leaves 20% of the rows completely observed). **(Add different percentages missingness??)** This study only considers only an MCAR missingness mechanism. Therefore, results may not be extrapolated to other missing data problems. Proper performance of the convergence diagnostics under MCAR is necessary but not sufficient to demonstrate appropriateness of  $\hat{R}$  and  $AC$  as convergence diagnostics. This is just a proof of concept.

Missing datapoints are imputed with the function `mice::mice()`. All MI procedures are performed with Bayesian linear regression imputation (`method =`

"norm"), and five imputation chains ( $m = 5$ ). The number of iterations varies between simulation conditions ( $\text{maxit} = 1, 2, \dots, 100$ ).

$\hat{R}$  (of what?? i.e., chain means and chain variances of each variable) is computed by implementing Vehtari et al.'s (2019) recommendations.  $AC$  (of what? and how is  $AC$  aggregated across imputations? i.e., mean in latest sim) is computed with function `stats::acf()`.

To estimate the quantity of scientific interest,  $Q$ , we perform multiple linear regression on each completed dataset with the function `stats::lm()`. We obtain an estimated regression coefficient per imputation, which are pooled into a single estimate,  $\bar{Q}$ . We use the function `mice::pool()` to estimate  $\bar{u}$  (check wat dat nou precies is!!) according to Rubin's (1987) rules, and subsequently implement finite population pooling conform Vink and van Buuren (2014).

We compute bias as the difference between  $Q$  and  $\bar{Q}$ . Confidence interval width (CIW) is defined as the difference between the lower and upper bound of the 95% confidence interval (CI95%) around  $\bar{Q}$ . We compute the CI95% bounds as

$$\bar{Q} \pm t_{(m-1)} \times SE_{\bar{Q}},$$

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. **Why track CIW?** Underestimating the variance of  $\bar{Q}$  may yield spurious inferences. From bias and CIW, we calculate empirical coverage rates. Coverage rate is the proportion of simulations in which  $Q$  is between the bounds of the CI95% around  $\bar{Q}$ .

## Results

(Add more info about figure legends and axes.)

### *Univariate estimates and convergence diagnostics (theta = mean and theta = variance)*

The bias in the estimates of the variable means show little to no difference between simulation conditions. It doesn't seem to matter how many iterations you use in the mice algorithm, the estimates are unbiased. Similarly, the bias in the estimated variances is more or less stable across simulation conditions. These univariate quantities appear to be unaffected by the number of iterations.

When applied to the imputation chain means,  $\hat{R}$  indicates that the mice algorithm does not reach a converged state ( $\hat{R} = 1$ ) in any of the simulation conditions. Neither is the most recent recommended threshold reached ( $\hat{R} < 1.01$ ). The conventional  $\hat{R}$  threshold of 1.2 is reached in simulation conditions  $T = 3$  and  $T > 6$ . (With the default number of iterations ( $\text{maxit} = 5$ ), this dip in  $\hat{R}$  values would be spotted, so it is no problem.) The point at which an extra iteration does not seem to improve the  $\hat{R}$  value is around  $T = 30$ .

Auto-correlations indicate no sign of trending within imputation chains. In most simulation conditions  $AC$  is smaller than or about equal to zero. Across simulation

conditions, however, the autocorrelation curve does not trend towards 0. Autocorrelation values plateau off at a value of around .1. This is a small positive autocorrelation, which would indicate some trending within chains.

The  $\hat{R}$  and  $AC$  values for the imputation chain variances show equal trends and are therefore not discussed separately. Taken together, univariate estimates seem robust across simulation conditions. There is no clear effect of the number of iterations on the bias in these estimates, while the convergence diagnostics indicate that the algorithm did not reach a completely converged state (yet).

### *Multivariate estimates*

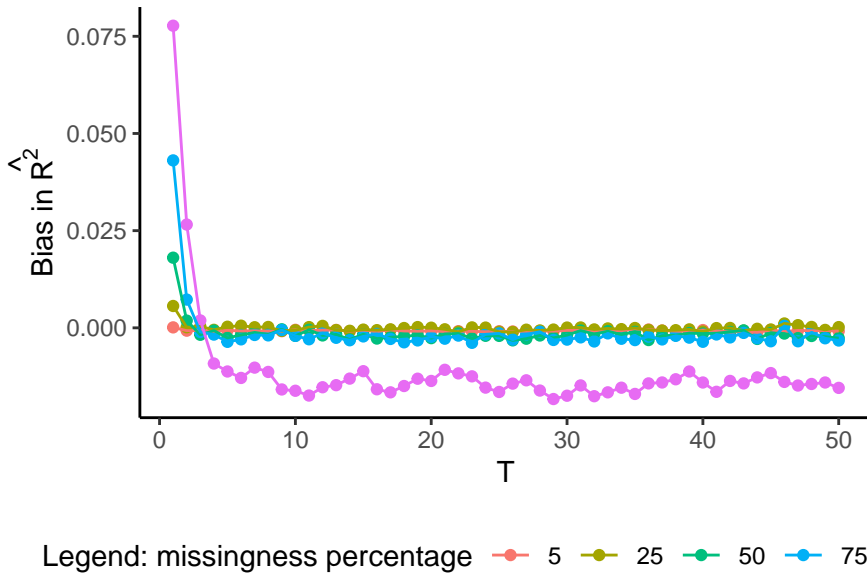
There is a clear bias in the regression estimates in simulation conditions where the number of iterations is smaller than four. In simulation conditions where  $T > 5$  there is little to no bias in the estimated regression coefficients. In most simulation conditions nominal coverage is obtained (i.e., coverage rates of .95) for the confidence intervals around the regression coefficients. Conditions with only one or two iterations show some undercoverage. Since confidence interval width is stable across conditions, the undercoverage may be attributed to the bias in the estimated regression coefficients.

