



# Semiparametric Bayesian multiple imputation for regression models with missing mixed continuous–discrete covariates

Ryo Kato<sup>1</sup> · Takahiro Hoshino<sup>2,3</sup>

Received: 17 July 2018 / Revised: 10 December 2018 / Published online: 11 March 2019  
© The Institute of Statistical Mathematics, Tokyo 2019

## Abstract

Issues regarding missing data are critical in observational and experimental research. Recently, for datasets with mixed continuous–discrete variables, multiple imputation by chained equation (MICE) has been widely used, although MICE may yield severely biased estimates. We propose a new semiparametric Bayes multiple imputation approach that can deal with continuous and discrete variables. This enables us to overcome the shortcomings of MICE; they must satisfy strong conditions (known as compatibility) to guarantee obtained estimators are consistent. Our simulation studies show the coverage probability of 95% interval calculated using MICE can be less than 1%, while the MSE of the proposed can be less than one-fiftieth. We applied our method to the Alzheimer’s Disease Neuroimaging Initiative (ADNI) dataset, and the results are consistent with those of the previous works that used panel data other than ADNI database, whereas the existing methods, such as MICE, resulted in inconsistent results.

**Keywords** Full conditional specification · Missing data · Multiple imputation · Probit stick-breaking process mixture · Semiparametric Bayes model

---

**Electronic supplementary material** The online version of this article (<https://doi.org/10.1007/s10463-019-00710-w>) contains supplementary material, which is available to authorized users.

---

✉ Takahiro Hoshino  
bayesian@jasmine.ocn.ne.jp  
Ryo Kato  
kato.ryo@keio.jp

<sup>1</sup> Research Institute for Economics and Business Administration, Kobe University, 2-1 Rokkodai-cho, Nada-ku, Kobe, Japan

<sup>2</sup> Department of Economics, Keio University, 2-15-45 Mita, Minato-ku, Tokyo, Japan

<sup>3</sup> RIKEN Center for Advanced Intelligence Project, 1-4-1 Nihonbashi, Chuo-ku, Tokyo, Japan

## 1 Introduction

Issues regarding missing data are critical in observational and experimental research, as they induce loss of information and biased result, and at times, lead to wrong decisions. [The National Research Council \(2010\)](#) published a report including recommendations on treating missing data in medical science research. According to the recommendation, researchers should employ as many confounders as possible in order to obtain valid estimates from analysis such as logistic regressions or Cox proportional hazards model. However, when they employ more covariates, the number of observations with at least one missing component increases. Also, if a researcher is interested in using a regression model containing missing components in covariates, a complete case analysis results in biased estimates even when the missing mechanism is missing at random (MAR) ([Ibrahim et al. 2005](#)).

In this case, if conditional distributions of incompletely observed covariates, given completely observed covariates, are correctly specified, we can obtain consistent estimators using the expectation-maximization (EM) algorithm or Bayesian estimation with the Markov chain Monte Carlo (MCMC) method. However, it is usually difficult to specify such a distribution because both incompletely and completely observed covariates generally have large dimensions, and the distributional form is not expressed by well-known distributional families owing to the mixed-scale variables.

For datasets with mixed continuous and discrete variables in various study areas, multiple imputation by chained equation (MICE), in which missing variables are iteratively imputed based on full conditional specification (FCS), has been cited numerous times by researchers from several fields including medical statistics ([van Buuren 2007](#); [White et al. 2011](#); [Paton et al. 2014](#)). This is because the researchers, especially the imputers, are not required to construct an explicit joint multivariate model with mixed-scale variables (continuous, categorical, ordinal, and so on). More specifically, the MICE-FCS approach specifies a multivariate imputation model using a sequence of seemingly “appropriate” univariate regression models corresponding to the types of missing variables; namely, one only needs to assign a univariate linear regression with a normally distributed error term for an incomplete continuous variable, a logistic regression for an incomplete binary variable, an ordered logistic regression for an incomplete ordinal variable, and so on. Moreover, researchers can easily implement MICE-FCS using several existing statistical software packages, such as the `mice` package in R and S-plus, `proc mi` with the FCS option in SAS, and `mi impute` in STATA.

In spite of the widespread use of MICE-FCS, recent studies showed that it leads to severely biased estimates in various setups. [Liu et al. \(2014\)](#) proved that using MICE-FCS does not guarantee that the asymptotic distribution is equivalent with the existing Bayesian MI estimator when the families of the conditional models are “incompatible” [see Section 4 in [Liu et al. \(2014\)](#)]. If the parameters of the conditional distribution cannot be represented by the parameters of the joint distribution, the conditional models are said to be incompatible. In fact, simulation studies by [Bartlett et al. \(2015\)](#) showed that MICE yields biased estimates when treating incompatible conditional models. Unfortunately, violation of the compatibility assumption is not uncommon (the example of the violation of the compatibility assumption is provided

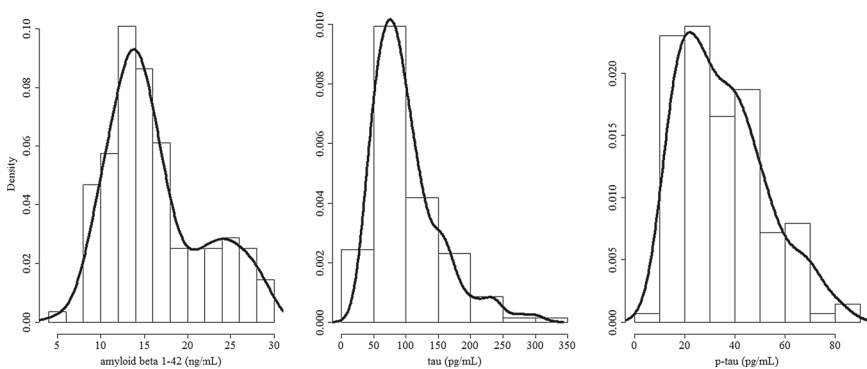
in Sect. 2.1). Therefore, although MICE-FCS is simple and convenient to use, it can provide statistically valid estimates in very limited cases.

### 1.1 Motivating example

We briefly introduce a motivating example of a real-world dataset in which it is very hard to properly impute the missing components using FCS approach. The data used in this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database ([adni.loni.usc.edu](http://adni.loni.usc.edu)). The ADNI was launched in 2003 as a public-private partnership, led by Principal Investigator Michael W. Weiner, MD. The primary goal of ADNI has been to test whether serial magnetic resonance imaging (MRI), positron emission tomography (PET), other biological markers, and clinical and neuropsychological assessment can be combined to measure the progression of mild cognitive impairment (MCI) and early Alzheimer's disease (AD). For up-to-date information, see <http://www.adni-info.org>.

Jack et al. (2010) used data from the ANDI to study baseline predictors that contribute to the progression of AD. Figure 1 shows the distributions of amyloid  $\beta_{1-42}$  ( $ABETA_{A1-42}$ ), tau (total tau protein), and p-tau (phosphorylated tau protein; P-tau181p) of MCI at baseline subjects, and Jack et al. (2010) show that they are associated with time to conversion to AD. Since the participants were not forced to contribute to the CSF (cerebrospinal fluid) measurement, around 50% of the data for  $ABETA_{A1-42}$ , tau, and p-tau are missing. Jack et al. (2010) employed the Cox proportional hazards model, in which the covariates include  $ABETA_{A1-42}$ , tau, and p-tau. The analysis was restricted to 218 samples whose  $ABETA_{A1-42}$  were available; nevertheless, the dataset contains approximately 400 subjects. If the researchers try to address the missing components using the MI method, EM algorithm, or Bayesian estimation with the MCMC technique, they must correctly specify the complex joint distribution of the covariates.

Figure 1 shows that they are not normally distributed and seem to be skewed, following a fat-tailed distribution. Accordingly, specifying the covariate distribution seems



**Fig. 1** Histograms of observed  $ABETA_{A1-42}$  (Amyloid beta 1–42), tau (total tau protein), and p-tau (phosphorylated tau protein). The bold lines represent the kernel densities of the data

to be very difficult in such cases. [Bartlett et al. \(2015\)](#) employed the SMC-FCS (substantive model compatible—fully conditional specification) approach to impute the missing covariates and estimated the Cox regression. They added information pertaining to the family history of Alzheimer’s patients, namely whether the subject’s father and mother had AD or not. However, their results seem to be biased owing to the violation of the model compatibility assumption required by the FCS approaches; several covariates do not follow normal distributions, as seen from [Fig. 1](#). Also, more than two binary missing covariates are employed. As shown in [Sect. 2](#), the FCS approaches result in biased estimates when missing covariates include two or more binary variables because of the model incompatibility. Therefore, it is evident that the existing FCS approaches should not be applied to these kinds of datasets, which are often found in the real world.

## 1.2 New contribution

In this paper, we propose a new flexible semiparametric Bayesian framework for MI, which is capable of treating mixed-scale incomplete variables. The model formulation is different from that seen in the existing literature in two ways.

