

Fully conditional specification in multivariate imputation

S. VAN BUUREN*†‡, J. P. L. BRAND§, C. G. M. GROOTHUIS-OUDSHOORN¶ and D. B. RUBIN $\|$

†Department of Statistics, TNO Quality of Life, Wassenaarseweg 56, P.O. Box 2215, 2301 CE Leiden, The Netherlands ‡University of Utrecht, The Netherlands \$Pennington Biomedical Research Center, Baton Rouge, LA 70808, USA ¶Roessingh Research and Development, Enschede, The Netherlands ||Harvard University, Cambridge, MA 02138, USA

(Received 15 January 2004; in final form 22 March 2005)

The use of the Gibbs sampler with fully conditionally specified models, where the distribution of each variable given the other variables is the starting point, has become a popular method to create imputations in incomplete multivariate data. The theoretical weakness of this approach is that the specified conditional densities can be incompatible, and therefore the stationary distribution to which the Gibbs sampler attempts to converge may not exist. This study investigates practical consequences of this problem by means of simulation. Missing data are created under four different missing data mechanisms. Attention is given to the statistical behavior under compatible and incompatible models. The results indicate that multiple imputation produces essentially unbiased estimates with appropriate coverage in the simple cases investigated, even for the incompatible models. Of particular interest is that these results were produced using only five Gibbs iterations starting from a simple draw from observed marginal distributions. It thus appears that, despite the theoretical weaknesses, the actual performance of conditional model specification for multivariate imputation can be quite good, and therefore deserves further study.

Keywords: Multivariate missing data; Multiple imputation; Distributional compatibility; Gibbs sampling; Simulation; Proper imputation

1. Introduction

Missing data often plague the statistical analysis of multivariate data. When confronted with incomplete data, the analyst can choose a variety of strategies: *ad hoc* methods (*e.g.*, analysis of the complete cases only, available case methods, use of some indicator variables with means filled in), likelihood-based approaches that allow for missing data (*e.g.*, EM algorithm, structural equations, or mixed models), weighting, or imputation-based methods. The relative merits of these approaches have been discussed elsewhere [1, 2]. Multiple imputation [3–5]

^{*}Corresponding author. Email: s.vanbuuren@pg.tno.nl

is a general and statistically valid method for dealing with missing data. This paper studies a particular method for creating imputations (single or multiple) in multivariate data.

Let y denote an $n \times k$ matrix with data from n individuals on k variables. Let Y_j be the jth variable, and y_j the jth column of $y(j=1,\ldots,k)$. We define y_j^{obs} as the observed part of y_j and y_j^{mis} as the missing part of y_j . Let $y^{\text{obs}} = (y_1^{\text{obs}}, \ldots, y_k^{\text{obs}})$ and $y^{\text{mis}} = (y_1^{\text{mis}}, \ldots, y_k^{\text{mis}})$ stand for the collection of all observed and missing data, respectively. Imputation of y_j^{mis} will typically be based on the relation between the incomplete variable Y_j and the k-1 predictors $Y_{-j} = (Y_1, \ldots, Y_{j-1}, Y_{j+1}, \ldots, Y_k)$, where the nature of the relation is primarily estimated from the units contributing to y_j^{obs} . For notational convenience, we suppress the notation for all variables that are fully observed, and so all distributions are implicitly conditional on the fully observed variables. Thus, each of the k columns in y has some missing values.

A number of practical problems can occur in general when k > 1:

- the predictors Y_{-i} themselves contain missing values;
- 'circular' dependence occurs, where y_j^{mis} depends on y_h^{mis} and y_h^{mis} depends on $y_j^{\text{mis}}(h \neq j)$, because in general Y_j and Y_h are correlated, even given other variables;
- especially with large k and small n, collinearity or empty cells occur;
- rows or columns can be ordered, e.g., as with longitudinal data;
- variables are of different types (e.g., binary, unordered, ordered, continuous), thereby making the application of theoretically convenient models, such as the multivariate normal, theoretically inappropriate;
- the relation between Y_j and predictors Y_{-j} is complex, e.g., nonlinear, or subject to censoring processes;
- imputation can create impossible combinations such as pregnant fathers.

This list is by no means exhaustive, and other complexities may appear for particular data. Two general strategies for imputing multivariate data have surfaced during the last decade: joint modeling and fully conditional specification (FCS).

The first common strategy, joint modeling, begins by specifying a parametric multivariate density $P(Y|\theta)$ for the data Y given the model parameters θ . Given this specification and appropriate prior distributions for θ , one can use the Bayesian framework to generate imputations as draws from the posterior predictive distribution $P(y^{\text{mis}}|y^{\text{obs}})$, usually under the assumption of an ignorable missing data mechanism. Using this approach, Schafer [2] described sophisticated methods for creating multivariate imputations under the multivariate normal, the log-linear, and the general location model. These methods are available as tools in S-Plus 6.2 and SAS V8.2.

The other common approach FCS does not start from an explicit multivariate density, but instead implicitly defines $P(Y|\theta)$ by specifying a separate conditional density $P(Y_j|Y_{-j},\theta_j)$ for each Y_{-j} . This density is used to impute y_j^{mis} given y_{-j} , for example, by linear or logistic regression applied to the cases in y_j^{obs} , where y_{-j} refers to the columns of y excluding y_j . Imputation under FCS is done by iterating over all conditionally specified imputation models, each iteration consisting of one cycle through all Y_i .

FCS has several important practical advantages over joint modeling. First, FCS allows for the creation of flexible multivariate models because it splits a k-dimensional problem into k one-dimensional problems. One can easily specify models that are outside any standard multivariate density $P(Y|\theta)$. Secondly, FCS may help to preserve investments in specialized imputation methods that are difficult to formulate as a part of a multivariate density $P(Y|\theta)$. For example, it is easy to incorporate imputation methods that preserve unique features in the data, e.g., bounds, skip patterns, interactions, bracketed responses, and so on. Also, it is relatively straightforward to maintain constraints between different variables. Such

constraints might be needed to avoid logical inconsistencies in the imputed data. Gelman and Raghunathan [6] observed that in such situations 'separate regressions often make more sense than joint models'. Thirdly, generalization to models under nonignorable missing data mechanisms might be easier. Finally, the idea of specifying a separate imputation model for each variable is easy to communicate to users.

However, FCS is not without drawbacks. First, each conditional density has to be specified separately, so substantial modeling effort can be needed for data sets with many variables. Second, FCS is often computationally more intensive than joint modeling. Typical computational shortcuts (*e.g.*, using the sweep operator) [1] may not apply. Last, and very importantly, relatively little is known about the quality of the resulting imputations because the implied joint distributions may not exist theoretically and that convergence criteria are ambiguous.

Variations of the FCS idea have appeared before. Buck [7] computed estimates for all missing entries by multiple regression, where the regression coefficients are computed using the complete cases, i.e., all individuals with fully complete data. The observed data for the individual constitute the independent variables in the equation predicting the missing variables for that individual. Gleason and Staelin [8] extended Buck's method to include multivariate regression and noted that their ad hoc method could also be derived more formally from the multivariate normal distribution. These authors also brought up the possibility of an iterated version of Buck's method, suggesting that missing entries from one iteration could be used to form an improved estimate of the correlation matrix for use in a subsequent iteration. Variations of iterated regression imputation have been studied later by Finkbeiner [9], Raymond and Roberts [10], Jinn and Sedransk [11], and Gold and Bentler [12]. Systems for creating multiple imputations by iterated regressions have been developed including FRITZ [13], IVEWARE [14], HERMES missing data engine [15], and MICE [16, 17]. Royston [18] created a version of MICE in Stata. Rubin [19] pioneers a technique that attempts to take the best of both worlds by combining an FCS model for some missing values with joint modeling on other missing data. Other applications of FCS can be found in Kennickell [20], Van Buuren et al. [21], Raghunathan and Siscovick [22], Oudshoorn et al. [23], Heeringa et al. [24], Gelman and Raghunathan [6], and Faris et al. [25].