If we look at the estimated proportion of explained variance in outcome variable  $Y$  we see that the coefficient of determination is underestimated in conditions where the number of iterations is equal to two or less, and slightly overestimated in conditions where the number of iterations is equal to three or more.

In short, we see that the minimum number of iterations required to obtain unbiased, confidence valid regression estimates is 5. This value, however, is dependent on the percentage of missing values. E.g., with 95% of cases having missing data we need at least seven iterations to obtain unbiased results.



## The effect of missingness percentage, an example



## Results [OLD, JUST FOR REFERENCE!]

### Convergence Diagnostics

(Pair convergence diagnostics with respective estimates, not split between convergence and simulation diagnostics)

It is apparent that there is a relation between the number of iterations per simulation condition ( $T$ ) and the convergence diagnostics. Generally speaking, conditions with longer imputation chains (higher  $T$ ) coincide with less signs of non-convergence ( $\hat{R}$ -values approach one, and  $AC$ -values approach zero).

Figure 2A shows that  $\hat{R}$ -values generally decrease with increasing imputation chain lengths. The decline stabilizes somewhere between the simulation conditions  $T = 30$  and  $T = 50$ . The downward trend is most pronounced for  $T = 3$ , and between  $T = 5$  and  $T = 10$ . In the intervening conditions ( $3 \leq T \leq 5$ ), however, we observe a steep increase in  $\hat{R}$ -values. This increase implies that non-convergence is under-estimated, and convergence should not be diagnosed at this point (**is that true???**). Based on the conventional threshold  $\hat{R} < 1.2$ , we would falsely diagnose convergence at  $T = 3$ . According to the widely used threshold  $\hat{R} < 1.1$ , convergence would be diagnosed for conditions where  $T > 9$ . If we use the recently recommended threshold  $\hat{R} < 1.01$ , we would conclude that convergence is reached in none of the simulation conditions.

The  $AC$ -values displayed in Figure 2B are almost uniformly increasing as a function of  $T$ . The only  $AC$ -value that deviates from this observation is for  $T = 2$ . The lowest

$AC$ -value is obtained for  $T = 3$ . %At  $T = 5$ , the  $AC$ -value reaches the level observed at  $T = 2$ . The gradual increase plateaus between  $T = 10$  and  $T = 30$ .  $AC$ -values in conditions where  $T > 70$  are indifferentiable from zero, indicating stationarity. We see an initial decrease between  $T = 2$  and  $T = 3$ . Simulation conditions where  $T > 5$  have  $AC$ -values greater than  $T = 2$  (**weird wording**).

According to this diagnostic, none of the simulation conditions show signs of non-convergence, since we only observe negative or zero  $AC$ -values.

Taken together, we see that  $T > 3$  is the minimal requirement to diagnose convergence ( $\hat{R} < 1.2$ ;  $AC \leq 0$ ). This threshold is, however, not sufficient, since we overlook the increase in  $\hat{R}$ -values up-to conditions where  $T > 5$ , and the convergence diagnostics only reach stability at  $T > 20$ .

### Simulation Diagnostics

We use average bias, average confidence interval width, and coverage rate as performance measures to evaluate  $\hat{R}$  and  $AC$  (**not the goal! we want to show the effects of non-conv, not just see how to diagnose it**). We make the general observation that the simulation diagnostics behave as theorized for most simulation conditions (bias around zero, stable confidence interval widths, nominal coverage rate at 95%).

Figure 3A shows that bias is fairly stable across simulation conditions. The condition  $T = 1$  clearly deviates from this trend with a negative bias. The bias for  $T = 2$  is below average, but within the range of fluctuations. Conditions where  $T > 3$  can be diagnosed as unbiased, but the average bias across iterations (as shown by a flat Loess line) only reaches stability at  $T = 20$  (**Loess might not be relevant, check and remove accordingly!**).

CIW is only clearly divergent in the simulation condition  $T = 1$ , see Figure 3B. Only conditions where  $T > 3$  are therefore considered sufficient. (**Possibly irrelevant Loess:** These conditions have similar CIWs, but across iterations stability is not established (i.e., the Loess line is never completely flat within the 100 simulation conditions).)