First, we express the full model as the product of the covariate distribution (conditional distribution of incompletely observed covariates given completely observed covariates) and the substantive model (the regression model researchers are interested in). We assume the parametric model to the substantive model since the researchers conducting applied research are generally concerned with the parameters of the functions in the substantive model, which should be built upon the existing theories or the previous literature in the field of study. Examples of the parametric substantive model are the Cox regression and the logistic regression in epidemiological and clinical research. On the other hand, with regard to the covariate distribution, we specify a joint distribution of the missing variables using the probit stick-breaking process mixture (PSBPM) model proposed by [Chung and Dunson \(2009\)](#), whose model specification is based on the Dirichlet process mixture (DPM) model. [Ibrahim et al. \(2005\)](#) also pointed out that one of the caveats of treating missing covariates lies in specifying the parametric model of the covariate distribution. However, it is nearly impossible to correctly prespecify the covariate distribution based on existing theories or some inferences, because the relationships of the missing variable and the complete variables are often “multivariate-to-multivariate,” they can be nonlinear relationships, or they may be non-normally distributed. Therefore, we employ the nonparametric Bayesian specification; specifically, we use PSBPM modeling instead of DPM since the stick-breaking weights can vary depending on the predictors. Since our approach does not rely on FCS approach, we do not have to consider the compatibility assumption holding.

Second, we express mixed-scale variables through the transformation of the latent continuous variables for probit modeling. This underlying continuous variables approach is used in the context of the DPM model, as in [Kottas et al. \(2005\)](#) for ordinal variables; in [Canale and Dunson \(2011\)](#) for count variables; and in [Kim and Ratchford \(2013\)](#) for ordinal variables. This approach enables us to deal in a straightforward man-

ner with many types of variables in the joint covariate distribution without specifying the complicated conditional joint distribution of mixed-scale variables.

Semiparametric model development is also motivated by the previous frequentist works about semiparametric model for missing variables. For example, [Robins et al. \(1995\)](#) developed semiparametric efficient regression model where conditional distribution of missing covariates can be nonparametric. [Lawless et al. \(1999\)](#) proposed semiparametric likelihood-based model when missingness only depends on a stratification in parametric regression model. [Zhang and Rockette \(2006\)](#) proposed semiparametric maximum likelihood inference with missing covariates, where while they assume parametric regression model, marginal distribution of the covariates can be nonparametric. These frequentist methods assume the substantive regression model to be parametric, and other distributions, such as missing probability or marginal distribution of covariates, are nonparametric. Good statistical property and usefulness toward the real data of semiparametric model are implied by these studies.

One of the Bayesian-related works to our research is [Murray and Reiter \(2016\)](#), which developed fully nonparametric multiple imputation methods for mixed-scale variable (hereafter, NP-MI). They consider [Rubin \(1987\)](#)'s-type multiple imputation, that is, "two step" procedure of imputation stage and analysis stage with Rubin's rules. Multiple imputation approaches like [Murray and Reiter \(2016\)](#) (as well as MICE-FCS) requires strong condition that the analysis model is congenial to the imputation model in order to be proper multiple imputation ([Meng 1994](#)). In addition, proper multiple imputation should satisfy self-efficiency condition to correctly estimate the variance of interested parameter with Rubin's rules. On the contrary, our model incorporates substantive model of interest such as proportional hazard or logistic regression models, and is different from uncongenial multiple imputation. Therefore, we can obtain better estimates even when the regression model of interest is thought to be uncongenial. As [Murray and Reiter \(2016\)](#) states, theoretically evaluating whether the substantive regression model (analysis model) is congenial to their nonparametric imputation model is generally difficult. However, the simulation studies below, in fact, show that their method resulted in about 2–60 times larger MSE compared with our proposed method even in the case of very simple proportional hazard or logistic regression models.

### 1.3 Organization

The rest of the paper is organized as follows. In the next section, we propose and formulate a semiparametric Bayesian multiple imputation (SB-MI) algorithm that can overcome the drawbacks of the existing methods. In Sect. 3, we describe the model specification, imputation procedure and posterior computation of the proposed model in detail. The simulation studies illustrating the performance of the proposed method compared with the MICE-FCS, SMC-FCS, NP-MI, and missForest approaches are presented in Sect. 4. In Sect. 5, we apply our proposed method to the real dataset described in the motivating example in Sect. 1.1. Section 6 concludes after providing a short discussion. The detailed descriptions of the simulation design and some additional analysis appear in the Appendix of ESM.

## 2 Model setup

In this paper, we consider a dataset consisting of  $N$  ( $i = 1, \dots, N$ ) cases, where the interest of the researchers lies in a model with outcomes  $\mathbf{y} : j \times 1$ , completely observed covariates  $\mathbf{v} : p \times 1$ , and incompletely observed covariates  $\mathbf{w} : q \times 1$ . We consider the case where some components of  $\mathbf{y}$  can be missing. Let  $\mathbf{r}$  be the vector of observation indicators whose element equals 1 if the corresponding element of the dataset is observed and 0 otherwise. Throughout this paper, we consider that the data are MAR. To be more precise, if  $\mathbf{y}$  has missing components, we assume  $p(\mathbf{r}|\mathbf{y}, \mathbf{w}, \mathbf{v}) = p(\mathbf{r}|\mathbf{v})$ . Otherwise if  $\mathbf{y}$  is completely observed,  $p(\mathbf{r}|\mathbf{y}, \mathbf{w}, \mathbf{v}) = p(\mathbf{r}|\mathbf{y}, \mathbf{v})$  is assumed. Additionally, we assume that all the observations are independent and identically distributed.

Let  $\boldsymbol{\vartheta}_s$  be the parameter vectors of the substantive model  $p(\mathbf{y}|\mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)$ . Note that the researcher's prime target of inference lies in  $\boldsymbol{\vartheta}_s$ , even in the context of missing data analysis. We propose a SB-MI algorithm expressed by the following imputation model as the product of two submodels:

$$p(\mathbf{y}, \mathbf{w}|\mathbf{v}, \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) = p(\mathbf{y}|\mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)p(\mathbf{w}|\mathbf{v}, \boldsymbol{\vartheta}_m) \quad (1)$$

where  $p(\mathbf{w}|\mathbf{v}, \boldsymbol{\vartheta}_m)$  represents the covariate distribution with parameters of lower interest  $\boldsymbol{\vartheta}_m$ , and  $p(\mathbf{y}|\mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)$  represents the substantive model with parameters of higher interest  $\boldsymbol{\vartheta}_s$ . We assume the parametric model for the substantive model  $p(\mathbf{y}|\mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)$  since the researchers' concern generally lies in the parameters of the substantive model  $\boldsymbol{\vartheta}_s$ . Additionally, the substantive model should be built upon the existing knowledge corresponding to the purpose of the study. Accordingly, the researcher may assume a linear regression with an interaction term, a Cox regression, or a logistic regression for  $p(\mathbf{y}|\mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)$ . While we assume a parametric structure for the substantive model, we do consider Bayesian nonparametric form rather than a parametric form for the covariate distribution  $p(\mathbf{w}|\mathbf{v}, \boldsymbol{\vartheta}_m)$ , because researchers generally have no interest in  $\boldsymbol{\vartheta}_m$  and parametric modeling of a large number of covariates can result in misspecification bias (Chib 2007). Moreover, we express mixed-scale variables through the transformation of latent continuous variables for probit modeling in order to deal with many types of continuous and discrete variables in the joint covariate distribution in a straightforward manner. This transformation enables us to avoid considering the compatibility assumption, which is required in MICE-FCS or SMC-FCS.

### 2.1 Existing method

In this situation,  $\boldsymbol{\vartheta}_s$  can be estimated by the existing MI methodology, EM algorithm, or Bayesian MCMC estimation. For example, MI uses the MCMC approach, and researchers iteratively draw the parameter of the joint model  $\boldsymbol{\psi}$  from  $p(\boldsymbol{\psi}|\mathbf{a})$ , and then draw  $\mathbf{a}^{\text{mis}}$  from  $p(\mathbf{a}^{\text{mis}}|\mathbf{a}^{\text{obs}}, \boldsymbol{\psi})$   $D$  times, where  $\mathbf{a}^{\text{obs}}$  and  $\mathbf{a}^{\text{mis}}$  are the observed and missing subsets of variable in the dataset, respectively. However, as is the case with the maximum likelihood estimation of the EM algorithm or the Bayesian MCMC estimation, it is difficult to correctly specify the joint distribution of all the variables

that have missing elements  $p(\mathbf{a}^{\text{mis}}|\mathbf{a}^{\text{obs}}, \boldsymbol{\psi})$ , especially when both continuous and discrete variables are missing.

The MICE-FCS method has become a more widely used methodology as researchers seek to avoid the difficulty in specifying the conditional joint distribution  $p(\mathbf{a}^{\text{mis}}|\mathbf{a}^{\text{obs}}, \boldsymbol{\psi})$ . The MICE-FCS approach specifies a multivariate covariate distribution by a sequence of univariate regressions for each missing variable. More specifically, MICE-FCS iterates drawing  $\boldsymbol{\psi}_j$  from  $p(\boldsymbol{\psi}_j|\mathbf{a})$  and  $\mathbf{a}^{\text{mis}}$  from  $p(\mathbf{a}_j^{\text{mis}}|\mathbf{a}_{-j}^{\text{mis}}, \mathbf{a}^{\text{obs}}, \boldsymbol{\psi}_j)$  for each  $\mathbf{a}_j^{\text{mis}}$ , where  $\mathbf{a}_{-j}$ , denoting the components of  $\mathbf{a}$  with  $\mathbf{a}_j$  removed. Because of the simplicity of its covariate distribution specification, MICE is popularly used to deal with missing data.

