

Multivariate Multiple Imputation

Utrecht University Winter School: Missing Data in R



**Utrecht
University**

Kyle M. Lang

Department of Methodology & Statistics
Utrecht University

Outline

Flavors of MI

Fully Conditional Specification

Joint Modeling

Mixed Data Types



Joint Modeling vs. Fully Conditional Specification

When imputing with *Joint Modeling* (JM) approaches, the missing data are replaced by samples from the joint posterior predictive distribution.

- To impute X , Y , and Z , we draw:

$$X, Y, Z \sim P(X, Y, Z | \theta)$$

With *Fully Conditional Specification* (FCS), the missing data are replaced with samples from the conditional posterior predictive distribution of each incomplete variable.

- To impute X , Y , and Z , we draw:

$$X \sim P(X | Y, Z, \theta_X)$$

$$Y \sim P(Y | X, Z, \theta_Y)$$

$$Z \sim P(Z | Y, X, \theta_Z)$$



Joint Modeling: Strengths

When correctly implemented, JM approaches are guaranteed to produce *Bayesianly proper* imputations.

- A sufficient condition for *properness* is that the imputations are randomly sampled from the correctly specified joint posterior predictive distribution of the missing data.
 - This is the defining characteristic of JM methods.

When using the correct distribution, imputations produced by JM methods will be the best possible imputations.

- Unbiased parameter estimates
- Well-calibrated sampling variability



Joint Modeling: Weaknesses

JM approaches don't scale well.

- The computational burden increases with the number of incomplete variables.

JM approaches are only applicable when the joint distribution of all incomplete variables follows a known form.

- Mixes of continuous and categorical variables are difficult to accommodate.

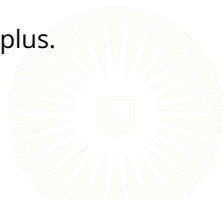


Software Implementations

In R, MI via JM is available from several packages.

- **Amelia** (Honaker, King, & Blackwell, 2011)
 - Bootstrapped EM algorithm
- **norm** (Schafer, 2013)
 - Classic data augmentation.
- **mice** (Van Buuren & Groothuis-Oudshoorn, 2011)
 - Data augmentation for block updating.

JM imputation is also available in SAS, Stata, SPSS, and Mplus.



Fully Conditional Specification: Strengths

FCS scales much better than JM.

- FCS only samples from a series of univariate distributions, not large joint distributions.

FCS approaches can create imputations for variables that don't have a sensible joint distribution.

- FCS can easily treat mixes of continuous and categorical variables.



Fully Conditional Specification: Weaknesses

FCS will usually be slower than JM.

- Each variable gets its own fully parameterized distribution, even if that granularity is unnecessary.

When the incomplete variables don't have a known joint distribution, FCS doesn't have theoretical support.

- There is, however, a large degree of empirical support for the tenability of the FCS approach.
- In practice, we usually choose FCS since real data rarely arise from a known joint distribution.



Software Implementations

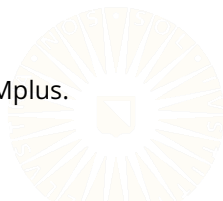
The **mice** package is the most popular R implementation of FCS. (Van Buuren & Groothuis-Oudshoorn, 2011).

- Mature implementation
- Well integrated into the larger R ecosystem
- Very active development

The **mi** package (Su, Yajima, Gelman, & Hill, 2011) offers another option.

- More focus on diagnostics
- Object oriented flavor
- Not very actively developed

FCS imputation is also available in SAS, SPSS, Stata, and Mplus.



FULLY CONDITIONAL SPECIFICATION



Procedure: Fully Conditional Specification

1. Fill the missing data with reasonable guesses.
2. For each incomplete variable, do a single iteration of univariate Bayesian MI (e.g., as seen in the last set of slides).
 - After each variable on the data set is so treated, we've completed one iteration.
3. Repeat Step 2 many times.
4. After the imputation model parameters stabilize, save M imputed data sets.