The empirical coverage rate across repetitions seems more or less stable for conditions where  $T > 1$ , see Figure 3C. We see some over-coverage at  $T = 2$ , but that is better than under-estimating the variance of  $\bar{Q}$ . On average, the coverage rate is somewhat higher than the expected nominal coverage of 95% (**Neyman 1934**) (**not anymore??**). (**Loess:** Similar to the CIWs, there is some trending across iterations (i.e., the Loess line is never flat)).

From the simulation diagnostics, we observe that unbiased estimates with nominal coverage rates were obtained in conditions where  $T > 3$ . This suggests that as little as four iterations may be sufficient for the MI algorithm under the current circumstances.

In short, there is a discrepancy between what the convergence diagnostics and what the performance measures indicate. While  $T = 4$  seems sufficient with respect to simulation quantities, the condition where  $T = 4$  resulted in the 'one-but-worst' values of both convergence diagnostics. Complete algorithmic convergence as indicated by  $\hat{R}$  and autocorrelation is not reached in conditions where  $T < 20$ .

## Discussion

This note shows that convergence diagnostics  $\hat{R}$  and  $AC$  may diagnose convergence of multiple imputation algorithms, but their performance differs from conventional applications to iterative algorithmic procedures. **(nope! it shows that MICE can lead to correct outcomes when they have not converged according to two common conv diags. This may be due to the measures (e.g., assumption of overdisp) or due to the Qs (lm reg coeff, not higher dimensional/more complex RQs). Add what %miss has to do with it.)**

$\hat{R}$  and autocorrelation indicate that algorithmic convergence may only be reached after twenty or even forty iterations, while unbiased, confidence valid estimates may be obtained with as little as four iterations. These results are in agreement with the simulation hypothesis:  $\hat{R}$  over-estimates the severity of non-convergence when applied to MI procedures. %This may be due to the quantity of scientific interest chosen. More 'complicated'  $Q$ s (e.g., higher order effects or variance components) might show bias, under- or over-coverage at higher  $T$ .

According to this simulation study, the recently proposed threshold of  $\hat{R} < 1.01$  may be too stringent for MI algorithms. **(This is only one of the goals: to give applied researchers a diagnostic to indicate that they should keep iterating. The other is the default in mice and other software packages, and yet another is ... i forgot)** Under the relatively easy missing data problem of the current study, the threshold was not reached. The other extreme of the  $\hat{R}$ -thresholds, the conventionally acceptable  $\hat{R} < 1.2$ , may be too lenient for MI procedures. Applying this threshold to the current data, lead to falsely diagnosing convergence at  $T = 3$  **(because it goes up after, not because it is not converged enough)**. It appears that the widely used threshold of  $\hat{R} < 1.1$  suits MI algorithms the best. We might, however, also formulate a new threshold, specifically for the evaluation of MI algorithms. The current study suggests that  $\hat{R} < 1.05$  may be implemented, since that is the level at which the  $\hat{R}$  stabilize (around  $T = 20$ ) **(Not necessary for this Q, but maybe for more complicated Qs)**.

The negative  $AC$ -values obtained in this study show no threat of non-stationarity. However, initial dip in  $AC$ -values may have implications for the default number of iterations in mice (`maxit = 5`). Terminating the algorithm at  $T = 5$  may not be the most appropriate, since this lead to the worst convergence **(nope, only for Rh, not AC)**, as indicated by  $\hat{R}$  and  $AC$ . Under the current specifications,  $T > 20$  would be more appropriate.

The observed dip in  $AC$  implies that default `maxit` value of five iterations is the worst possible number of iterations. Moreover, the results of this study imply that assessing the stationarity component of convergence with  $AC$  might be redundant.

Further research is needed to investigate their performance under clear violation of convergence, e.g. dependency between predictors (predictors with very high correlations). Until then, we have only shown that the convergence diagnostics can diagnose non-convergence of MI algorithms that trend towards a converged state. Also for future research, look at developing a convergence diagnostic for substantive models, and implement a Wald test for  $AC = 0$ .

## References

- 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.
- Geweke J (1992) Evaluating the accuracy of sampling-based approaches to the calculations of posterior moments. *Bayesian statistics* 4: 641–649.
- 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.
- Neyman J (1934) On the Two Different Aspects of the Representative Method: The Method of Stratified Sampling and the Method of Purposive Selection. *Journal of the Royal Statistical Society* 97(4): 558–625. DOI:10.2307/2342192.
- R Core Team (2020) *R: A language and environment for statistical computing*. Vienna, Austria. URL <https://www.R-project.org/>. Tex.organization: R Foundation for Statistical Computing.
- 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. URL <https://www.jstor.org/stable/2335739>. Publisher: [Oxford University Press, Biometrika Trust].
- 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. URL <https://www.jstatsoft.org/index.php/jss/article/view/v045i03>.

- 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>. ArXiv: 1903.08008.
- Vink G and van Buuren S (2014) Pooling multiple imputations when the sample happens to be the population. *arXiv:1409.8542 [math, stat]* URL <http://arxiv.org/abs/1409.8542>. ArXiv: 1409.8542 version: 1.