In spite of the widespread use of MICE-FCS due to its convenience, it was recently proved that the asymptotic distribution drawn using MICE-FCS is not equivalent to the existing Bayesian simulation in several settings. Liu et al. (2014) showed that the MICE-FCS algorithm does not guarantee that the asymptotic distributions are consistent with the existing Bayesian joint model MI estimator when the family of conditional models and their joint distributions are incompatible. According to Liu et al. (2014), the compatibility assumption is satisfied when a parameter set of conditional models  $f_j(a_j^{\text{mis}}|\mathbf{a}_{-j}^{\text{mis}}, \mathbf{a}^{\text{obs}}, \boldsymbol{\psi}_j)$  is represented by surjective mapping of a collection of the joint model  $p(\mathbf{a}^{\text{mis}}, \mathbf{a}^{\text{obs}}|\boldsymbol{\psi})$  parameter  $\boldsymbol{\psi}$ , that is,  $g_j(\boldsymbol{\psi}) = \boldsymbol{\psi}_j$ , and hence,  $p(a_j^{\text{mis}}|\mathbf{a}_{-j}^{\text{mis}}, \mathbf{a}^{\text{obs}}, \boldsymbol{\psi}) = f_j(a_j^{\text{mis}}|\mathbf{a}_{-j}^{\text{mis}}, \mathbf{a}^{\text{obs}}, \boldsymbol{\psi}_j)$ ; otherwise, they are said to be incompatible. Put simply, compatibility purports that the parameters of the conditional distribution can be expressed by the parameters of the joint distribution of the model. Liu et al. (2014) also showed that the MICE-FCS algorithm generates a consistent estimator using Rubin's rules, but the variance of the parameters cannot be applied to Rubin's rules if the family of the conditional models is semicompatible as a special case of incompatibility. On the other hand, if the model is compatible, MICE-FCS is asymptotically equivalent to the existing Bayesian simulation; hence, one can apply Rubin's rule to calculate the mean and variance of the parameters of interest.

In what kinds of cases this compatibility assumption holds? If the all variables in the datasets are consist only of continuous variables that follow an i.i.d. multivariate normal distribution and the conditional models are linear regressions with normally distributed error terms, the conditionals and joint model are compatible, and the estimators of MICE-FCS are applicable to Rubin's rules. If one variable is binary variable and the rest are continuous, one can also apply Rubin's rules to the substantive model described in the form of a linear regression with normally distributed error terms. However, when the researcher is interested in binary outcome modeling with a logistic regression, wherein there exist binary covariates in the datasets, and even if all the other covariates are continuous, the conditionals and substantive model (logistic model specification) are incompatible, and the MICE-FCS estimators are not equivalent to those corresponding to Gibbs sampling. In epidemiological and clinical research, researchers often assume nonlinear models such as the Cox proportional hazards model, regression models with quadratic terms, or regression models with interaction terms. Yet, unfortunately, these are examples of model incompatibility. In addition, the conditionals of MICE-FCS are under the immediate control of the researcher, and hence, the joint distribution is only implicitly known and may not

exist (van Buuren 2012). Therefore, although MICE-FCS is simple and convenient, the estimators are valid in a very limited number of cases only.

Bartlett et al. (2015) recently developed SMC-FCS in order to relax the compatibility assumption, which assigns the imputer compatibility of the joint distribution of covariates only and not all the variables. However, it has hardly solved the problem of the compatibility assumption holding because the number of covariates is generally larger than the outcome variable. For example, if the missing covariates contain two or more binary variables, a parameter set of conditional models cannot be represented by onto mapping of a collection of joint model parameters. Hence, the model compatibility assumption among covariates is violated, and the estimates from the SMC-FCS do not guarantee the consistency. Additionally, if some covariates do not follow a normal or Bernoulli distribution (say, they follow a log-normal distribution or a mixture of normal distributions), the compatibility assumption of joint distribution among covariates required in the SMC-FCS method cannot be generally satisfied, as seen in the motivating example in Sect. 1.1.

### 3 Semiparametric Bayes multiple imputation

#### 3.1 Semiparametric model formulation

As stated in Eq. (1), we propose a SB-MI algorithm, expressed by imputation model  $p(\mathbf{y}, \mathbf{w}|\mathbf{v})$  as the product of two submodels  $p(\mathbf{y}|\mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)p(\mathbf{w}|\mathbf{v}, \boldsymbol{\vartheta}_m)$ . Assuming the independence of the priors  $p(\boldsymbol{\vartheta}_m)$  and  $p(\boldsymbol{\vartheta}_s)$ , the posteriors are

$$p(\boldsymbol{\vartheta}_m, \boldsymbol{\vartheta}_s | \mathbf{y}, \mathbf{v}, \mathbf{w}) \propto p(\boldsymbol{\vartheta}_m)p(\boldsymbol{\vartheta}_s)p(\mathbf{y}|\mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)p(\mathbf{w}|\mathbf{v}, \boldsymbol{\vartheta}_m)$$

$\boldsymbol{\vartheta}_m$  and  $\boldsymbol{\vartheta}_s$  can be drawn from  $p(\boldsymbol{\vartheta}_m)p(\mathbf{w}|\mathbf{v}, \boldsymbol{\vartheta}_m)$  and  $p(\boldsymbol{\vartheta}_s)p(\mathbf{y}|\mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)$ , respectively. Given these parameters, the missing values are drawn from the density proportional to  $p(\mathbf{y}, \mathbf{w}|\mathbf{v})$ .

The specification of the substantive model  $p(\mathbf{y}|\mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)$  varies by the purpose of the analysis and the properties of the outcome  $\mathbf{y}$ . One may specify the linear regression to the continuous outcome or the logistic regression to the discrete outcome. Besides, one must employ specified model forms on  $p(\mathbf{y}|\mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)$ , such as the Cox proportional hazards model or quadratic models. Our proposed methodology, in any case, can properly estimate  $\boldsymbol{\vartheta}_s$  even when MICE-FCS or SMC-FCS cannot because of model incompatibility.

#### 3.2 Covariate distribution

On the other hand, we have to specify the complicated covariate distribution  $p(\mathbf{w}|\mathbf{v}, \boldsymbol{\vartheta}_m)$ . Usually,  $\mathbf{w}$  will include continuous and discrete variables. In order to deal with mixed-scale covariates, we employ a transformation of the latent continuous variables for probit modeling. Thus, we rewrite  $\mathbf{w}$  as  $\mathbf{w} = (\mathbf{w}_c, \mathbf{w}_d, \mathbf{w}_n)$  where  $\mathbf{w}_c$  denotes the continuous variable component,  $\mathbf{w}_d$  denotes the ordered variable com-



ponent with  $Q_d (= 1, \dots, q_d, \dots, Q_d)$  orders, and  $\mathbf{w}_n$  denotes the nominal variable component with  $Q_n (= 0, \dots, q_n, \dots, Q_n)$  choices. To deal with discrete variables simply, we introduce continuous latent variables  $\mathbf{u}_d$  and  $\mathbf{u}_n$  where

$$w_d = q_d \quad \text{if } \zeta_{q_d-1} < u_d \leq \zeta_{q_d}$$

$$w_n = \begin{cases} 0 & \text{if } \max(u_n) < 0 \\ q_n & \text{if } \max(u_n) = u_{nq_n} > 0 \end{cases}$$

We assume the following structure on the covariate distribution.

$$\mathbf{w}_i^* = f(\mathbf{v}_i) + \boldsymbol{\eta}_i$$

with  $\mathbf{w}_i = g(\mathbf{w}_i^*)$ , where  $f$  is an unknown function,  $\boldsymbol{\eta}_i$  is  $q$ -dimensional i.i.d. random errors with  $E(\boldsymbol{\eta}_i) = \mathbf{0}$ ,  $\mathbf{w}^* = (\mathbf{w}_c, \mathbf{u}_d, \mathbf{u}_n)$ , and  $g(\cdot)$  represents the function converting the latent continuous variables  $\mathbf{w}_i^*$  to  $\mathbf{w}_i$ . This enables us to deal with many types of continuous and discrete variables in the covariate distribution in a straightforward manner.

More concretely, we employ DPM modeling to represent the covariate distribution. DPM modeling is frequently utilized in applied statistical modeling when researchers intend to avoid making assumptions about parameter distribution within the Bayesian framework. For example, Hirano (2002) developed autoregressive models with individual effects where the disturbances are not restricted to a parametric class. Rodriuez et al. (2009) used DPM to develop a Bayesian semiparametric approach for functional data analysis. Kunihamu et al. (2016) developed a nonparametric Bayes model with DPM to incorporate sample survey weights. The theoretical properties of DPM were investigated by Shen et al. (2013).

According to Sethuraman (1994), the Dirichlet process as a prior for a random distribution  $G$  can be represented by the stick-breaking process. Let  $\xi_1, \xi_2, \dots$  be an independent draw from a beta distribution  $Be(1, \gamma)$ . If  $G$  follows the Dirichlet process prior with concentration parameter  $\gamma$  and base distribution  $G_0$ , that is,  $G \sim DP(\gamma, G_0)$ ,  $G$  can be represented as

$$G = \sum_{l=1}^{\infty} \kappa_l \delta_{\theta_l}, \quad \theta_l \sim G_0$$

where  $\kappa_l = \xi_l \prod_{h < l} (1 - \xi_h)$  and  $\delta_{\theta}$  is a point mass at  $\theta$  [refer to Walker et al. (1999) for a detailed description of DPM].

