

Imputation Diagnostics

Utrecht University Winter School: Missing Data in R



**Utrecht
University**

Kyle M. Lang

Department of Methodology & Statistics
Utrecht University

2022-02-03

Outline

Imputation Model Convergence

Plausibility of Imputed Values



Example Data

The example data are synthesized from questionnaire data collected by Lang, Salter, and Adams (2009).

- $N = 87$
- $P = 33$ Likert-type variables assessing:
 - Perceptions and definitions of racism
 - Political affiliation
 - Support for affirmative action policies
 - Belief in meritocratic ideals

The data synthesis involved:

1. Resampling the original data to produce a new sample of 250 cases
2. Adding Gaussian noise
3. Imposing 25% MAR missing
 - MAR Predictors = Political Affiliation, Definition of Racism



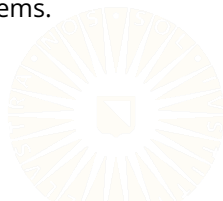
Imputation Diagnostics

After we run an MI routine, we need to make sure that the procedure has performed as expected.

Problems can arise to two different places:

1. The imputation model may fail to converge.
2. The imputed values may be invalid.

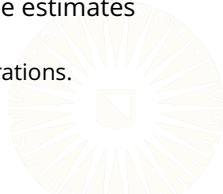
We need to examine our results to check for these problems.



Imputation Model Convergence

The imputation model is usually estimated through some form of Bayesian simulation.

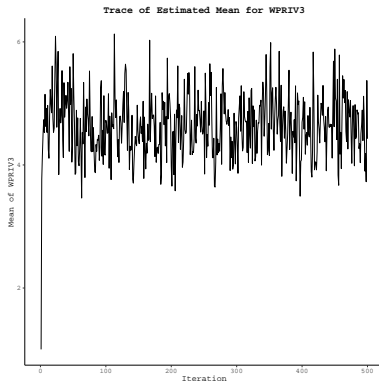
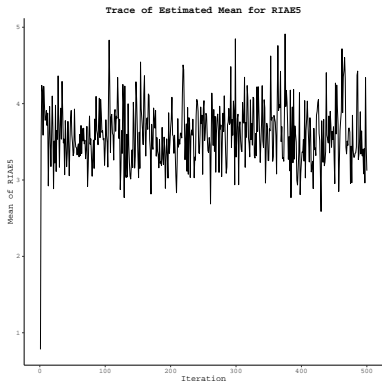
- Gibbs sampled parameters form a *Markov Chain*.
 - Each draw is dependent on only its immediate predecessor in the chain.
 - $\theta^{(t)} | \theta^{(t-1)} \perp \theta^{(t-j)} \forall j > 1$
- Early elements of a Markov chain are similar to the starting values.
 - Samples are poor approximations of the true posterior.
- We must let the sampler iterate for a while to allow the estimates time to separate from their starting values.
 - We call these initial iterations “burn-in” or “warm-up” iterations.



Traceplots

Once converged, each sampled parameter should “bounce” around some equilibrium point.

- The draws will never converge to a single point.
- Deterministic convergence would defeat the purpose of simulation.



Potential Scale Reduction Factor

Suppose we have two Markov chains for the same parameter.

- If these chains have converged, the average distance between any two points on separate chains should be the same as the average distance between two points on the same chain.
- The *between-chain* variance should, on average, equal the *within-chain* variance.

The Gelman and Rubin (1992) *Potential Scale Reduction Factor*, \hat{R} , quantifies this concept:

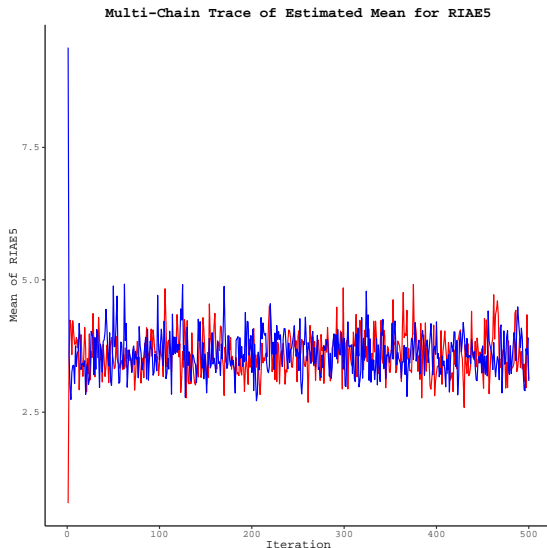
$$\hat{R} = \frac{\hat{\sigma}_{between}^2}{\hat{\sigma}_{within}^2}$$

\hat{R} will approach 1.0 at convergence.