Example: Data Generation

First we'll simulate some synthetic data.

```
## Simulate some data:
simData <-
  simCovData(nObs = 1000, sigma = 0.25, nVars = 4)

head(simData, 10)
```

	x1	x2	x3	x4
1	0.06313632	-1.4057704	-0.01709217	1.47929405
2	-1.31592547	1.2970920	-0.83500777	-0.44528158
3	-0.30997023	0.9782580	0.02731853	0.35507390
4	0.06927787	0.1836032	0.68794409	0.08049987
5	-0.99354894	-0.3038956	0.80918329	-1.72143555
6	0.36828016	-0.9423245	1.05155348	-0.11078496
7	1.33333163	2.3089780	1.47203000	0.85877495
8	-0.02759718	0.1714383	0.18927909	0.28627771
9	2.37929433	2.5080935	2.18344726	1.56980951
10	0.31841502	0.7886025	0.73658136	0.39445970

Example: Data Generation

Next, we impose some missing values on the simulated data.

```
targets <- paste0("x", 1:3)
missData <- imposeMissData(data    = simData,
                           targets = targets,
                           preds   = "x4",
                           pm      = 0.3,
                           types  = c("low", "center", "high")
                           )
```

Example: Data Visualization

Check the results.

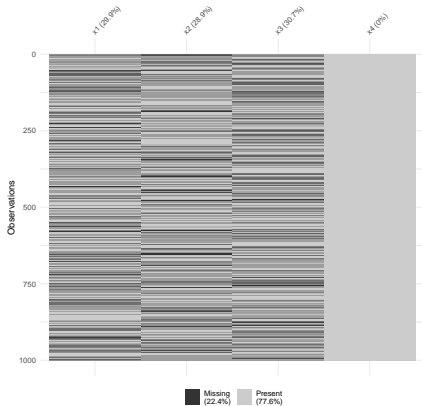
```
head(missData, 10)
```

	x1	x2	x3	x4
1	0.06313632	NA	-0.01709217	1.47929405
2	NA	1.2970920	NA	-0.44528158
3	-0.30997023	NA	0.02731853	0.35507390
4	NA	NA	0.68794409	0.08049987
5	NA	NA	0.80918329	-1.72143555
6	0.36828016	NA	NA	-0.11078496
7	1.33333163	2.3089780	NA	0.85877495
8	NA	NA	0.18927909	0.28627771
9	2.37929433	2.5080935	NA	1.56980951
10	0.31841502	0.7886025	0.73658136	0.39445970

Example: Data Visualization

Use the `naniar::vis_missing()` function to visualize the pattern of missing values.

```
vis_miss(missData)
```



JOINT MODELING

Aside: Definition of Regression Parameters

So far, we've been using the least-squares estimates of α , β , and σ^2 to parameterize our posterior distributions.

- We can also define the parameters in terms of sufficient statistics.

Given μ and Σ , we can define all of our regression moments as:

$$\begin{aligned}\beta &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= \text{Cov}(\mathbf{X})^{-1} \text{Cov}(\mathbf{X}, \mathbf{Y}) \\ \alpha &= \mu_Y - \beta^T \mu_X \\ \Sigma_\epsilon &= \Sigma_Y - \beta^T \Sigma_X \beta\end{aligned}$$

These definitions are crucial for JM approaches.

- Within the subset of data define by a given response pattern, the outcome variables will be entirely missing.

Multivariate Bayesian Regression

Previously, we saw examples of univariate Bayesian regression which used the following model:

$$\sigma^2 \sim \text{Inv-}\chi^2(N - P, \text{MSE})$$

$$\beta \sim \text{MVN}(\hat{\beta}_{ls}, \sigma^2(\mathbf{X}^T\mathbf{X})^{-1})$$

We can directly extend the above to the multivariate case:

$$\Sigma^{(i)} \sim \text{Inv-W}(N - 1, (N - 1)\Sigma^{(i-1)})$$

$$\mu^{(i)} \sim \text{MVN}(\mu^{(i-1)}, N^{-1}\Sigma^{(i)})$$

We get α , β , and Σ_ϵ via the calculations on the preceding slide

Procedure: Joint Modeling

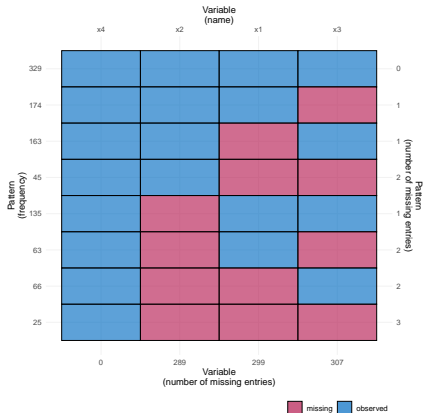
In JM imputation, we estimate the imputation model via the Tanner and Wong (1987) *data augmentation* algorithm.