Although DPM is used to flexibly express a variety of parameters or distributions, they are greatly restricted because probability weight  $\kappa_l$  is a constant (Dunson et al. 2007). If the stick-breaking weights  $\pi_l$  are constant and independent of predictor  $\mathbf{x}_i$ , as in DPM and other nonparametric Bayesian models, the mean regression structure is reduced to a linear one, namely  $\sum_{l=1}^{\infty} \pi_l \boldsymbol{\beta}_l^T \mathbf{x}_i \approx \bar{\boldsymbol{\beta}}^T \mathbf{x}_i$ , where  $\bar{\boldsymbol{\beta}} = (\sum_{l=1}^{\infty} \pi_l \boldsymbol{\beta}_l^T)$ . Therefore, in this paper, we apply the PSBPM model proposed by Chung and Dunson (2009) since the algorithm allows for greater flexibility through predictor-dependent stick-breaking weights  $\pi_l(\mathbf{x}_i)$ . In addition, PSBPM results in a conjugate structure, and

hence simpler posterior calculation. The statistical properties of PSBPM are described in [Pati et al. \(2013\)](#). For example, [Hoshino \(2013\)](#) proposed a semiparametric Bayesian model for causal inference where the parameters of no interest for researchers are modeled using PSBPM.

We apply PSBPM modeling to the covariate distribution specification. The resulting regression function of  $\mathbf{w}^*$  on  $\mathbf{v}$  can be represented as

$$f(\mathbf{w}_i^* | \mathbf{v}_i) = \sum_{l=1}^{\infty} \pi_l(\mathbf{v}) N(\Gamma_l \mathbf{v}_i, \Phi_l)$$

with the probability weights

$$\pi_l(\mathbf{v}) = \Phi(\xi_l(\mathbf{v})) \prod_{h < l} \{1 - \Phi(\xi_h(\mathbf{v}))\},$$

where  $\Phi(\cdot)$  is the cumulative distribution function of the standard normal distribution. To make the probability weights  $\pi_l(\mathbf{v})$  vary with covariates  $\mathbf{v}$ , we let  $\xi_l(\mathbf{v}) = \alpha_l + f_l(\mathbf{v})$ ,  $\alpha_l \sim N(\mu_\alpha, 1)$ ,  $\mu_\alpha \sim N(\mu_{\alpha_0}, \sigma_{\alpha_0}^2)$  and we introduce the following regression function as in [Chung and Dunson \(2009\)](#):

$$f_l(\mathbf{v}) = - \sum_{k=1}^q \psi_{lk} |v_k - \Omega_{lk}|$$

with  $\psi_{lk} \sim N(\mu_{\psi_k}, \sigma_{\psi_k}^2) \mathbf{1}_{R^+}$ ,  $\Omega_{lk} \sim \sum_{m=1}^{M_k} \frac{1}{M_k} \delta_{\Omega_{km}^*}(\Omega_{lk})$ , where  $N(\mu_U, \sigma_U^2) \mathbf{1}_U$  denotes a normal distribution with mean  $\mu_U$  and variance  $\sigma_U^2$  truncated to the set  $U$  and  $\Omega_{km}^*$  are discrete points over a reasonable range of the  $k$ -th covariate  $v_k$ .

### 3.3 Imputation and analysis procedure

Data imputation and analysis procedures are as follows.

1. Impute missing component of  $\mathbf{y}$  and  $\mathbf{w}$  if missing as starting values.
2. Assign each case  $i$  to any class  $l$  of the Dirichlet mixture.
3. Generate  $\boldsymbol{\vartheta}_m^{(t)}$  from the posterior distribution based on the likelihood of  $p(\mathbf{w}^{(t-1)} | \mathbf{v}, \boldsymbol{\vartheta}_m)$  calculated on the complete case dataset or pseudo-complete dataset.
4. If  $\mathbf{w}$  is missing, generate the missing component of  $\mathbf{w}$  proportional to  $p(\mathbf{y}^{(t-1)}, \mathbf{w} | \mathbf{v}, \boldsymbol{\vartheta}_s^{(t-1)}, \boldsymbol{\vartheta}_m^{(t)})$ .
5. If  $\mathbf{y}$  is missing, generate the missing component of  $\mathbf{y}$  from  $p(\mathbf{y} | \mathbf{w}^{(t)}, \mathbf{v}, \boldsymbol{\vartheta}_s^{(t-1)})$ . If case  $i$  has missingness both on  $\mathbf{w}_i$  and  $\mathbf{y}_i$ , generate  $\mathbf{y}_i$  from  $p(\mathbf{y} | \mathbf{w}^{(t)}, \mathbf{v}, \boldsymbol{\vartheta}_s^{(t-1)})$ , where missing  $\mathbf{w}_i^{(t)}$  is imputed in step 3.
6. Given the pseudo-complete dataset, generate  $\boldsymbol{\vartheta}_s^{(t)}$  from the posterior distribution based on the likelihood of  $p(\mathbf{y}^{(t)} | \mathbf{w}^{(t)}, \mathbf{v}, \boldsymbol{\vartheta}_s)$ .

Steps 2–6 are repeated for  $t = 1, 2, \dots$  until convergence. Then, the sequence of  $\vartheta_s^{(t)}$  obtained in Step 6 is used to the posterior inference of the substantive model.

The starting values in Step 1 can be imputed based on single imputation such as mean imputation or regression imputation using complete case.

### 3.4 Posterior computation of the proposed model

Let  $\mathbf{K}_i$  be the indicator denoting where case  $i$  belongs, and  $\mathbf{K}_i = l$  if case  $i$  belongs to class  $l$ . Recall that  $\vartheta_m$  and  $\vartheta_s$  are the parameter vectors for the covariate distribution  $p(\mathbf{w}|\mathbf{v})$  and the substantive model  $p(\mathbf{y}|\mathbf{w}, \mathbf{v})$ , respectively. This yields the following hierarchical representation of the finite-dimensional PSBPM model:

$$\begin{aligned} \mathbf{y}_i | \mathbf{w}_i, \mathbf{v}_i, \vartheta_s &\sim p(\mathbf{y}_i | \mathbf{w}_i, \mathbf{v}_i, \vartheta_s), \\ \mathbf{w}_i | \mathbf{v}_i, \vartheta_m, \mathbf{K}_i &\sim p(\mathbf{w}_i | \mathbf{v}_i, \vartheta_m, \mathbf{K}_i), \\ \mathbf{K} | \phi &\sim \sum_{l=1}^{\infty} \pi_l(\mathbf{v}_i | \phi_l) \delta_l(\cdot) \quad (i = 1, \dots, N), \\ \phi^* &\sim p(\phi^* | \tau_\phi), \\ \vartheta_s &\sim p(\vartheta_s | \tau_{\vartheta_s}), \quad \vartheta_m \sim p(\vartheta_m | \tau_{\vartheta_m}), \\ \tau &\sim p(\tau), \end{aligned}$$

where  $\phi^* = (\alpha_l, \phi_{11}, \dots, \phi_{Lq}, \Omega_{11}, \dots, \Omega_{Lq})$  and  $\tau = (\tau_\phi^T, \tau_{\vartheta_s}^T, \tau_{\vartheta_m}^T)^T$ .

The blocked Gibbs sampler (Ishwaran and James 2001) is applied to the posterior computation of the PSBPM parameters  $\vartheta_m$ . The blocked Gibbs sampler is very similar to the Gibbs sampler except for the assignment of samples to each class. Since we employ PSBPM modeling, we can directly apply the estimation algorithm of Chung and Dunson (2009) for the simulation of  $\vartheta_m$ . Each case  $i$  is assigned to one of the  $L$  classes in the blocked Gibbs sampling, where  $L$  denotes the maximum number or truncation of classes. As stated above,  $\mathbf{w}$  may include continuous and discrete variables. We employ the transformation of the latent continuous variables  $\mathbf{w}^*$  for probit modeling through a function  $g$  such that  $\mathbf{w} = g(\mathbf{w}^*)$ . Given the draw of  $\vartheta_m$ , the missing components of  $\mathbf{y}$  and  $\mathbf{w}$  are imputed, and then, the substantive model parameter  $\vartheta_s$  is simulated. We obtain the detailed posterior computation using the MCMC estimation as follows.

1. Update  $K_i$  ( $i = 1, \dots, N$ )

To assign samples to each class, generate  $K_i$  by  $\sum_{l=1}^L \pi_{li} \delta_l(\cdot)$ , where  $\pi_{li}$  is

$$\pi_{li} = \frac{\pi_l(\mathbf{v}_i) N(\mathbf{\Gamma}_l \mathbf{v}_i, \mathbf{\Phi}_l)}{\sum_{l=1}^L \pi_l(\mathbf{v}_i) N(\mathbf{\Gamma}_l \mathbf{v}_i, \mathbf{\Phi}_l)}$$

with  $\pi_l(\mathbf{v}_i) = \Phi(\xi_l(\mathbf{v}_i)) \prod_{h < l} \{1 - \Phi(\xi_h(\mathbf{v}_i))\}$ .

2. Update  $Z_{il}^*$

We introduce latent variable  $Z_{il}^*$  where  $Z_{il} = \mathbf{1}(Z_{il}^* > 0)$  and

$$Z_{il}^* \sim \begin{cases} N\left(\alpha_l - \sum_{k=1}^q \psi_{lk} |v_{ik} - \Omega_{lk}|, 1\right) \mathbf{1}_{R^+} & \text{for } l = K_i \\ N\left(\alpha_l - \sum_{k=1}^q \psi_{lk} |v_{ik} - \Omega_{lk}|, 1\right) \mathbf{1}_{R^-} & \text{for } l < K_i. \end{cases}$$

3. Update  $\alpha_l (l = 1, \dots, L - 1)$

Draw  $\alpha_l$  from the following normal distribution.

$$\alpha_l \sim N\left(\frac{\sum_{i:K_i \geq l} W_{il}^* + \mu_\nu}{N_l + 1}, \frac{1}{N_l + 1}\right),$$

where  $N_l = \sum_{i=1}^N \mathbf{1}(K_i \geq l)$  and  $W_{il}^* = Z_{il}^* + \sum_{k=1}^q \psi_{lk} |v_{ik} - \Omega_{lk}|$ .

