

Lecture 7: Imputation Diagnostics

EPSY 6349: Modern Missing Data Analysis

Kyle M. Lang

Institute for Measurement, Methodology, Analysis & Policy
Texas Tech University
Lubbock, TX

October 13, 2015



TEXAS TECH

UNIVERSITY.

- Discuss imputation diagnostics
 - Assessing imputation model convergence
 - Checking the imputations' plausibility
- Look at graphical and numeric options for both

Today, our example data will be from a real study.

These data were analyzed by Lang, Salter, and Adams (2009).

- $N = 87$
- $V = 33$
 - Variables assessing:
 - Perceptions of and Definitions of Racism
 - Political Affiliation
 - Support for Affirmative Action Policies
 - Belief in meritocratic ideals
- Almost no missing data
 - I've artificially imposed 30% MAR missing data on all variables (except political affiliation) using political affiliation as the MAR predictor.

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:

- ❶ The imputation model may fail to converge.
- ❷ The imputed values may not be plausible.

We need to examine our results to see if either of these problems are present.

Imputation Model Convergence

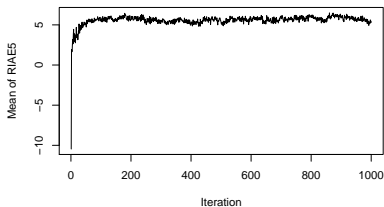
The imputation model is almost always estimated through some form of Gibbs sampling.

- Gibbs sampled parameters form a *Markov Chain*.
 - Each draw is dependent on only its immediate predecessor in the chain.
 - $\theta^{(t)} | \theta^{(t-1)} \perp \theta^{(t-j)} \quad \forall j > 1$
- Early elements of a Markov chain are similar to the starting values
 - Early Gibbs 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

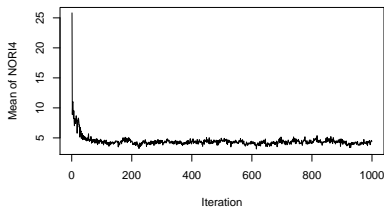
Once converged, each sampled imputation model parameter should “bounce” around an equilibrium point.

- The draws will never converge onto a single point
- That would defeat the purpose of simulation-based inference

Trace of RIAE5's Estimated Mean



Trace of NORI4's Estimated Mean



Potential Scale Reduction Factor

Suppose we have two independent Markov chains of 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, asymptotically, equal the *within-chain* variance.

The *Potential Scale Reduction Factor* \hat{R} quantifies this concept:

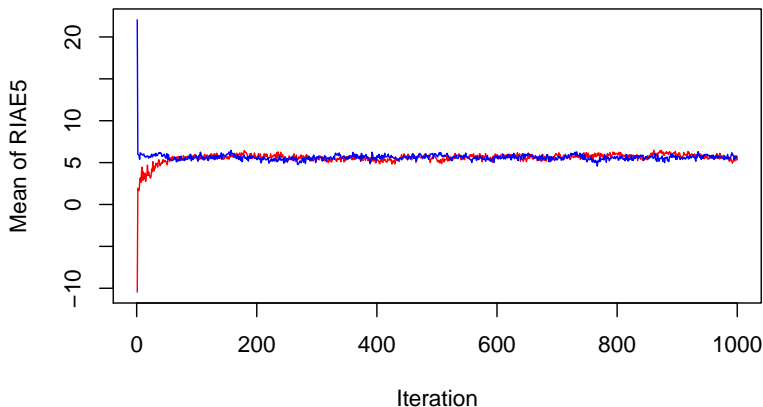
$$\hat{R} = \frac{\hat{\sigma}_{between}^2}{\hat{\sigma}_{within}^2}$$

The \hat{R} will approach 1.0 at convergence.

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

Example: Potential Scale Reduction Factor

Multi-Chain Trace of RIAE5's Estimated Mean



Example: Potential Scale Reduction Factor

```
## Full Chains:
iterMat ← cbind(chain1, chain2)
## Excluding Burn-In:
burntMat ← iterMat[201 : 1000, ]
## Full Chain R-Hat:
wVar1 ← mean(apply(iterMat, 2, var))
bVar1 ← mean(apply(iterMat, 1, var))
rHat1 ← bVar1 / wVar1
rHat1
```

```
[1] 1.682921
```

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

```
[1] 1.104803
```

More Imputation Model Convergence

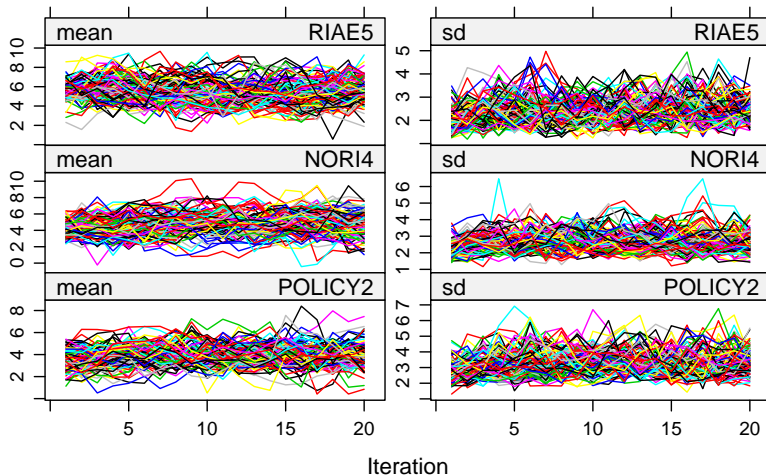
A convergent imputation model will produce imputed values that fluctuate around an equilibrium point.