There appears to be yet no really satisfactory theory, but in many examples, FCS seems to work well, is of great importance in practice, and is easily applied. Some simulation work is available [26–28], but this is relatively limited in scope and complexity. This paper presents a more extensive simulation-based evaluation of FCS.

2. Imputation by FCS

2.1 Definitions and introduction

Suppose $Y=(Y_1,Y_2,\ldots,Y_k)$ is a vector of k random variables with k-variate distribution $P(Y|\theta)$. We assume that the joint multivariate distribution of Y is completely specified by θ , a vector of unknown parameters. For example, if Y is multivariate normally distributed, $\theta=(\mu,\Sigma)$, with μ a k-dimensional mean vector and Σ a $k\times k$ covariance matrix. Let the matrix $y=(y_1,\ldots,y_n)$ with $y_i=(y_{i1},y_{i2},\ldots,y_{ik}), i=1,\ldots,n$ be an i.i.d. sample of the vector Y. The matrix y is partially observed, in the sense that each column in y has missing data. The standard procedure [3] for creating multiple imputations y^* of y^{mis} is as follows:

- 1. calculate the posterior distribution $p(\theta|y^{\text{obs}})$ of θ based on the observed data y^{obs} ;
- 2. draw a value θ^* from $p(\theta|y^{\text{obs}})$;

3. draw a value y^* from $p(y^{\text{mis}}|y^{\text{obs}}, \theta = \theta^*)$, the conditional posterior distribution of y^{mis} given $\theta = \theta^*$.

Repeat steps 2 and 3 for more imputations, e.g., 5–10. Appendix A gives algorithms for the cases where Y is univariate normally distributed, dichotomous, or polytomous.

For multivariate Y, the central problem is how to get the multivariate distribution of θ , either explicitly or implicitly. FCS proposes to obtain a posterior distribution of θ by sampling iteratively from conditional distributions of the form

$$P(Y_1|Y_{-1}, \theta_1),$$

$$\vdots$$

$$P(Y_k|Y_{-k}, \theta_k).$$
(1)

The parameters $\theta_1, \ldots, \theta_k$ are treated as specific to the respective conditional densities and are not necessarily the product of some factorization of the 'true' joint distribution $P(Y|\theta)$. More precisely, the *t*th iteration of the method consists of the following successive draws of the Gibbs sampler:

$$\theta_{1}^{*(t)} \sim P(\theta_{1}|y_{1}^{\text{obs}}, y_{2}^{(t-1)}, \dots, y_{k}^{(t-1)})$$

$$y_{1}^{*(t)} \sim P(y_{1}^{\text{mis}}|y_{1}^{\text{obs}}, y_{2}^{(t-1)}, \dots, y_{k}^{(t-1)}, \theta_{1}^{*(t)})$$

$$\vdots$$

$$\theta_{k}^{*(t)} \sim P(\theta_{k}|y_{k}^{\text{obs}}, y_{1}^{(t)}, y_{2}^{(t)}, \dots, y_{k-1}^{(t)})$$

$$y_{k}^{(t)} \sim P(y_{k}^{\text{mis}}|y_{k}^{\text{obs}}, y_{1}^{(t)}, \dots, y_{k-1}^{(t)}, \theta_{k}^{*(t)})$$

$$(2)$$

Observe that no information about y_j^{mis} is used to draw $\theta_j^{*(t)}$, which differs from the Markov Chain Monte Carlo approaches to joint modeling. Our method is just a concatenation of univariate procedures applied to the cases with complete y_j , and deviates from MCMC theory at this point. The iterations of equation (2) are executed m times in parallel to generate m multiple imputations. The number of iterations is fixed to a small number, say 5 or 10. This procedure implicitly assumes that the joint distribution is specified by equation (1) and that the Gibbs sampler in equation (2) provides draws from it. With k incomplete variables, the vector parameters $\theta_1, \ldots, \theta_k$ will generally depend on each other, and so the sampler can be overparametrized. For example, the space spanned by $\theta_1, \ldots, \theta_k$ generally has more dimensions than appropriate. If this occurs, the implicit joint distribution does not exist. This issue is known as the problem of compatibility of conditionals [29].

2.2 Compatibility

Bhattacharryya [30] observed that the combination of two conditional normal densities with linear regressions and constant variance defines a joint bivariate normal density. Two conditional densities are compatible if a joint distribution exists that has the given densities as its conditional density. In general, two conditional densities f(x|y) and g(y|x) are compatible if and only if (apart from a technical condition on the support of the densities) their density ratio f(x|y)/g(y|x) factorizes into u(x)v(y) for some integrable functions u and v [31]. So, the joint distribution either exists and is unique, or does not exist. If f(x|y) and g(y|x) are known functions, we can calculate f(x|y)/g(y|x) on a grid of x and y values and infer the compatibility if the matrix of f(x|y)/g(y|x) values is of rank 1.

In actual data analysis, compatibility is not an all-or-nothing phenomenon. For example, even simple rounding errors can destroy exact compatibility. Measures have been proposed to measure the amount of compatibility [32]. The effects of near compatibility and clear incompatibility on the quality of statistical inference are yet unknown, except in special cases. Section 6 therefore addresses the robustness of FCS under clearly incompatible models.

3. Evaluation of univariate imputation

This section provides simulation results for both compatible (univariate and multivariate) and incompatible FCS imputation methods. Gibbs sampling is not actually needed in univariate models, but forms an important ingredient for the multivariate case because of the use of univariate distributions in equation (1).

3.1 Setup: data and simulations

We study the quality of univariate linear and logistic imputation methods using a data set of Irish wind speeds [33], which contains the average daily wind speed measured at 12 meteorological stations in Ireland during the years 1961–1978 (6574 time points). The correlations among these stations are high, ranging from 0.59 to 0.84, thus enabling the use of MAR [34] missing data mechanisms that generate large differences between the complete and incomplete records. A random sample of n = 400 was taken. No attempts were made here to model the temporal variation between the measures. To investigate the linear imputation method, the following five locations from the Irish wind speed data were selected: $y_1 = \text{Rosslare}$, $x_1 = \text{Roche's Pt}$, $x_2 = \text{Shannon}$, $x_3 = \text{Dublin}$, and $x_4 = \text{Clones}$. The original data of y_1 were replaced by new data that were generated according to the conditional probability of y_1 given x_1, \ldots, x_4 under a linear model, which was done to avoid any issues of inaccuracy of model fit. Missing data in y_1 were subsequently created, where the response probability possibly depended on x_1, \ldots, x_4 using methods that are described subsequently.

For the logistic method, we selected y_1 = Valentia, x_1 = Roche's Pt, x_2 = Rosslare, x_3 = Shannon, and x_4 = Dublin, dichotomized y_1 in equally sized groups, and replaced the original data in y_1 by data generated according to a logistic regression model with linear predictors x_1, \ldots, x_4 . Missing entries were then created in these substitutes.

We used data from Hosmer and Lemeshow [35, p. 265] to study the polytomous model. This data set contains six responses from a survey of 412 women on knowledge, attitude, and behavior towards mammography. The target variable y_1 was mammographic experience (ME) with three response categories (0 = never, 1 = during past year, 2 = over a year ago). As before, the original values of y_1 were first replaced by data conforming to the polytomous logistic model conditional on the other data, and subsequently made missing.

Simulations were done using 1000 replications. In every replication, \sim 50% missingness in y_1 was generated under four different MAR missing data mechanisms: MCAR (missing completely at random), MARRIGHT, MARTAIL, and MARMID. Mechanism MCAR deletes observations in a completely random fashion, MARRIGHT creates more missingness for larger values, MARTAIL deletes more cases from both tails, and MARMID introduces more nonresponse in the center of the distribution. The latter three mechanisms introduce biases in the statistical analysis based on complete cases. Appendix B describes the methodology for generating the missing entries in detail. Figure 1 graphs the resulting distributions of the complete and incomplete target variable for one replication.