- $\hat{R} < 1.1$ or 1.2 suggests acceptable convergence.



Example: Potential Scale Reduction Factor



Example: Potential Scale Reduction Factor

```
## Pool chains:
iterMat  <- cbind(chain1, chain2)
burntMat <- tail(iterMat, nSams - 100)

## R-Hat from all iterations:
wVar1 <- apply(iterMat, 2, var) %>% mean()
bVar1 <- apply(iterMat, 1, var) %>% mean()
rHat1 <- bVar1 / wVar1
rHat1

[1] 1.242591

## Burnt-In R-Hat:
wVar2 <- apply(burntMat, 2, var) %>% mean()
bVar2 <- apply(burntMat, 1, var) %>% mean()
rHat2 <- bVar2 / wVar2
rHat2

[1] 1.067016
```

More Imputation Model Convergence

A convergent imputation model will produce imputed values that fluctuate around an equilibrium point.

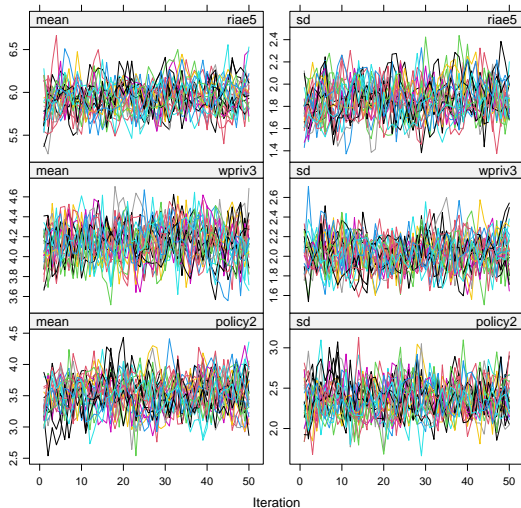
- Imputation model convergence can be assessed indirectly by looking at plots of the item-level sufficient statistics for each imputation.

This approach is automated for **mice** via `plot.mice()`.

```
## Impute missing values:
miceOut <- mice(data = dat3,
               m      = 25,
               maxit   = 50,
               method  = "norm",
               seed    = 235711)

## Create diagnostic traceplots:
plot(miceOut1, c("riae5", "wpriv3", "policy2"))
```

More Imputation Model Convergence



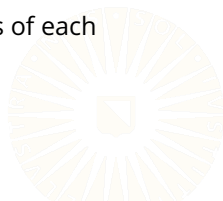
Imputed Value Plausibility

We need to ensure that the imputations are sensible.

- Imputed values shouldn't be *too* dissimilar from their observed counterparts.
 - What constitutes *too* much dissimilarity is subjective and problem-specific.

We can assess dissimilarity graphically or through summary statistics.

- Out-of-bounds values for the imputations are perfectly acceptable.
 - MI is *NOT* designed to maintain the range.
 - We don't want wildly extreme values, though.
- The means of the observed and imputed components of each variable shouldn't differ too much.
 - Again, how much is *too* much is subjective.