- 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()`.

```
miceOut1 ← readRDS("examples/miceOut1.rds")  
plot(miceOut1, c("RIAE5", "NORI4", "POLICY2"))
```

More Imputation Model Convergence



Imputed Value Plausibility

We want to ensure that the values imputed for the missing data are sensible.

- Imputed values shouldn't be *too* dissimilar to 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

```
rawMeans ← colMeans(missData, na.rm = TRUE)
impMeans ← colMeans(do.call("rbind", impList))
rawSds ← apply(missData, 2, sd, na.rm = TRUE)
impSds ← apply(do.call("rbind", impList), 2, sd)
rawRanges ← apply(missData, 2, range, na.rm = TRUE)
impRanges ← apply(do.call("rbind", impList), 2, range)
```

Numeric Imputation Checks

```
round(rawMeans[1 : 5], 3)
```

```
RIAE2 RIAE3 RIAE7 RIAE8 RIAE9  
3.677 3.108 3.774 3.092 3.726
```

```
round(impMeans[1 : 5], 3)
```

```
RIAE2 RIAE3 RIAE7 RIAE8 RIAE9  
3.697 3.122 3.340 3.134 3.010
```

```
round(rawSds[1 : 5], 3)
```

```
RIAE2 RIAE3 RIAE7 RIAE8 RIAE9  
1.696 1.522 2.060 1.693 1.700
```

```
round(impSds[1 : 5], 3)
```

```
RIAE2 RIAE3 RIAE7 RIAE8 RIAE9  
2.281 1.924 2.994 2.175 2.327
```

Numeric Imputation Checks

```
round(rawRanges[ , 1 : 5], 3)
```

	RIAE2	RIAE3	RIAE7	RIAE8	RIAE9
[1,]	1	1	1	1	1
[2,]	7	6	7	7	7

```
round(impRanges[ , 1 : 5], 3)
```

	RIAE2	RIAE3	RIAE7	RIAE8	RIAE9
[1,]	-7.819	-8.486	-17.790	-9.473	-10.184
[2,]	21.724	15.941	15.869	16.551	11.201

We can also construct plots of the imputed vs. observed values.

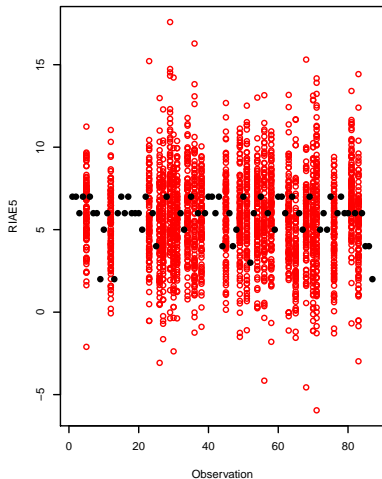
I tend to use two different flavors:

- 1 Scatterplots
- 2 Overlaid Density plots

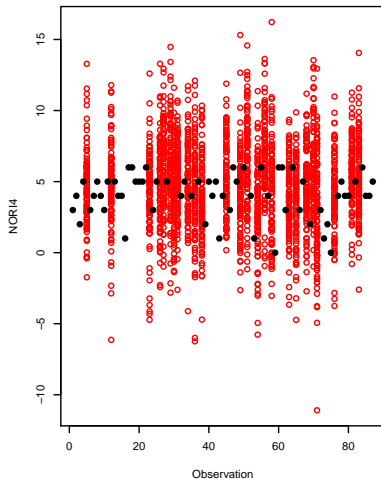
```
## Scatterplots of imputed vs. observed values:  
par(mfrow = c(1, 2), cex = 0.5)  
plotImps(impList = impList,  
         targetVar = c("RIAE5", "NORI4"))
```


Graphical Imputation Checks

Imputed (Red) vs. Observed (Black) Values
for RIAE5



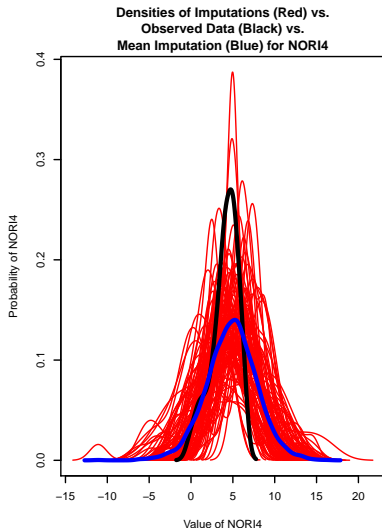
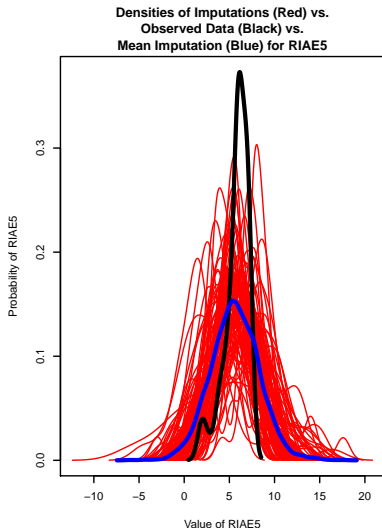
Imputed (Red) vs. Observed (Black) Values
for NORI4



Graphical Imputation Checks

```
## Overlaid density plots of imputed vs. observed values:  
par(mfrow = c(1, 2), cex = 0.5)  
plotImps(impList = impList,  
         targetVar = c("RIAE5", "NORI4"),  
         type = "density")
```

Graphical Imputation Checks



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*.