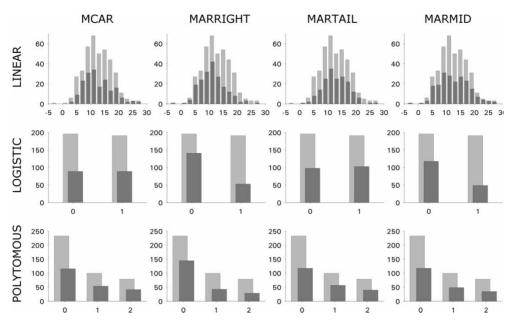


Figure 1. Distribution of the target variable before (light) and after (dark) missing data under four missing data mechanisms. MARRIGHT deletes more from the right tail, MARTAIL from both tails, and MARMID from the middle values. Data are the Irish wind speed sample (n = 400), [33] and ME data (n = 412) [35].

3.2 Results

Let Q be the quantity of interest, and let \hat{Q} be the associated complete-data estimate with variance U. For each replication $i=1,\ldots,1000$, multiple imputation with m=5 is applied to the incomplete data. This results in the pooled estimates $\bar{Q}_m^{(i)}$ of \hat{Q} , $\bar{U}_m^{(i)}$ of U, and $B_m^{(i)}$ as the variance between the m-complete data estimates [3, p. 76]. Validity conditions for proper imputation similar to those presented by Rubin [5] are: $\hat{E}[\bar{Q}_m] \approx \hat{Q}$, $\hat{E}[\bar{U}_m] \approx U$, and $\text{var}(\bar{Q}_m) \approx (1+m^{-1})\hat{E}[B_m]$, where $\hat{E}[]$ is the mean and var[] the variance over the replications. Computations were done in SAS-IML.

Table 1 reports the following statistics based on the simulation:

- \hat{Q} , the complete data statistic based on the underlying values of all cases (Pop);
- the fraction of information (Fmi) about Q missing due to nonresponse;
- $\hat{E}[\hat{Q}_{ac}]$, the average \hat{Q} computed from the available cases;
- the coverage of the 95% c.i. of \hat{Q} for Q for the available cases;
- $\hat{E}[\bar{Q}_m]$, the average \hat{Q} after multiple imputation;
- the coverage of the 95% c.i. under multiple imputation.

Choices for \hat{Q} reflect aspects of the distribution of the incomplete variable (e.g., mean, probability of a category, quantiles) or quantify the relation with the predictors (e.g., correlations, conditional means).

The results for $\hat{E}[\hat{Q}_{ac}]$ indicate that available case analysis is often severely biased under MARRIGHT, MARTAIL, and MARMID. Note that, unlike MARRIGHT and MARTAIL, MARMID mechanism generally increases the correlations with the predictors. Almost everywhere, the difference $\hat{E}[\bar{Q}_m] - \hat{Q}$ is much smaller than $\hat{E}[\hat{Q}_{ac}] - \hat{Q}$, so multiple imputation corrects for the biases introduced by MARRIGHT, MARTAIL, and MARMID. For example, the bias of the median (P50) estimated by available case analysis under MARRIGHT

Table 1.	Properties of multiple imputation in a u	nivariate y_1 .
	MARRIGHT	MARTAII
MI	AC MI	1.0

			MCAR				MARRIG	GHT MARTAIL			IL	MARMID		
	Statistic	Pop	Fmi	AC (coverage)	MI (coverage)	Fmi	AC (coverage)	MI (coverage)	Fmi	AC (coverage)	MI (coverage)	Fmi	AC (coverage)	MI (coverage)
Linear	$E(y_1)$	11.66	0.32	11.65 (95)	11.66 (95)	0.44	9.82 (00)	11.65 (95)	0.32	11.53 (91)	11.66 (95)	0.32	11.85 (96)	11.65 (96)
	$P25(y_1)$	8.16	0.35	8.16 (94)	8.15 (97)	0.23	6.77 (08)	8.13 (97)	0.36	8.50 (86)	8.15 (97)	0.36	7.60 (84)	8.13 (98)
	$P50(y_1)$	11.42	0.32	11.39 (95)	11.42 (97)	0.36	9.53(01)	11.39 (97)	0.27	11.35 (92)	11.42 (97)	0.38	11.55 (98)	11.41 (98)
	$P75(y_1)$	14.90	0.30	14.89 (96)	14.92 (98)	0.48	12.53 (00)	14.92 (98)	0.31	14.35 (77)	14.91 (97)	0.31	15.89 (65)	14.92 (98)
	$r(y_1 \cdot x_1)$	0.72	0.43	0.72(95)	0.72(95)	0.54	0.65 (47)	0.72(95)	0.51	0.65 (51)	0.72(95)	0.35	0.79(37)	0.72(95)
	$r(y_1 \cdot x_2)$	0.59	0.39	0.59 (96)	0.59 (96)	0.48	0.50(55)	0.58 (95)	0.42	0.50 (56)	0.59(95)	0.36	0.67 (46)	0.59 (95)
	$r(y_1 \cdot x_3)$	0.66	0.41	0.66 (94)	0.66(95)	0.52	0.59 (59)	0.66(94)	0.47	0.58 (52)	0.66(95)	0.36	0.74(41)	0.66 (96)
	$r(y_1 \cdot x_4)$	0.61	0.40	0.61 (94)	0.60 (95)	0.48	0.52 (53)	0.60 (94)	0.43	0.51 (53)	0.60 (95)	0.37	0.70(40)	0.61 (96)
Logistic	$P(y_1 = 0)$	0.50	0.32	0.50 (95)	0.50 (94)	0.43	0.71(00)	0.51 (95)	0.21	0.51 (91)	0.50 (95)	0.54	0.50 (99)	0.50 (95)
	$P(y_1 = 1)$	0.50	0.32	0.50(95)	0.50(94)	0.43	0.29(00)	0.49(95)	0.21	0.49(91)	0.50(95)	0.54	0.50(99)	0.50(95)
	$E(x_1 y_1=0)$	8.66	0.29	8.67 (96)	8.68 (96)	0.51	8.02 (27)	8.72 (97)	0.26	9.59(19)	8.70 (96)	0.42	7.30(02)	8.69 (94)
	$E(x_1 y_1 = 1)$	16.10	0.23	16.13 (95)	16.09 (95)	0.14	14.16(20)	16.11 (95)	0.19	14.81 (21)	16.07 (96)	0.37	17.91 (04)	16.09 (95)
	$E(x_2 y_1=0)$	9.29	0.24	9.30 (95)	9.31 (96)	0.33	8.91 (75)	9.35 (97)	0.17	9.89 (73)	9.32 (95)	0.42	8.43 (38)	9.32 (95)
	$E(x_2 y_1=1)$	14.05	0.21	14.07 (96)	14.04 (94)	0.19	12.85 (53)	14.04 (95)	0.14	13.26 (60)	14.04 (95)	0.41	15.19 (28)	14.04 (95)
	$E(x_3 y_1=0)$	7.17	0.29	7.18 (96)	7.19 (95)	0.51	6.60(25)	7.21 (97)	0.26	8.00(16)	7.20 (96)	0.43	5.94(02)	7.19 (95)
	$E(x_3 y_1 = 1)$	13.78	0.23	13.80 (96)	13.77 (96)	0.14	12.09(20)	13.78 (95)	0.19	12.67 (21)	13.75 (96)	0.36	15.35 (04)	13.76 (95)
	$E(x_4 y_1=0)$	7.18	0.26	7.19 (95)	7.20(97)	0.41	6.71 (60)	7.21 (96)	0.21	7.80 (58)	7.19 (96)	0.43	6.25 (23)	7.19 (95)
	$E(x_4 y_1=1)$	12.45	0.20	12.47 (93)	12.44 (94)	0.18	10.99 (39)	12.43 (96)	0.15	11.56 (51)	12.44 (96)	0.36	13.74 (23)	12.43 (95)
Polytome	$P(y_1 = 0)$	0.57	0.52	0.57 (96)	0.55 (96)	0.61	0.68 (08)	0.56 (95)	0.54	0.56 (95)	0.55 (95)	0.64	0.58 (96)	0.56 (95)
	$P(y_1 = 1)$	0.25	0.54	0.25 (96)	0.26(97)	0.70	0.17(13)	0.26(96)	0.58	0.25 (95)	0.26(97)	0.65	0.25 (95)	0.26(95)
	$P(y_1 = 2)$	0.18	0.55	0.18 (95)	0.19(97)	0.69	0.15 (75)	0.18 (96)	0.58	0.19 (95)	0.19 (96)	0.64	0.17 (93)	0.18 (96)
	$E(x_1 y_1=0)$	8.06	0.29	388.05 (98)	8.03 (99)	0.38	8.56 (09)	8.05 (99)	0.36	8.12 (97)	8.02 (100)	0.29	7.96 (98)	8.05 (100)
	$E(x_1 y_1=1)$	6.71	0.52	6.70 (95)	6.79 (97)	0.48	7.53 (30)	6.78 (98)	0.54	7.16 (56)	6.82 (98)	0.68	6.02(11)	6.72 (95)
	$E(x_1 y_1=2)$	7.19	0.49	7.19 (96)	7.22 (97)	0.50	8.00 (41)	7.23 (96)	0.54	7.53 (82)	7.25 (96)	0.56	6.63 (58)	7.17 (96)