4. Update  $\psi_{lk} (l = 1, \dots, L - 1, k = 1, \dots, q)$

Draw  $\psi_{lk}$  from the following left-truncated normal distribution.

$$\psi_{lk} \sim N\left(\frac{\sigma_{\psi_k}^2 \mu_{\psi_k} + \sum_{i:K_i \geq l} |v_{ik} - \Omega_{lk}| U_{il}^*}{\sigma_{\psi_k}^2 + \sum_{i:K_i \geq l} |v_{ik} - \Omega_{lk}|^2}, \frac{1}{\sigma_{\psi_k}^2 + \sum_{i:K_i \geq l} |v_{ik} - \Omega_{lk}|^2}\right) \mathbf{1}_{R^+},$$

where  $U_{il}^* = \alpha_l - Z_{il}^* - \sum_{s=1, s \neq k}^q \psi_{ls} |v_{is} - \Omega_{ls}|$ .

5. Update  $\Omega_{lk} (l = 1, \dots, L - 1, k = 1, \dots, q)$

Draw  $\Omega_{lk}$  from the following probability.

$$\begin{aligned} & \Pr(\Omega_{lk} = \Omega_{km}^*) \\ &= \frac{\frac{1}{M_k} \prod_{i:K_i \geq l} N\left(Z_{il}^*; \alpha_l - \sum_{s=1, s \neq k}^q \psi_{ls} |v_{is} - \Omega_{ls}| - \psi_{lk} |v_{ik} - \Omega_{km}^*|, 1\right)}{\sum_{m=1}^{M_k} \frac{1}{M_k} \prod_{i:K_i \geq l} N\left(Z_{il}^*; \alpha_l - \sum_{s=1, s \neq k}^q \psi_{ls} |v_{is} - \Omega_{ls}| - \psi_{lk} |v_{ik} - \Omega_{km}^*|, 1\right)}. \end{aligned}$$

6. Update  $\boldsymbol{\vartheta}_m (= \boldsymbol{\Gamma}_l, \boldsymbol{\Phi}_l)$

Draw  $\boldsymbol{\vartheta}_m (= \boldsymbol{\Gamma}_l, \boldsymbol{\Phi}_l)$  from the following multivariate normal and inverted Wishart distribution.

$$\begin{aligned} \boldsymbol{\Gamma}_l | rest &\sim N\left(\text{vec}(\widehat{\boldsymbol{\Gamma}}), \boldsymbol{\Phi}_l \otimes (\mathbf{V}_l^T \mathbf{V}_l)^{-1}\right), \\ \boldsymbol{\Phi}_l | rest &\sim IW\left(f_0 + N, \mathbf{G}_0^{-1} + (\mathbf{W}_l^* - \boldsymbol{\Gamma}_l \mathbf{V}_l)^T (\mathbf{W}_l^* - \boldsymbol{\Gamma}_l \mathbf{V}_l)\right), \end{aligned}$$

where  $\widehat{\boldsymbol{\Gamma}} = (\mathbf{V}_l^T \mathbf{V}_l)^{-1} \mathbf{V}_l^T \mathbf{W}_l$ ,  $\mathbf{V} = (\mathbf{v}_1^T, \dots, \mathbf{v}_N^T)^T$ ,  $\mathbf{W}^* = (\mathbf{w}_1^{*T}, \dots, \mathbf{w}_N^{*T})^T$ , and  $\mathbf{V}_l$  and  $\mathbf{W}_l^*$  denote the subset of  $\mathbf{V}$  and  $\mathbf{W}^*$  whose case  $i$  belong to class  $l$ .  $f_0$  and  $\mathbf{G}_0^{-1}$  denotes the parameter of the prior distribution of  $\boldsymbol{\Gamma}_l$ ;  $\boldsymbol{\Gamma}_l \sim IW(f_0, \mathbf{G}_0^{-1})$ .

7. Update  $\mu_\alpha$

Draw  $\mu_\alpha$  from the following normal distribution

$$\mu_\alpha \sim N \left( (L - 1 + \sigma_{\alpha_0}^{-2})^{-1} \left[ \sum_{l=1}^{L-1} \alpha_l + \sigma_{\alpha_0}^{-2} \mu_{\alpha_0} \right], (L - 1 + \sigma_{\alpha_0}^{-2})^{-1} \right).$$

8. Update the missing components

Draw the missing component of  $\mathbf{w}$  from a density proportional to  $p(\mathbf{y}, \mathbf{w} | \mathbf{v}, \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ . Since it is difficult to draw the missing  $\mathbf{w}$ , we employ the Metropolis–Hastings algorithm and use  $p(\mathbf{w} | \mathbf{v}, \boldsymbol{\vartheta}_m)$  as a proposal density in order to draw a candidate of  $\mathbf{w}_i, \mathbf{w}_i^c$ . Note that the candidates are obtained after the transformation  $\mathbf{w} = g(\mathbf{w}^*)$ . We accept the candidates with the following probability:

$$\min \left( \frac{p(\mathbf{y} | \mathbf{w}^c, \mathbf{v}, \boldsymbol{\vartheta}_s)}{p(\mathbf{y} | \mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)}, 1 \right). \tag{2}$$

If any component of  $\mathbf{y}$  is missing, we draw  $\mathbf{y}^{\text{mis}}$  from

$$p(\mathbf{y}^{\text{mis}} | \mathbf{w}, \mathbf{v}, \boldsymbol{\vartheta}_s)$$

9. Update  $\boldsymbol{\vartheta}_s$

Draw  $\boldsymbol{\vartheta}_s$  from the density proportional to

$$p(\boldsymbol{\vartheta}_s) \prod_{i=1}^N p(y_i | \mathbf{w}_i, \mathbf{v}_i, \boldsymbol{\vartheta}_s).$$

For example, if we are interested in inferring the binary logistic regression coefficients, the acceptance probability in Eq. (2) can be written as

$$\min \left( \frac{\left\{ \frac{\exp(\boldsymbol{\Gamma}^T \mathbf{x}_i^c)}{1 + \exp(\boldsymbol{\Gamma}^T \mathbf{x}_i^c)} \right\}^{y_i} \left\{ \frac{1}{1 + \exp(\boldsymbol{\Gamma}^T \mathbf{x}_i^c)} \right\}^{1-y_i}}{\left\{ \frac{\exp(\boldsymbol{\Gamma}^T \mathbf{x}_i)}{1 + \exp(\boldsymbol{\Gamma}^T \mathbf{x}_i)} \right\}^{y_i} \left\{ \frac{1}{1 + \exp(\boldsymbol{\Gamma}^T \mathbf{x}_i)} \right\}^{1-y_i}}, 1 \right)$$

where  $\boldsymbol{\Gamma}$  is a vector of coefficients,  $\mathbf{x}_i = (\mathbf{w}_i^T, \mathbf{v}_i^T)^T$ , and  $\mathbf{x}_i^c$  is the vector  $\mathbf{x}_i$  whose missing components are replaced by the candidate value.

In addition to the above steps, if any  $\mathbf{w}$  is ordered variable, cutting points  $\zeta_{q_d}$  must be estimated. We can employ Bayesian-ordered probit MCMC method such as [Albert and Chib \(1993\)](#) or [Albert and Chib \(2001\)](#) before Step 8. Furthermore, we have to make some restriction on  $\boldsymbol{\Phi}_l$  when  $\mathbf{w}$  includes categorical variables. The diagonal elements of  $\boldsymbol{\Phi}_l$  corresponding to  $u_d$  and one of  $u_{nq_n}$  (e.g.,  $u_{nq_N}$ ) are restricted to be 1. We can employ many kinds of methods to generate  $\boldsymbol{\Phi}_l$  such as [McCulloch and Rossi \(1994\)](#) or [Zhang et al. \(2008\)](#).

## 4 Simulation study

We conduct the following four simulation studies in order to illustrate the performance of the proposed method when MICE-FCS cannot draw from a Bayesian joint model: (i) linear regression with a quadratic term, (ii) linear regression with an interaction term, (iii) proportional hazards model with a binary covariate, and (iv) logistic regression with a binary covariate. Through the simulation study, we consider the case  $N = 400$ , and 30% of the incomplete covariates are set to be missing. We generate the missing values based on  $p(\mathbf{r}|\mathbf{y}, \mathbf{w}, \mathbf{v}) = p(\mathbf{r}|\mathbf{y})$ . We compare the following imputation methods with SB-MI: MICE-FCS, SMC-FCS, NP-MI, and missForest. Stekhoven and Bühlmann (2012) proposed the missForest algorithm, which imputes missing values from the random forest predictors and they are reported to provide lower imputation errors than the FCS method (Liao et al. 2014; Waljee et al. 2013). The detailed simulation design, results, and detailed discussion appear in Appendix of ESM.