1. Partition the incomplete data by response pattern.
 - Produce S subsets wherein each row shares the same response pattern.
2. Provide initial guesses for μ and Σ .
3. Within each subset, use the current guesses of μ and Σ to generate imputations via multivariate Bayesian regression.
4. Use the filled-in data matrix to updated the sufficient statistics.
5. Repeat Steps 3 and 4 many times.
6. After the imputation model parameters have stabilized, save M imputed data sets produced in Step 3.

Example: Data Visualization

Use `ggmice::plot_pattern()` to visualize the response patterns.

```
library(ggmice)
plot_pattern(missData)
```



MIXED DATA TYPES

FCS for Mixed Data Types

FCS imputation can easily accommodate incomplete data that contain both continuous and categorical/non-normal variables.

- Replace the normal-theory elementary imputation model described above with an appropriate model for the distribution of each incomplete variable
 - Logistic regression (various flavors)
 - Donor-based methods
 - Tree-based methods

The FCS framework can essentially accommodate any data for which you can define an appropriate supervised model.

- Many useful methods are already implemented in the **mice** package.

JM for Mixed Data Types

When applying JM to incomplete data with mixed variable types, we have to general options.

1. Impute under the multivariate normal model and round, coarsen, or truncate the continuous imputations to “match” the original data.
 - This was the old-school recommendation from the days before FCS (e.g., Allison, 2002; Schafer, 1997).
 - The **Amelia** package implements this approach.
 - This approach tends to perform poorly in methodological evaluations (e.g., Lang & Wu, 2017; Wu, Jia, & Enders, 2015).
2. Impute under an appropriate joint model for the data.
 - This approach is only available when a suitable joint model exists.
 - The **mix** package (Schafer, 2017) implements this approach for the general location model (Little & Schluchter, 1985).
 - This approach also doesn't do very well, in practice (Lang & Wu, 2017).

References

- Allison, P. D. (2002). *Missing data*. Thousand Oaks, CA: Sage Publications.
- Honaker, J., King, G., & Blackwell, M. (2011). Amelia II: A program for missing data. *Journal of Statistical Software*, 45(7), 1–47. Retrieved from <https://www.jstatsoft.org/v45/i07/>
- Lang, K. M., & Wu, W. (2017). A comparison of methods for creating multiple imputations of nominal variables. *Multivariate Behavioral Research*, 52(3), 290-304. doi: 10.1080/00273171.2017.1289360
- Little, R. J. A., & Schluchter, M. D. (1985). Maximum likelihood estimation for mixed continuous and categorical variables with missing values. *Biometrika*, 72, 497–512.
- Schafer, J. L. (1997). *Analysis of incomplete multivariate data* (Vol. 72). Boca Raton, FL: Chapman & Hall/CRC.
- Schafer, J. L. (2013). norm: Analysis of multivariate normal datasets with missing values [Computer software manual]. Retrieved from <https://CRAN.R-project.org/package=norm> (R package version 1.0-9.5, Ported to R by Alvaro A. Novo)

References

- Schafer, J. L. (2017). *mix: Estimation/multiple imputation for mixed categorical and continuous data* [Computer software manual]. Retrieved from <https://CRAN.R-project.org/package=mix> (R package version 1.0-10)
- Su, Y.-S., Yajima, M., Gelman, A. E., & Hill, J. (2011). Multiple imputation with diagnostics (mi) in r: opening windows into the black box. *Journal of Statistical Software*, 45(2), 1–31.
- Tanner, M. A., & Wong, W. H. (1987). The calculation of posterior distributions by data augmentation. *Journal of the American Statistical Association*, 82(398), 528–540.
- Van Buuren, S., & Groothuis-Oudshoorn, K. (2011). mice: Multivariate imputation by chained equations in R. *Journal of Statistical Software*, 45(3), 1–67.

References

Wu, W., Jia, F., & Enders, C. (2015). A comparison of imputation strategies for ordinal missing data on likert scale variables. *Multivariate Behavioral Research*, 50(5), 484-503. doi: 10.1080/00273171.2015.1022644