Note: Given are the population value (Pop), the fraction of missing information (Fmi), the mean estimate under available case analysis and its 95% c.i. coverage (AC), and the mean estimate for MI and its 95% c.i. coverage (MI) under four MAR missing data mechanisms: MCAR, MARRIGHT, MARTAIL, and MARMID (On the basis of m = 5 imputations). P25(y), P50(y), and P75(y) represent the first, second, and third quartiles of y, respectively. Notation $P(\cdot)$ is the marginal probability of observing the argument and $E(\cdot|\cdot)$ the conditional expectation. Notation P(y,x) is used for the Pearson correlation between y and x.

is quite large (11.42 - 9.53 = 1.89), although it is negligible (11.42 - 11.39 = 0.03) after imputation.

The coverages under available case analysis are low. The worst case occurs when the mean of Y_1 is estimated from the available cases under the MARRIGHT mechanism. Here, the 95% confidence interval (c.i.) covered the true value only once (!) in 1000 simulations. In general, the c.i.'s of the available cases have acceptable coverage only under MCAR. In all other cases, c.i.'s are much too short and lead to incorrect statistical inferences.

In contrast, the actual coverage of the 95% c.i. for multiple imputation is generally close to the nominal levels and is nowhere below 93%. Coverage under MARMID is usually even larger than the nominal level. All in all, the results show that the linear, logistic, and polytomous multiple imputation methods for univariate data are on target, well calibrated, and adhere to the validity conditions for proper imputation under a variety of missing data mechanisms.

4. Evaluation of FCS for multivariate imputation of continuous data

This section studies the performance of FCS for multivariate missing values with continuous data.

4.1 Setup: data and simulations

Six locations from the Irish wind speed data are selected: $y_1 = \text{Roche's Pt}$, $y_2 = \text{Rosslare}$, $y_3 = \text{Shannon}$, $y_4 = \text{Dublin}$, $x_1 = \text{Clones}$, and $x_2 = \text{Malin Head}$. Two complete data sets, one simulated and one real, are created. The simulated data consist of approximately 400 cases drawn from the multivariate normal distribution with mean vector and covariance matrix equal to that estimated from the raw wind speed data. This simulated data set presents an idealized case where there are no issues of inaccuracy of model fit. The real data set is a random sample of about 400 observations from the raw wind speed data. Imputing this data yields insight into the robustness of the imputation model in more practical situations. Conditional on the observed data, nonmonotone multivariate missing data were created in Y_1, \ldots, Y_4 using the method described in Appendix B. The percentage of cases that were made incomplete was 62.5.

The incomplete data were multiply imputed using m=10, a relatively high value chosen to account for larger fractions of missing information. The FCS consists of the set of linear regressions of each Y_j on all other variables, Y_{-j} . The number of Gibbs sampling iterations is set to 5. This is a low value in a Gibbs sampling context, but we found it to work well in this type of application using starting values of y^{mis} drawn from each variable's observed marginal distribution. The execution time for generating, imputing, and analyzing the $1000 \times 10 = 10,000$ incomplete data sets was approximately 30 min on a Intel 1.7 GHz processor running SAS 6.2.

4.2 Results

Table 2 provides estimates of the bias and coverage of multiple imputation for many descriptive statistics estimating various quantities of the six-variate distribution $(Y_1, Y_2, Y_3, Y_4, X_1, X_2)$. Under MAR, the available case analysis is often biased, but multiple imputation consistently moves in the right direction, and nearly always repairs the damage done by the missing data mechanism. For example, the correlation between Y_2 (Rosslare) and Y_3 (Shannon) drops from 0.59 to 0.33 for the available cases, but is 0.59 after imputation.

Table 2. Results for multiple imputation of multivariate missing data in y_1, \ldots, y_4 by means of iterated linear regressions (Gibbs sampling) for two data sets (simulated and raw) constructed from the Irish windspeed data. Given are the population value (Pop), the mean estimate under available case analysis and the coverage of its 95% c.i. (AC), the mean estimate under multiple imputation MI and its 95% c.i. (MI). Based on m = 10 imputations. P25(y), P50(y) and P75(y) represent the first, second and third quartile of y respectively. Notation r(y, x) stands for the Pearson correlation between y and x. Variables x_1 and x_2 are complete.