This section summarizes the results of simulation study (iv). We specify the substantive model of the logistic function as follows:

$$\text{logit}(y = 1) = \Gamma_0 + \Gamma_1 w_1 + \Gamma_2 w_2 + \Gamma_3 v$$

with  $\Gamma_0 = 1$ ,  $\Gamma_1 = 2$ ,  $\Gamma_2 = -2$ , and  $\Gamma_3 = 3$ . In this simulation, we consider three scenarios where the missing covariates follow (a) a multivariate normal distribution, (b) a multivariate log-normal distribution, and (c) a multivariate normal mixture distribution. We consider a case where one of the incompletely observed covariates  $w_1$  is binary, where  $w_{i,1} = 1$  if the latent variable, which is simulated based on the above three process,  $w_{i,1}^* > 0$  and  $w_{i,1} = 0$  if the latent variable  $w_{i,1}^* \leq 0$ . The detailed data generating process is described in Appendix of ESM.

Table 1 describes the results of the simulation, including the empirical mean, standard deviation, the coverage of nominal 95% confidence intervals (CIs) of the estimate, and the mean squared error (MSE) from the true value of  $\Gamma$ . The last row for each scenario shows each sum of the MSE ratio for MICE-FCS. With Scenario (a), namely the normally distributed covariates, SB-MI gives the most accurate estimates, with an empirical CI coverage of approximately 0.95. However, MICE-FCS, SMC-FCS, NP-MI, and missForest result in biased estimates, and their CI coverage are also considerably poor. With Scenario (b), namely the log-normally distributed missing covariates, MICE-FCS and SMC-FCS continue to be biased with incorrect empirical CI coverage. missForest provides relatively good estimates in terms of MSE, but, once again, the CIs are poor. SB-MI gives relatively correct estimates, with the CI coverage closer to 0.95 compared with MICE-FCS, SMC-FCS, NP-MI, and missForest. In Scenario (c), namely the mixture of normally distributed covariates, MICE-FCS, SMC-FCS, NP-MI, and missForest once again result in biased estimates, and end with poor empirical CI coverage for some  $\Gamma$ s. On the other hand, MSE of SB-MI provides the best result among the three, and the CI coverage are very close to 0.95 for all  $\Gamma$ s. For all three scenarios, MICE-FCS and NP-MI provide very similar and biased results since both are thought to be uncongenial multiple imputation method.

**Table 1** Empirical mean, standard deviation, coverage of nominal 95% CIs, and mean squared error of the estimates, and sum of MSE Ratio of MICE-FCS basis from 1000 simulations are described

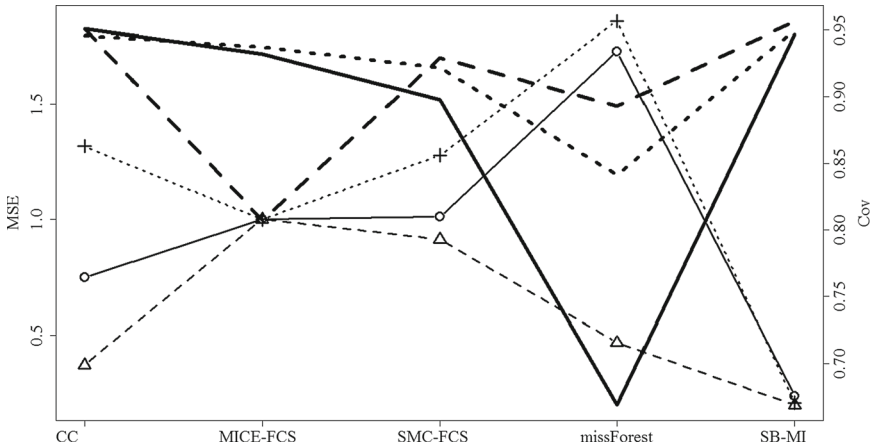
True	Complete case				MICE-FCS				SMC-FCS			
	Mean	(SD)	Cov	MSE	Mean	(SD)	Cov	MSE	Mean	(SD)	Cov	MSE
<b>(a) Normal</b>												
$\Gamma_0 = 1$	1.023	(0.286)	0.950	0.086	0.767	(0.257)	0.830	0.313	0.721	(0.260)	0.793	0.243
$\Gamma_1 = 2$	2.058	(0.517)	0.955	0.283	1.929	(0.485)	0.959	0.286	2.180	(0.502)	0.943	0.300
$\Gamma_2 = -2$	-2.081	(0.391)	0.946	0.165	-1.944	(0.355)	0.960	0.212	-2.161	(0.368)	0.944	0.169
$\Gamma_3 = 3$	3.127	(0.500)	0.952	0.265	3.118	(0.412)	0.979	0.254	3.388	(0.443)	0.911	0.365
MSE ratio	0.750				1.000				1.011			
<b>(b) Log-normal</b>												
$\Gamma_0 = 1$	1.026	(0.312)	0.950	0.098	0.554	(0.302)	0.673	0.492	0.873	(0.309)	0.917	0.259
$\Gamma_1 = 2$	2.055	(0.401)	0.954	0.167	1.748	(0.385)	0.885	0.350	2.195	(0.397)	0.960	0.293
$\Gamma_2 = -2$	-2.056	(0.297)	0.952	0.091	-1.706	(0.276)	0.774	0.275	-2.190	(0.287)	0.928	0.297
$\Gamma_3 = 3$	3.090	(0.386)	0.947	0.157	2.794	(0.353)	0.896	0.264	3.356	(0.370)	0.913	0.414
MSE ratio	0.371				1.000				0.914			
<b>(c) Mixture of normals</b>												
$\Gamma_0 = 1$	1.035	(0.487)	0.951	0.254	0.805	(0.324)	0.884	0.281	1.021	(0.471)	0.942	0.251
$\Gamma_1 = 2$	2.143	(0.886)	0.935	0.915	2.143	(0.582)	0.962	0.395	1.983	(0.851)	0.943	0.794
$\Gamma_2 = -2$	-2.105	(0.404)	0.947	0.192	-2.125	(0.396)	0.969	0.266	-2.137	(0.380)	0.912	0.181
$\Gamma_3 = 3$	3.151	(0.515)	0.950	0.312	3.379	(0.451)	0.934	0.329	3.388	(0.465)	0.890	0.398
MSE ratio	1.318				1.000				1.279			

**Table 1** continued

True	NP-MI			missForest			SB-MI		
	Mean	(SD)	MSE	Mean	(SD)	MSE	Mean	(SD)	MSE
<b>(a) Normal</b>									
$\Gamma_0 = 1$	0.775	(0.260)	0.127	0.600	(0.240)	0.235	1.027	(0.223)	0.053
$\Gamma_1 = 2$	1.938	(0.493)	0.227	2.529	(0.449)	0.664	1.965	(0.414)	0.066
$\Gamma_2 = -2$	-1.991	(0.365)	0.124	-2.452	(0.349)	0.413	-1.989	(0.294)	0.056
$\Gamma_3 = 3$	3.193	(0.426)	0.210	3.542	(0.410)	0.530	2.985	(0.353)	0.078
MSE ratio	0.645			1.728			0.237		
<b>(b) Log-normal</b>									
$\Gamma_0 = 1$	0.782	(0.339)	0.144	1.070	(0.299)	0.1105	1.024	(0.259)	0.062
$\Gamma_1 = 2$	1.825	(0.408)	0.174	2.224	(0.376)	0.218	1.980	(0.341)	0.065
$\Gamma_2 = -2$	-1.900	(0.342)	0.091	-2.268	(0.276)	0.156	-2.006	(0.228)	0.066
$\Gamma_3 = 3$	2.855	(0.384)	0.147	3.194	(0.345)	0.167	3.019	(0.302)	0.081
MSE ratio	0.403			0.467			0.198		
<b>(c) Mixture of normals</b>									
$\Gamma_0 = 1$	0.767	(0.323)	0.155	0.681	(0.442)	0.326	1.023	(0.235)	0.046
$\Gamma_1 = 2$	2.221	(0.595)	0.399	2.675	(0.826)	1.319	1.978	(0.439)	0.082
$\Gamma_2 = -2$	-2.118	(0.415)	0.175	-2.384	(0.357)	0.321	-1.989	(0.299)	0.069
$\Gamma_3 = 3$	3.323	(0.465)	0.307	3.393	(0.431)	0.393	2.997	(0.354)	0.067
MSE ratio	0.816			1.859			0.208		

CC, complete case analysis; MICE-FCS, multiple imputation by chained equation-fully conditional specification; NP-MI, nonparametric multiple imputation; missForest, random forest approach for missing value prediction; SMC-FCS, substantive model compatible-fully conditional specification; SB-MI, semiparametric Bayesian multiple imputation (Proposed)





**Fig. 2** The thinner lines with markers and the thicker lines correspond to the left vertical axis (MSE) and the right vertical axis (Cov), respectively. The solid, dashed, and dotted lines represent the results of the following scenarios, respectively: (a) multivariate normal distribution, (b) multivariate log-normal distribution, and (c) multivariate normal mixture distribution. CC, complete case analysis; MICE-FCS, multiple imputation by chained equation—fully conditional specification; SMC-FCS, substantive model compatible—fully conditional specification; SB-MI, semiparametric Bayes multiple imputation (proposed method)

Figure 2 compares the results of MSE and the coverage of nominal 95% CIs. We observe that the proposed method gives the smallest MSE and the best coverage.