Numeric Imputation Checks

```
## Fill the missing values with imputations:
impList <- complete(miceOut, "all")

## Computes means:
rawMeans <- colMeans(missData, na.rm = TRUE)
impMeans <- do.call("rbind", impList) %>% colMeans()

## Compute standard deviations:
rawSds <- sapply(missData, sd, na.rm = TRUE)
impSds <- lapply(impList, function(x) sapply(x, sd)) %>%
  do.call(rbind, .) %>%
  colMeans()

## Compute ranges:
rawRanges <- sapply(missData, range, na.rm = TRUE)
impRanges <- do.call("rbind", impList) %>% sapply(range)
```

Numeric Imputation Checks

Compare observed and imputation-based means:

```
vars <- grep("policy\\d", colnames(missData))  
  
round(rawMeans[vars], 3)  
  
policy1 policy3 policy4 policy5 policy6 policy2  
3.299 3.784 3.528 3.024 4.722 3.700  
  
round(impMeans[vars], 3)  
  
policy1 policy3 policy4 policy5 policy6 policy2  
3.343 4.010 3.223 3.413 4.435 3.666
```

Numeric Imputation Checks

Compare observed and imputation-based standard deviations:

```
round(rawSds[vars], 3)
```

```
policy1 policy3 policy4 policy5 policy6 policy2  
2.275   2.238   1.996   2.008   1.999   2.114
```

```
round(impSds[vars], 3)
```

```
policy1 policy3 policy4 policy5 policy6 policy2  
2.259   2.309   2.115   2.167   2.093   2.187
```

Numeric Imputation Checks

Compare observed and imputation-based ranges:

```
round(rawRanges[ , vars], 3)
```

| | policy1 | policy3 | policy4 | policy5 | policy6 | policy2 |
|------|---------|---------|---------|---------|---------|---------|
| [1,] | -1.383 | -1.342 | -0.826 | -1.106 | 0.126 | -1.116 |
| [2,] | 8.975 | 9.551 | 9.471 | 8.700 | 9.608 | 10.237 |

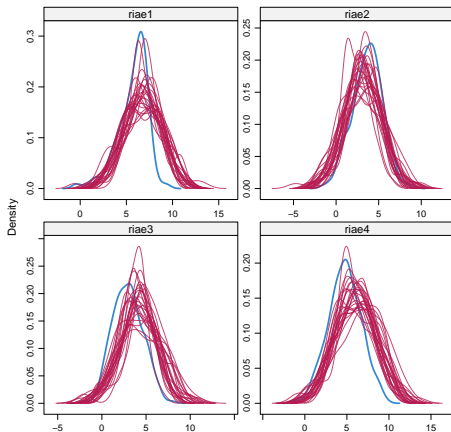
```
round(impRanges[ , vars], 3)
```

| | policy1 | policy3 | policy4 | policy5 | policy6 | policy2 |
|------|---------|---------|---------|---------|---------|---------|
| [1,] | -3.290 | -4.887 | -4.717 | -2.696 | -3.450 | -4.241 |
| [2,] | 10.751 | 13.251 | 9.471 | 12.595 | 11.724 | 12.047 |

Graphical Imputation Checks

Overlaid density plots of imputed vs. observed values:

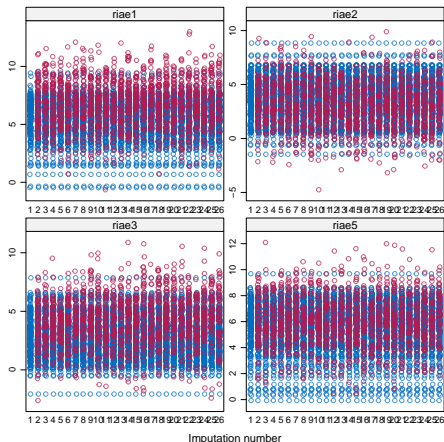
```
densityplot(miceOut, ~ riae1 + riae2 + riae3 + riae4, layout = c(2, 2))
```



Graphical Imputation Checks

Stripplots of imputed vs. observed values:

```
stripplot(miceOut, riae1 + riae2 + riae3 + riae5 ~ .imp, layout = c(2, 2))
```



References

- Gelman, A., & Rubin, D. B. (1992). Inference from iterative simulation using multiple sequences. *Statistical science*, 7(4), 457–472.
- Lang, K. M., Salter, P. S., & Adams, G. (2009, April). What drives the relationship between conservatism and racism? A mediation analysis. In *Proceedings of the annual meeting of the Southwestern Psychological Association*.