			Simula			Raw data		
Statistic	Pop	Fmi	AC (coverage)	MI (coverage)	Pop	Fmi	AC (coverage)	MI (coverage)
$E(y_1)$	12.38	0.15	11.38 (15)	12.38 (97)	12.36	0.15	11.36 (14)	12.37 (95)
$P25(y_1)$	8.61	0.11	7.67 (49)	8.60 (98)	8.15	0.11	7.39 (46)	8.19 (97)
$P50(y_1)$	12.40	0.17	11.32 (30)	12.40 (98)	11.72	0.17	10.54 (17)	11.79 (96)
$P75(y_1)$	16.25	0.24	15.04 (30)	16.24 (99)	15.88	0.22	14.45 (27)	15.93 (98)
$E(y_2)$	11.69	0.22	10.95 (31)	11.68 (97)	11.65	0.22	10.94 (31)	11.67 (95)
$P25(y_2)$	8.26	0.17	7.56 (58)	8.26 (97)	7.97	0.16	7.36 (43)	7.93 (97)
$P50(y_2)$	11.68	0.23	10.86 (42)	11.64 (98)	10.91	0.21	10.09 (40)	11.04 (97)
$P75(y_2)$	15.16	0.29	14.26 (51)	15.12 (98)	14.64	0.26	13.70 (54)	14.83 (97)
$E(y_3)$	10.48	0.13	9.55 (11)	10.47 (95)	10.45	0.13	9.54 (11)	10.45 (95)
$P25(y_3)$	7.10	0.09	6.26 (40)	7.09 (97)	6.76	0.10	6.09 (43)	6.81 (96)
$P50(y_3)$	10.51	0.16	9.44 (19)	10.48 (97)	9.94	0.14	8.89 (16)	9.98 (97)
$P75(y_3)$	13.83	0.24	12.74 (29)	13.82 (98)	13.52	0.21	12.21 (24)	13.54 (98)
$E(y_4)$	9.81	0.12	8.89 (12)	9.81 (96)	9.78	0.12	8.88 (12)	9.80 (95)
$P25(y_4)$	6.54	0.09	5.70 (42)	6.52 (97)	6.01	0.09	5.35 (45)	6.06 (98)
$P50(y_4)$	9.75	0.15	8.68 (18)	9.77 (98)	9.16	0.14	8.13 (15)	9.24 (96)
$P75(y_4)$	13.14	0.23	12.06 (33)	13.15 (98)	12.88	0.21	11.53 (24)	12.94 (98)
$r(y_1, y_2)$	0.73	0.39	0.70 (82)	0.73 (95)	0.73	0.43	0.69 (74)	0.74 (86)
$r(y_1, y_3)$	0.83	0.37	0.81 (81)	0.83 (94)	0.83	0.40	0.81 (79)	0.85 (86)
$r(y_1, y_4)$	0.74	0.43	0.50(01)	0.74 (95)	0.74	0.53	0.38 (00)	0.77 (79)
$r(y_1, x_1)$	0.75	0.22	0.75 (93)	0.75 (95)	0.75	0.24	0.74 (91)	0.75 (93)
$r(y_1, x_2)$	0.62	0.18	0.62 (94)	0.62 (96)	0.62	0.20	0.61 (92)	0.63 (92)
$r(y_2, y_3)$	0.59	0.43	0.33 (02)	0.59 (97)	0.59	0.50	0.22(00)	0.61 (90)
$r(y_2, y_4)$	0.66	0.37	0.63 (86)	0.66 (96)	0.66	0.39	0.61 (79)	0.68 (92)
$r(y_2, x_1)$	0.61	0.27	0.60 (95)	0.61 (95)	0.61	0.29	0.58 (87)	0.60 (93)
$r(y_2, x_2)$	0.47	0.24	0.46 (95)	0.47 (95)	0.48	0.25	0.46 (89)	0.47 (91)
$r(y_3, y_4)$	0.79	0.35	0.76 (80)	0.79 (95)	0.79	0.40	0.74 (60)	0.78 (93)
$r(y_3, x_1)$	0.82	0.23	0.82 (94)	0.82 (97)	0.82	0.25	0.81 (93)	0.82 (94)
$r(y_3, x_2)$	0.67	0.17	0.66 (95)	0.67 (94)	0.67	0.19	0.66 (90)	0.67 (93)
$r(y_4, x_1)$	0.84	0.24	0.83 (94)	0.84 (96)	0.84	0.24	0.84 (94)	0.85 (91)
$r(y_4, x_2)$	0.77	0.20	0.76 (95)	0.77 (96)	0.77	0.21	0.76 (90)	0.77 (90)

Figure 2 illustrates some properties of the imputation process for Rosslare and Shannon in more detail. The figure on the left is a scatter plot of a sample of about 400 cases from the original Irish wind speed data. The middle figure portrays a subset of this data set. The slope of the regression line is biased downwards, and both variable ranges are limited to about half of the scale. The panel on the right is the first imputed data set after imputation. Note that imputation 'restores' the slope of the regression, the ranges of the variables, and the correlation between them.

Coverages are often close to the nominal value. The average coverage percentage over all statistics is 96.5 for the simulated data and 93.5 for the real data, and thus both are remarkably close to the nominal value of 95. Coverage is occasionally lower than 90, especially for statistics with large fractions of missing information. We observed this also at other runs of the simulations, but at different places. The results provide firm evidence that for both simulated and real data, the Gibbs sampling imputation algorithm is on target and well calibrated under the studied conditions.

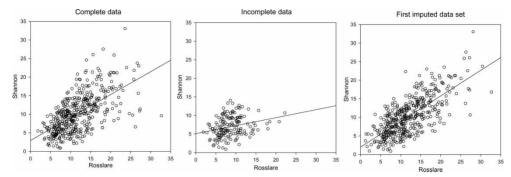


Figure 2. Scatter plots of the locations Rosslare and Shannon from the Irish wind speed data [33]. The incomplete data (middle panel) contain a subset of the complete data (left-hand panel). Data are missing at random (MAR). The right-hand panel illustrates the scatter of cases of the first multiply imputed data set.

5. Evaluation of FCS for multiple imputation of mixed data

5.1 Setup: data and simulations

Multivariate missing data were created in the ME data set [35, p. 265]. As before, imputation of real and simulated data is studied. Nonmonotone missing data under an MAR mechanism were created in four missing data patterns. Each of these patterns was characterized by missing data on one of the following pair of variables: (SYMPT, BSE), (ME, SYMPT), (BSE, DETC) and (SYMPT, DETC). Each pattern occurs in \sim 15.6% of all cases, so the total percentage of incomplete cases is $4 \times 15.6 = 62.5\%$. Incomplete data in SYMPT and BSE are imputed by logistic regression, whereas that in ME and DETC are imputed by polytomous logistic regression. Imputation was always done conditionally on the five other variables, with m = 10 and using five Gibbs sampling iterations from the marginal starting values.

5.2 Results

The results in table 3 indicate that the performance of multiple imputation for multivariate data is quite satisfactory for both data sets. The point estimates are nearly always closer to the true values than under available case analysis, and the empirical coverage of the 95% intervals

Table 3. Bias and coverage of multiple imputation (m = 10) after Gibbs sampling applied to multivariate categorical data when compared with available case analysis.

	Si	imulated data (r	a = 412	Raw data $(n = 412)$			
Statistic	Pop	Fmi	AC (coverage)	MI (coverage)	Fmi	AC (coverage)	MI (coverage)
P(ME = 0) $P(ME = 1)$ $P(ME = 2)$	0.57	0.51	0.59 (87)	0.58 (95)	0.43	0.61 (65)	0.57 (95)
	0.25	0.63	0.23 (83)	0.23 (95)	0.55	0.22 (75)	0.25 (98)
	0.19	0.44	0.19 (96)	0.19 (96)	0.48	0.17 (88)	0.18 (95)
E[PB SYMPT = 0] $E[PB SYMPT = 1]$	8.24	0.22	8.68 (61)	8.29 (95)	0.26	8.51 (76)	8.19 (96)
	7.29	0.09	7.66 (22)	7.28 (96)	0.09	7.71 (20)	7.32 (95)
OR(SYMPT, HIST)	0.37	0.33	0.23 (93)	0.23 (93)	0.29	0.25 (94)	0.27 (95)
OR(SYMPT, BSE)	0.51	0.41	0.28 (93)	0.59 (96)	0.42	0.33 (78)	0.72 (96)
OR(HIST, BSE)	0.50	0.36	0.57 (96)	0.51 (96)	0.38	0.68 (97)	0.60 (98)

Note: Mammography Experience Study (n = 412), where SYMPT was recoded into two categories. Notation $P(\cdot)$ is the marginal probability of observing the argument and $E(\cdot|\cdot)$ the conditional expectation. OR(x, y) is the odds ratio of x and y.

is close to the nominal value. The line labeled 'E[PB|SYMPT=1]' illustrates that multiple imputation is clearly superior to available case even if the amount of missing information is small.

6. Evaluation of FCS for imputation based on incompatible models

This section addresses the potential consequences of incompatibility by means of simulation.