For the other three simulations, the summarized main results are as follows. First, in terms of the MSE, the proposed method gives estimates equivalent to those found with SMC-FCS when the latter gives consistent estimates, but MICE-FCS results in biased estimates because of the violation of the model compatibility assumption. NP-MI also gives biased results which are very similar to MICE-FCS. This indicates that nonparametric imputation stage is uncongenial to the analysis models that we assumed. The coverage of nominal 95% CIs for the proposed method is very close to that of SMC-FCS. Note that missForest shows unbiased estimates in some situations, but produces underestimated standard deviations, and hence, poor CIs. This indicates that missForest is occasionally good at inferring unbiased estimates, but it should not be applied in fields such as medical or epidemiological research, where the results of statistical significance (or hypothesis testing) are crucial.

Second, the proposed method shows smaller MSEs when the model compatibility assumptions of FCS approaches (MICE-FCS and SMC-FCS) are not satisfied. Our simulation study includes the missing covariate that follows a log-normal or mixture of normal distribution. Although MICE-FCS and SMC-FCS result in larger MSEs in these simulation settings, our proposed method gives considerably smaller MSEs. Even in these situations, the coverage of the CIs is better compared to that under the imputation methods. NP-MI (as well as MICE-FCS) also provides very large MSE compared with the proposed since the analysis model such as simple proportional hazards model is thought to be uncongenial to the imputation model.

These results indicate that the proposed method can deal with more complicated covariate distributions that the researcher cannot prespecify. Therefore, these results suggest that SB-MI approach is very practical for treating missing datasets in the real world.

## 5 Real data analysis

In this section, we apply our proposed algorithm to the real dataset with missing components. The data used in this implementation are sourced from the ADNI dataset described in the motivating example in Sect. 1.1. The substantive model in this example is the Cox proportional hazards model, which helps us study the time to conversion to AD. The samples comprise 382 observations of MCI in baseline subjects who had at least one follow-up after the first diagnosis. Of these subjects, 167 participants converted to AD during the data period. The dataset contains missing covariates  $ABETA_{1-42}$ , the square of  $ABETA_{1-42}$ , tau, p-tau, the dummy variable of whether or not the subject's mother had AD, and the dummy variable of whether or not the subject's father had AD. The dataset also contains the completely observed dummy variable APOE4, which equals 1 if the subject has the APOE4 gene, and 0 otherwise. Jack et al. (2010) found evidence that  $ABETA_{1-42}$  is positively associated, p-tau is positively associated after controlling the effect of tau, and APOE4 is positively associated with the hazard of converting to AD. Bartlett et al. (2015) showed that contrary to expectations, "mother had AD" and "father had AD" are negatively associated with the hazard of converting to AD. We should note that  $ABETA_{1-42}$  of 190 observations, tau of 193 observations, p-tau of 189 observations, "mother had AD" of 77 observations, and "father had AD" of 93 observations are missing.

We employ a gamma process prior to the cumulative baseline hazard proposed by Kalbfleisch (1978), namely  $H_0 \sim GP(c_0 H^*, c_0)$ . We specify the hyperparameters to be  $c_0 = 0.01$ , and  $H^*$  follows an exponential distribution with parameter  $\lambda$ . Note that Kalbfleisch (1978) and Sinha et al. (2003) showed that the estimates from the gamma process prior to the cumulative baseline hazard are equivalent to the non-Bayesian estimates based on Cox's partial likelihood (Cox 1975) when  $c_0 \rightarrow 0$  [see Chen et al. (2006) for a detailed description of the MCMC method for a Bayesian Cox regression]. We draw 25,000 MCMC iterations after 25,000 burn-in phases. We confirmed the convergence using a diagnostic proposed by Geweke (1992).

Table 2 shows the results of the coefficients estimated using the Cox proportional hazards model. A positive coefficient indicates that the variable is associated with hazard of the subject converting to AD. We compare these results with those obtained using the other imputation methods, namely MICE-FCS, SMC-FCS, NP-MI, and missForest. Note that the results of MICE-FCS and SMC-FCS can be biased because the model compatibility assumption is violated, as noted in Sect. 2. NP-MI, as well as MICE-FCS, can also be biased due to the uncongeniality as indicated from the simulation. The estimated results of our proposed method (SB-MI) are different from those of MICE-FCS, SMC-FCS, NP-MI, missForest, as well as the complete case analysis in some ways. The results of SB-MI suggest that the effect of increasing  $ABETA_{1-42}$  to the hazard of conversion is nonlinear as in Bartlett et al. (2015),

whereas the coefficient of  $ABETA_{1-42}$  is not statistically significant for the estimates of the complete case samples, MICE-FCS, SMC-FCS, NP-MI, and missForest at the 5% level. The coefficients of  $ABETA_{1-42}$  and  $(ABETA_{1-42})^2$  from MICE-FCS are statistically significant at the 10% and 5% level, respectively, but they show signs opposite to those of SB-MI. In addition, the estimated coefficient of p-tau from SB-MI is much larger and closer to that of the complete case analysis compared to the other methods. Like the other existing methods, the hazard of the presence of parents who had AD is not statistically significant. The presence of the APOE4 gene, which is suspected to be associated with the development of AD, is positively associated with the hazard of converting to AD for SB-MI. This relationship cannot be found when we use complete case samples only. In addition, compared with the complete case analysis, several coefficients from SB-MI are statistically significant because they avoid the restrictions posed by the complete case sample. We conduct a logistic regression and find that data availability of each biomarker is not related to time to conversion to AD at the 1% level, which suggests that the missing data do not depend on the outcome, and assuming unbiased results from the complete case analysis is reasonable. However, the analysis based on the restricted sample fails to detect statistically significant variables related to the hazard of converting to AD.

On the other hand, the results from SB-MI are consistent with the findings of the previous studies, such as [Hansson et al. \(2006\)](#) and [Okello et al. \(2009\)](#), in terms of their signs. These studies employed separate longitudinal datasets and not the ADNI dataset. In conclusion, our method is very practical when applied to a real dataset, which often contains non-normally distributed and mixed-scale missing variables.

## 6 Discussion and conclusions

In this study, we proposed a SB-MI approach for regression models with missing mixed continuous and discrete covariates, in which the substantive model of the researcher's interest is a parametric formulation, and the covariate distributions are nonparametric formulations employing PSBPM modeling.

If the covariate model can be correctly specified, the EM algorithm or Bayesian MCMC approach can estimate unbiased results. However, prespecifying a covariate distribution is generally impossible, especially in cases where the missing variables are continuous and discrete. The FCS approach including MICE and SMC has been widely used because researchers are not required to specify the covariate distribution among mixed-scale variables. However, these methods yield severely biased estimates if the compatibility assumption of the model is violated.

On the other hand, the SB-MI framework, which is proposed in this paper, is capable of easily dealing with incompletely observed mixed-scale variables in the covariate distribution without using FCS. Therefore, we do not have to consider whether the compatibility assumption holds, and thus, we can assume a nonlinear regression, such as a linear regression with quadratic terms, Cox proportional hazards model, or logistic regression model on the substantive model even when the variables include discrete and non-normal continuous variables. The simulation studies show that the proposed method gives the best estimator in terms of MSE in cases where MICE-FCS and

**Table 2** Estimated results of Cox proportional hazard model

	CC ( <i>n</i> = 139)		MICE-FCS ( <i>n</i> = 382)		SMC-FCS ( <i>n</i> = 382)		NP-MI ( <i>n</i> = 382)		missForest ( <i>n</i> = 382)		SB-MI ( <i>n</i> = 382)	
	Coef.	(SD)	Coef.	(SD)	Coef.	(SD)	Coef.	(SD)	Coef.	(SD)	Coef.	(SD)
ABFTA <sub>1-42</sub> (ng/mL)	0.228	(0.178)	-0.048	(0.029)*	0.203	(0.162)	-0.024	(0.018)	-0.012	(0.128)	0.191	(0.092)**
(ABHTA <sub>1-42</sub> ) <sup>2</sup> (ng <sup>2</sup> /ml <sup>2</sup> )	-0.009	(0.005)*	0.001	(0.000)**	-0.008	(0.005)*	-0.027	(0.028)	-0.002	(0.004)	-0.009	(0.005)*
tau (pg/mL)	-0.028	(0.032)	-0.065	(0.020)***	-0.025	(0.027)	-0.002	(0.001)*	-0.023	(0.024)	-0.013	(0.011)
p-tau (pg/mL)	0.282	(0.105)***	0.241	(0.067)**	0.131	(0.076)*	0.168	(0.090)*	0.163	(0.087)*	0.250	(0.094)***
Mother had AD	-0.259	(0.293)	-0.187	(0.187)	-0.125	(0.194)	-0.224	(0.197)	-0.130	(0.188)	-0.098	(0.214)
Father had AD	-0.194	(0.485)	0.028	(0.241)	-0.126	(0.285)	-0.177	(0.290)	-0.202	(0.301)	-0.053	(0.343)
APOE4 positive	0.001	(0.289)	0.642	(0.176)***	0.333	(0.195)*	0.068	(0.194)	0.455	(0.177)**	0.455	(0.178)***

CC, complete case analysis; MICE-FCS, multiple imputation by chained equation—fully conditional specification; SMC-FCS, substantive model compatible—fully conditional specification; NP-MI, nonparametric multiple imputation; missForest, random forest approach for missing value prediction; SB-MI, semiparametric Bayesian multiple imputation (proposed)

\*, \*\*, and \*\*\* represent statistical significance at 10%, 5%, and 1% level, respectively

SMC-FCS result in biased estimates due to the violation of the model compatibility assumption. Furthermore, the proposed model is more robust when the distributions of the missing variables are non-normal. Since it is usual that some variables do not follow the normal distribution, and hence, the compatibility assumption is not satisfied, the results suggest that the SB-MI approach is very practical. NP-MI is a method imputing missing values using nonparametric specification, but they suffer from uncongenial analysis model. Although missForest, which accommodates the random forest approach, can sometimes give estimates closer to the “true” value compared to the proposed and existing methods, it underestimates the variance of the estimates, resulting in poor CIs. Consequently, it should not be applied to fields where the results of statistical significance (or hypothesis testing) are concern to researchers. The results of the real data analysis show that our proposed method can provide new insights that cannot be obtained from existing statistically improper methods.