6.1 Setup: data and simulations

For each replication, 1000 draws were made from the bivariate normal distribution $P(Y_1, Y_2)$ with $\mu_1 = \mu_2 = 5$, $\sigma_1^2 = \sigma_2^2 = 1$, and $\rho_{12} = 0.6$. All values generated were positive. Missing data in Y_1 and Y_2 were in three ways:

```
MARRIGHT: logit(Pr(Y_1 = missing)) = -1 + Y_2/5å logit(Pr(Y_2 = missing)) = -1 + Y_1/5;
MARTAIL: logit(Pr(Y_1 = missing)) = -1 + 0.4|Y_2|å logit(Pr(Y_2 = missing))
= -1 + 0.4|Y_1|;
MARMID: 1 - Pr(MARTAIL).
```

Figure 3 plots the probability to be missing under each mechanism as a function of the data. When taken together, these specifications led to zero, one or two missing observations in the pair (Y_1, Y_2) . Under MARRIGHT, there were \sim 50% missing entries, 75% incomplete cases, and \sim 25% of the cases for which both Y_1 and Y_2 were missing. There were proportionally more missing data for the higher values of Y_1 and Y_2 , like in the univariate MARRIGHT mechanism in table 1. The multivariate missing data were not entirely MAR because the cases where Y_1 or Y_2 (or both) is (are) missing were more frequent for the higher values. The regression lines are, however, not affected because the nonresponse is generated symmetrically around the regression lines.

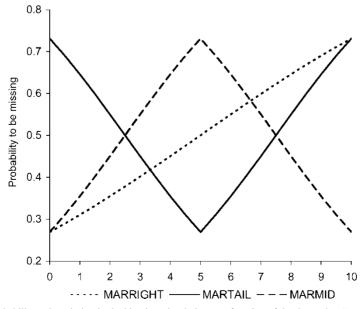


Figure 3. Probability to be missing in the bivariate simulation as a function of the data value (i.e. Y_1 or Y_2) under three missing missing data mechanisms.

Multiple imputations Y_1^* and Y_2^* using m=5 were created using five iterations of the Gibbs sampler. One iteration of the compatible Gibbs sampler consisted of a chain of two univariate imputation models $Y_2^*|Y_1\sim N(\mu_1^*+\beta_1^*Y_1,\sigma_1^{2*})$ and $Y_1^*|Y_2\sim N(\mu_2^*+\beta_2^*Y_2,\sigma_2^{2*})$, where μ_1^* , β_1^* , σ_1^{2*} , μ_2^* , β_2^* , and σ_2^{2*} were draws from the appropriate posterior distributions. The first incompatible conditionally specified model was formulated by replacing the imputation step for Y_2 by $Y_2^*|Y_1\sim N(\mu_1^*+\beta_1^*Y_1^2,\sigma_1^{2*})$, so imputation was conditional on Y_1^2 instead of Y_1 . A second incompatible model used $\log(Y_1)$ instead of Y_1 . The linear model $Y_1=\alpha+\beta Y_2+\varepsilon$ was taken as the complete-data model, where scientific interest focussed on β . The number of replications was set to 500. The average fraction of missing information about β was approximately 0.63, so the imputation problem was quite difficult.

6.2 Results

Table 4 contains the results. The equality $E(\beta_1^*/\sigma_1^{2*}) = E(\beta_2^*/\sigma_2^{2*})$ must hold for the compatible normal model. Note that this equality is empirically obtained for the conditionally specified linear model. For both incompatible models, we find $E(\beta_1^*/\sigma_1^{2*}) \neq E(\beta_2^*/\sigma_2^{2*})$, thus indicating serious incompatibility. Column 5 lists the marginal mean of Y_1 . Under MARRIGHT, CCA is biased because there are proportionally more missing data for higher Y_1 , whereas the imputation methods are closer to the theoretical value, though not completely unbiased. Note that the imputation methods assume MAR, and therefore cannot completely account for all selectivity induced by the missing data mechanism. Column 6 is the mean of the regression weight β over all replications. All methods produce essentially the same regression weight, except for CCA under MARMID and MARTAIL. The standard errors indicate that all multiple imputation models, even the incompatible ones, are more efficient than CCA, which is due to the fact that they use the incomplete data in a more efficient way. The last column is the coverage coefficient, which is equal to the percentage of cases in which the 95% c.i. includes the true value. Coverage is excellent in all cases, again except for CCA under MARMID and MARTAIL.

It appears that the forms of incompatibility as used here do not influence the statistical properties of multiple imputation in any major way. Somewhat surprisingly, we found that,

Table 4. Regression slopes, standard errors, and coverages (95% c.i.) under one compatible and two incompatible multiple imputation models (bivariate normal data, $\rho = 0.6$, n = 1000, m = 5, three symmetric missing data mechanisms, 500 replications) when compared with complete case analysis.

		Compatibil						
Mechanism	Method	$E(\beta_1/\sigma_1^2)$	$E(\beta_2/\sigma_2^2)$	$E(Y_1)$	$E(\beta)$	$E(\operatorname{se}(\beta))$	Fmi	Cov
	Theoretical values			5.00	0.600			95
MARRIGHT	Complete case analysis			4.84	0.597	0.051		93
	MI compatible linear	0.94	0.94	4.91	0.595	0.046	0.63	95
	MI incompatible quadratic	0.09	0.92	4.91	0.589	0.046	0.63	95
	MI incompatible log	4.12	0.90	4.91	0.582	0.047	0.64	95
MARMID	Complete case analysis			5.00	0.678	0.060		79
	MI compatible linear	0.94	0.96	5.00	0.613	0.057	0.75	94
	MI incompatible quadratic	0.09	0.92	5.00	0.601	0.058	0.75	94
	MI incompatible log	4.13	0.90	5.00	0.579	0.058	0.75	94
MARTAIL	Complete case analysis			5.00	0.556	0.040		78
	MI compatible linear	0.95	0.95	5.00	0.596	0.038	0.50	94
	MI incompatible quadratic	0.09	0.93	5.00	0.590	0.038	0.50	94
	MI incompatible log	4.35	0.93	5.00	0.590	0.037	0.50	95

Note: Fmi = fraction of missing information; Cov = 95% c.i. coverage.

in the cases studied, applying a deliberately specified incompatible method gives less bias and more efficiency than CCA. Thus, imputation using the Gibbs sampler seems to be robust against incompatible-specified conditionals in terms of bias and precision, thus suggesting that incompatibility may be a relatively minor problem in multivariate imputation.

7. Discussion

FCS is a convenient and powerful approach for creating imputations in multivariate missing data. Its theoretical weakness is that convergence can only be guaranteed under compatibility of conditionals, a condition that is often difficult to verify in practice. Our simulations show that this weakness does not seem to affect the quality of imputation in the cases considered. In fact, even for clearly incompatible models, the method produces reasonable multiple imputations with appropriate coverage. Thus, FCS appears to be robust against incompatibility. We suspect that this result will hold more generally, but more work is needed to explore the boundaries of this conclusion.

The amount of work per iteration can be substantial, but only relatively few iterations appear to be needed when using well-chosen starting values. All simulations were done with just five iterations. Increasing the number of iterations did not result in a lower bias and a better coverage. This is unlike many Markov Chain Monte Carlo methods that often require thousands of iterations. Of course, we only studied simple models with not more than four incomplete variables, and therefore do not suggest that a small number of iterations would also be sufficient for larger or more complex models. We have also seen applications with large fractions of missing data and high correlations that required several hundreds of iterations before reaching some stability. In our limited experience, using 20 iterations for modest missing data problems (<10–15% missing data) is ample. In more demanding problems, convergence of critical parameters should be carefully monitored, for example, by the method of multiple sequences [36]. A particular difficulty here is how to specify overdispersed starting values in the context of multivariate missing data, which is an area for further research. The costs of drawing a longer chain are only computational, so if the implementation is efficient, the benefits of faster convergence will be small.

The major advantage of FCS is increased flexibility in model building. It is easy to incorporate constraints on the imputed values, work with different transformations of the same variable, account for skip patterns, rounding, and so on. We concentrated on models containing main effects only. It is, however, straightforward to build imputation models that preserve higher order interactions between variables. Following imputation of the main effect, all interaction terms containing this effect can be immediately updated, thus preserving consistency across the data. It would be worthwhile to study how robust multiple imputation along these lines would be in preserving higher-order interactions.

Throughout the paper, the imputation models assumed that MAR holds. It is relatively easy to adapt the method for models that are not MAR, but assumptions outside the data will be needed. Another extension is to impute vectors instead of scalars. This may be helpful if the relationships between the variables are difficult to model or to speed up the method. Of course, we do not know whether our results will hold in such more complex cases, but the evidence obtained thus far suggests that FCS might also work well in such situations. Using FCS in concert with monotone missing data methods, as in Rubin [19], appears to be particularly attractive because the potential for incompatibility is reduced relative to FCS. Of course, there is always some point at which the technique breaks down, but conditional specification in multivariate imputation seems to be remarkably robust and well worth investigating further.

Acknowledgements

The authors thank the anonymous reviewer for providing extremely helpful comments.

References

- [1] Little, R.J.A. and Rubin, D.B., 2002, Statistical Analysis with Missing Data. (2nd edn) (New York: Wiley).
- [2] Schafer, J., 1997, Analysis of Incomplete Multivariate Data (London: Chapman and Hall).
- [3] Rubin, D.B., 1987, Multiple Imputation for Nonresponse in Surveys (New York: John Wiley).
- [4] Rubin, D.B., 2004, Multiple Imputation for Nonresponse in Surveys (Reprint) (New York: John Wiley).
- [5] Rubin, D.B., 1996, Multiple imputation after 18+ years. *Journal of the American Statistical Association*, 91, 473–489.
- [6] Gelman, A. and Raghunathan, T.E., 2001, Discussion of Arnold et al. 'Conditionally specified distributions'. Statistical Science, 16, 249–274.
- [7] Buck, S.F., 1960, A method of estimation of missing values in multivariate data suitable for use with an electronic computer. *Journal of the Royal Statistical Society B*, **22**, 302–306.
- [8] Gleason, T.C. and Staelin, R., 1975, A proposal for handling missing data. Psychometrika, 40, 229-252.
- [9] Finkbeiner, C., 1979, Estimation for the multiple factor model when data are missing. *Psychometrika*, 44, 409–420.
- [10] Raymond, M.R. and Roberts, D.M., 1987, A comparison of methods for treating incomplete data in selection research. Educational and Psychological Measurement, 47, 13–26.
- [11] Jinn, J.-H. and Sedransk, J., 1989, Effect on secondary data analysis of common imputation methods. Sociological Methodology, 19, 213–241.
- [12] Gold, M.S. and Bentler, P.M., 2000, Treatments of missing data: a Monte Carlo comparison of RBHDI, iterative stochastic regression imputation, and expectation–maximization. *Structural Equation Modeling*, 7, 319–355.
- [13] Kennickell, A.B., 1991, Imputation of the 1989 survey of consumer finances: stochastic relaxation and multiple imputation. ASA 1991 Proceedings of the Section on Survey Research Methods, Alexandria, ASA, pp. 1–10.
- [14] Raghunathan, T.E., Solenberger, P. and van Hoewyk, J., 2000, IVEware: imputation and variance estimation software: installation instructions and user guide. Survey Research Center, Institute of Social Research, University of Michigan. Available online at: http://www.isr.umich.edu/src/smp/ive/
- [15] Brand, J.P.L., 1999, Development, implementation and evaluation of multiple imputation strategies for the statistical analysis of incomplete data sets. Dissertation, Erasmus University Rotterdam.
- [16] Van Buuren, S., van Rijckevorsel, J.L.A. and Rubin, D.B., 1993, Multiple imputation by splines. Bulletin of the International Statistical Institute, Contributed Papers II, 503–504.
- [17] Van Buuren, S. and Oudshoorn C.G.M., 2000, Multivariate imputation by chained equations: MICE V1.0 User's manual. Report PG/VGZ/00.038. Leiden, TNO Preventie en Gezondheid.
- [18] Royston, P., 2004, Multiple imputation of missing values. The Stata Journal, 4, 227–241.
- [19] Rubin, D.B., 2003, Nested multiple imputation of NMES via partially incompatible MCMC. Statistica Neerlandica, 57, 3–18.
- [20] Kennickell, A.B., 1999, Multiple imputation and disclosure control: the case of the 1995 Survey of Consumer Finances. Record Linkage Techniques 1997 (Washington, DC: National Academy Press), pp. 248–267.
- [21] Van Buuren, S., Boshuizen, H.C. and Knook, D.L., 1999, Multiple imputation of missing blood pressure covariates in survival analysis. Statistics in Medicine, 18, 681–694.
- [22] Raghunathan, T.E. and Siscovick, D.S., 1996, A multiple imputation analysis of a case–control study of the risk of primary cardiac arrest among pharmacologically treated hypertensives. *Applied Statistics*, 45, 335–352.
- [23] Oudshoorn, C.G.M., van Buuren, S. and van Rijckevorsel, J.L.A., 1999, Flexible multiple imputation by chained equations of the AVO-95 survey. Report PG/VGZ/99.045. Leiden, TNO Prevention and Health. Available online at: http://www.multiple-imputation.com
- [24] Heeringa, S.G., Little, R.J.A. and Raghunathan, T.E., 2002, Multivariate imputation of coarsened survey data on household wealth. In: R.M. Groves *et al.* (Eds), *Survey Nonresponse* (New York: Wiley), pp. 357–371.
- [25] Faris, P.D., Ghali, W.A., Brant, R., Norris, C.M., Galbraith, P.D. and Knudtson, M.L., 2002, Multiple imputation versus data enhancement for dealing with missing data in observational health care outcome analyses. *Journal* of Clinical Epidemiology, 55, 184–191.
- [26] Horton, N.J. and Lipsitz, S.R., 2001, Multiple imputation in practice: comparison of software packages for regression models with missing variables. *American Statistician*, 55, 244–254.
- [27] Raghunathan, T.E., Lepkowski, J.M., van Hoewyk, J. and Solenberger, P., 2001, A multivariate technique for multiply imputing missing values using a sequence of regression models. Survey Methodology, 27, 85–95.
- [28] Brand, J.P.L., van Buuren, S., Groothuis-Oudshoorn, C.G.M. and Gelsema, E.S., 2003, A toolkit in SAS for the evaluation of multiple imputation methods. *Statistica Neerlandica*, 57, 36–45.
- [29] Arnold, B.C. and Press, S.J., 1989, Compatible conditional distributions. *Journal of the American Statistical Association*, 84, 152–156.
- [30] Bhattacharryya, A., 1943, On some sets of sufficient conditions leading to the normal bivariate distribution. *Sankhya*, **6**, 399–406.
- [31] Besag, J., 1974, Spatial interaction and the statistical analysis of lattice systems (with discussions). Journal of the Royal Statistical Society, Series B, 36, 192–236.

- [32] Arnold, B.C., Castillo, E., and Sarabia, J.M., 1999, Conditional Specification of Statistical Models (New York: Springer-Verlag).
- [33] Hasslet, J. and Raftery, A.E., 1989, Space-time modeling with long-memory dependence: assessing Ireland's wind power resource. Applied Statistics, 38, 1-50.
- [34] Rubin, D.B., 1976, Inference and missing data, *Biometrika*, **63**, 581–592.
- [35] Hosmer, D.W. and Lemeshow, S., 2000, Applied Logistic Regression (2nd edn) (New York: John Wiley).
- [36] Gelman, A. and Rubin, D.B., 1992, Inference from iterative simulation using multiple sequences (with discussion). Statistical Science, 7, 457–511.
- [37] Clogg, C.C., Rubin, D.B., Schenker, N., Schultz, B. and Weidman, L., 1991, Multiple Imputation of industry and occupation codes in census public-use samples using Bayesian logistic regression. Journal of the American Statistical Association, 86 (413), 68-78.

Appendix A: algorithms for univariate imputation

Depending on the distribution of y given x, concrete algorithms are as follows.

For normally distributed y with mean βx and variance σ^2 , $x = (x^{\text{obs}}, x^{\text{mis}})$ is the $n \times p$ matrix of covariates and $n = n^{\text{obs}} + n^{\text{mis}}$ [3, p. 167]:

- 1. estimate β by $b = (x^{\text{obs}}/x^{\text{obs}})_{-1}(x^{\text{obs}})'y^{\text{obs}}$;
- 2. (A) draw a random variable $g \sim \chi^2(n^{\text{obs}} p)$;
 - (B) calculate $\sigma^{*^2} = (y^{\text{obs}} x^{\text{obs}}b)'(y^{\text{obs}} x^{\text{obs}}b)/g$;
- 3. (A) draw $w_1 \sim N(0, I_p)$, i.e., p independent N(0, 1) variates, where I_p is the identity matrix of order p;
 - (B) calculate $b^* = b + \sigma^* w_1 V^{1/2}$, where $V^{1/2}$ is the triangular square root of V = $(x^{\text{obs}}, x^{\text{obs}})^{-1}$ obtained by Cholesky factorization;
- 4. (A) draw $w_2 \sim N(0, I_n^{\text{mis}});$ (B) calculate $y^* = x^{\text{mis}}b^* + w_2\sigma^*.$

Closely related algorithms that account for deviations from the normal distribution are as follows.

Predictive mean matching. Replace step 4 by: calculate $y^{\text{mis}} = x^{\text{mis}}b^*$. For each missing value $i = 1, ..., n^{\text{mis}}$, find the respondent whose $y^{\text{obs}} = x^{\text{obs}}b^*$ is closest to y_i^{mis} and take y^{obs} of this case as the imputed value of i.

Hot-deck version. Replace in step 4 the draws of the n^{mis} normal deviates with draws of n^{mis} values with replacement from the set of n^{obs} observed standardized residuals $\{(y_i^{\text{obs}} - bx_i^{\text{obs}})(1 - p/n^{\text{obs}})^{-1/2}/s\}$, where $s = \sum_{\text{obs}} (y_i - x_i b)^2/(n^{\text{obs}} - p)$.

For dichotomous y, we assume $P(y_i|x_i, \beta) = (\exp(x_i\beta)/(1 + \exp(x_i\beta)))^{y_i}(1 - \exp(x_i\beta)/(1 + \exp(x_i\beta)))^{y_i}$ $(1 + \exp(x_i \beta))^{1-y_i}$. Imputations of y^{mis} are obtained as follows:

- 1. calculate by an iterative algorithm b, the MLE estimate of β and an estimate of the posterior variance of β (e.g., the Hessian matrix in $\beta = b$);
- 2. (A) draw $b^* \sim N(b, V(b))$;
- (B) calculate for $i=1,\ldots,n^{\text{mis}},$ $w_i=\exp(x_ib^*)/(1+\exp(x_ib^*);$ 3. draw $u_i\sim \text{unif}(0,1),$ $i=1,\ldots,n^{\text{mis}}.$ If $u_i>w_i,$ impute $y_i=0,$ otherwise impute $y_i=1.$ For small samples, this procedure can be improved by SIR, as in Clogg *et al.* [37].

A predictive mean-matching version of this algorithm replaces step 3 by: calculate $y^{mis} =$ $x^{\text{mis}}b^*$. For each missing value $i=1,\ldots,n^{\text{mis}}$, find the respondent whose $y^{\text{obs}}=x^{\text{obs}}b^*$ is closest to y_i^{mis} and take y_i^{obs} of this case as the imputed value of i.

For categorical y with unordered categories denoted by $0, \ldots, s-1$, suppose that the distribution of y can be characterized as $\ln(P(y=j|x)/P(y=0|x)) = \beta_j x$, for $j = 1, \dots, s - 1$, so the model for y is a series of separate logistic regression models of categories $1, \ldots, s-1$ against baseline category 0. An appropriate algorithm for this model is as follows.

- 1. Draw b^* from N(b, V(b)) where b is the MLE estimator of $\beta = (\beta_1, \dots, \beta_{s-1})$ and V(b) its estimated covariance matrix.
- 2. Calculate $\pi_{ij}^{\text{mis}} = \exp(-b_j^* x_i^{\text{mis}})/(1 + \sum_{v=1}^{s-1} \exp(-(b_v^* x_i^{\text{mis}}))$ for $i = 1, \dots, n^{\text{mis}}, j = 0, \dots, s-1$, and $b_0 = (0, \dots, 0)$.
- 3. Draw y_i^* from $\{0, \ldots, s-1\}$ with probabilities π_{ij}^{mis} , $i=1,\ldots,n^{\text{mis}}$, $j=0,\ldots,s-1$.

Appendix B: generation of the missing data

This appendix describes the method developed by Brand [15, pp. 110–113] for generating nonmonotone multivariate missing entries in J variables Y_1, \ldots, Y_J under MAR. We assume that Y_1, \ldots, Y_J are initially completely known. Additional complete covariates X_1, \ldots, X_L can be present for which no missing entries are sought. The method requires specification of the proportion of incomplete cases, the patterns of missing data that are allowed, the relative frequency of each pattern, and a specification of the way in which the observed information influence the response probability of each pattern.

More in particular, for a sample size n, let $\alpha(0 < \alpha < 1)$ denote the desired proportion of incomplete cases. Let there be P missing data patterns R_1, \ldots, R_P , chosen by the user, where $R_P = \{r_{p1}, \ldots, r_{pJ}\}$ is a 0-1 response indicator vector of length J, with $r_{pj} = 0$ if variable Y_j is missing and $r_{pj} = 1$ otherwise. All response patterns except $(0, 0, \ldots, 0)$ or $(1, 1, \ldots, 1)$ may occur. Furthermore, let the vector $f = (f_1, \ldots, f_P)$ specify the relative frequencies of patterns R_1, \ldots, R_P , with $\sum_p f_p = 1$, also specified by the user.

Each case is randomly allocated to one of P candidate blocks with probability f_p . Within each candidate block, a subgroup of $\alpha n f_p$ cases is made incomplete according to pattern R_p using a probability model as follows. First, calculate a linear score $s_i = \sum_j^J a_{pj} r_{pj} Y_{ij} + \sum_l^L b_{pl} X_{il}$ for each case in the block, where a_{pj} and b_{pl} are user weights specific to pattern p. A convenient choice for a_{pj} and b_{pl} is the set of regression weights from the linear regression of Y_j on $\{Y_{-j}, X_1, \ldots, X_L\}$ as computed from the initially complete data. Subsequently, divide the nf_p cases within the candidate block p into k_p subgroups using their value s_i . The user can control the composition of each candidate subgroup by specifying $k_p - 1$ break points q_{pk} for $k = 1, \ldots, k_p - 1$ (in the form of quantiles). In addition, specify for each subgroup $h_k(2 < k = k_p)$ the odds w_{pk} of having response pattern R_P relative to the reference subgroup h_1 . Together with α , these odds determine the probability on response pattern R_P for each case in the candidate block. For each case, a random draw from the uniform distribution is made. If this random draw does not exceed the probability on response pattern R_P , the data for that case are set to missing according to response pattern R_P . The procedure is repeated for every candidate block.

Choices for the subgroup size and the odds govern the properties of the incomplete data. For example, MARMID-like mechanisms have high missingness odds for cases with average s_i scores, whereas MARTAIL-like mechanisms are obtained by specifying higher missingness odds for extreme s_i scores. All of these are MAR, as the linear score s_i depends on the observed data only.

As an example, the simulations for table 2 generated missing data in the multivariate data $\{Y_1, \ldots, Y_4, X_1, X_2\}$ using the following settings: P = 4, $R = \{010111, 001111, 110011, 101011\}$, $f_1 = f_2 = f_3 = f_4 = 0.25$, $\alpha = 0.625$ $k_1 = k_2 = k_3 = k_4 = 2$, $q_{p1} = 0.5$, and $w_{pk} = 4$ with $p = 1, \ldots, P$ and $k = 1, \ldots, k_p - 1$.