Further study is needed to improve the efficiency of the SB-MI algorithm. Since our inference is based on the MCMC algorithm, the computation time required to obtain valid estimates is higher than that of the existing imputation method. In addition, our model can be extended to missing not at random (MNAR) by adding a submodel of the missing mechanism to our semiparametric specification using PSBPM modeling. However, it is very difficult to correctly specify the missing mechanism even if we assume a nonparametric formulation. Consequently, we did not consider the missing mechanism to be MNAR.

Another direction is to consider the case of high-dimensional covariates. Nonparametric Bayesian regression with Dirichlet process attains good performance when the number of the variables is large for its sample size, in general. For example, [Hannah et al. \(2011\)](#) conducted Monte Carlo simulation on varying sample size and number of covariates, showing that the Dirichlet process mixture regression results in smaller MSE than the existing regression model such as OLS or Gaussian process. It is also the case with the application to imputation. [Si and Reiter \(2013\)](#), which proposed nonparametric multiple imputation for categorical missing covariates, conducted simulation studies under the number of the variables 50 and  $N = 1000$ . They showed that nonparametric imputation using Dirichlet process mixture works well under such high dimensions.

Recently, nonparametric Bayes model including DPM for much more higher dimension or sparse datasets has been proposed. For example, [Tokdar et al. \(2010\)](#) and [Reich et al. \(2011\)](#) developed dimensionality reduction methods for Bayes nonparametric regression. [Li et al. \(2015\)](#) proposed nonparametric Bayes regression model for high dimensions using LASSO. Incorporating such models for missing data imputation may be useful for larger-dimensional dataset.

**Acknowledgements** This work was supported by JSPS KAKENHI Grant Numbers JP26285151, 18H03209, 16H02013, 16H06323. Data used in preparation of this article were obtained from the Alzheimer’s Disease Neuroimaging Initiative (ADNI) database ([adni.loni.usc.edu](http://adni.loni.usc.edu); <http://adni.loni.usc.edu>.) As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in analysis or writing of this report. A complete listing of ADNI investigators can be found at: [http://adni.loni.usc.edu/wp-content/uploads/how\\_to\\_apply/ADNI\\_Acknowledgement\\_List.pdf](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf).

## References

- Albert, J., Chib, S. (1993). Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, *88*, 669–679.
- Albert, J., Chib, S. (2001). Sequential ordinal modeling with applications to survival data. *Biometrics*, *57*, 829–836.
- Bartlett, J. W., Seaman, S. R., White, I. R., Carpenter, J. R. (2015). Multiple imputation of covariates by fully conditional specification: Accommodating the substantive model. *Statistics in Medicine*, *24*, 462–487.
- Canale, A., Dunson, D. B. (2011). Bayesian kernel mixtures for counts. *Journal of the American Statistical Association*, *106*, 1528–1539.
- Chen, M., Ibrahim, J. G., Shao, Q. (2006). Posterior propriety and computation for the Cox regression model with application to missing covariates. *Biometrika*, *93*, 791–807.
- Chib, S. (2007). Analysis of treatment response data without the joint distribution of potential outcomes. *Journal of Econometrics*, *140*, 401–412.
- Chung, Y., Dunson, D. B. (2009). Nonparametric Bayes conditional distribution modeling with variable selection. *Journal of the American Statistical Association*, *104*, 1646–1660.
- Cox, D. R. (1975). Partial likelihood. *Biometrika*, *62*, 269–276.
- Dunson, D. B., Pillai, N., Park, J. H. (2007). Bayesian density regression. *Journal of the Royal Statistical Society: Series B*, *69*, 163–183.
- Geweke, J. (1992). Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments. In J. M. Bernardo, J. O. Berger, A. P. Dawid, A. F. M. Smith (Eds.), *Bayesian Statistics*, Vol. 4. New York: Oxford University Press.
- Hannah, L. A., Blei, D. M., Powell, W. B. (2011). Dirichlet process mixtures of generalized linear models. *Journal of Machine Learning Research*, *12*, 1923–1953.
- Hansson, O., Zetterberg, H., Buchhave, P., Londos, E., Blennow, K., Minthon, L. (2006). Association between CSF biomarkers and incipient Alzheimer's disease in patients with mild cognitive impairment: A follow-up study. *The Lancet Neurology*, *5*, 228–234.
- Hirano, K. (2002). Semiparametric Bayesian inference in autoregressive panel data models. *Econometrica*, *70*, 781–799.
- Hoshino, T. (2013). Semiparametric Bayesian estimation for marginal parametric potential outcome modeling: Application to causal inference. *Journal of the American Statistical Association*, *108*, 1189–1204.
- Ibrahim, J. G., Chen, M. H., Lipsitz, S. R., Herring, A. H. (2005). Missing-data methods for generalized linear models: A comparative review. *Journal of the American Statistical Association*, *100*, 332–346.
- Ishwaran, H., James, L. F. (2001). Gibbs sampling methods for stick-breaking priors. *Journal of the American Statistical Association*, *96*, 161–173.
- Jack, C. R., Wiste, H. J., Vemuri, P., Weigand, S. D., Senjem, M. L., Zeng, G. (2010). Brain beta-amyloid measures and magnetic resonance imaging atrophy both predict time-to-progression from mild cognitive impairment to Alzheimer's disease. *Brain*, *133*, 3336–3348.
- Kalbfleisch, J. D. (1978). Non-parametric Bayesian analysis of survival time data. *Journal of the Royal Statistical Society: Series B*, *40*, 214–221.
- Kim, J. S., Ratchford, B. T. (2013). A Bayesian multivariate probit for ordinal data with semiparametric random-effects. *Computational Statistics and Data Analysis*, *64*, 192–208.
- Kottas, A., Muller, P., Quintana, F. (2005). Nonparametric Bayesian modeling for multivariate ordinal data. *Journal of Computational and Graphical Statistics*, *14*, 610–625.
- Kunihama, T., Herring, A. H., Halpern, C. T., Dunson, D. B. (2016). Nonparametric Bayes modeling with sample survey weights. *Statistics and Probability Letters*, *113*, 41–48.
- Lawless, J. F., Kalbfleisch, J. D., Wild, C. J. (1999). Semiparametric methods for response-selective and missing data problems in regression. *Journal of the Royal Statistical Society: Series B*, *61*, 413–438.
- Li, J., Wang, Z., Li, R., Wu, R. (2015). Bayesian group LASSO for nonparametric varying-coefficient models with application to functional genome-wide association studies. *The Annals of Applied Statistics*, *9*, 640–664.
- Liao, S. G., Lin, Y., Kang, D. D., Chandra, D., Bon, J., Kaminski, N., Sciruba, F. C., Tseng, G. C. (2014). Missing value imputation in high-dimensional phenomic data: Imputable or not, and how? *BMC Bioinformatics*, *15*, 346.
- Liu, J., Gelman, A., Hill, J., Su, Y. S., Kropko, J. (2014). On the stationary distribution of iterative imputations. *Biometrika*, *101*, 155–173.

- McCulloch, R., Rossi, P. E. (1994). An exact likelihood analysis of the multinomial probit model. *Journal of Econometrics*, *64*, 207–240.
- Meng, X. L. (1994). Multiple-imputation inferences with uncongenial sources of input. *Statistical Science*, *9*, 538–558.
- Murray, J. S., Reiter, J. P. (2016). Multiple imputation of missing categorical and continuous values via bayesian mixture models with local dependence. *Journal of the American Statistical Association*, *111*, 1466–1476.
- Okello, A., Koivunen, J., Edison, P., Archer, H. A., Turkheimer, F. E., Nagren, K. U., Bullock, R., Walker, Z., Kennedy, A., Fox, N. C., Rossor, M. N., Rinne, J. O., Brooks, D. J. (2009). Conversion of amyloid positive and negative MCI to AD over 3 years An 11C-PIB PET study. *Neurology*, *73*, 754–760.
- Pati, D., Dunson, D. B., Tokdar, S. T. (2013). Posterior consistency in conditional distribution estimation. *Journal of Multivariate Analysis*, *116*, 456–472.
- Paton, N. I., Kityo, C., Hoppe, A., Reid, A., Kambugu, A., Lugemwa, A. (2014). Assessment of second-line antiretroviral regimens for HIV therapy in Africa. *The New England Journal of Medicine*, *371*, 234–247.
- Reich, B. J., Bondell, H. D., Li, L. (2011). Sufficient dimension reduction via Bayesian mixture modeling. *Biometrics*, *67*, 886–895.
- Robins, J. M., Hsieh, F., Newey, W. (1995). Semiparametric efficient estimation of a conditional density with missing or mismeasured covariates. *Journal of the Royal Statistical Society. Series B*, *61*, 409–424.
- Rodriguez, A., Dunson, D. B., Gelfand, A. E. (2009). Bayesian nonparametric functional data analysis through density estimation. *Biometrika*, *96*, 149–162.
- Rubin, D. B. (1987). *Multiple imputation for nonresponse in surveys*. New York: Wiley.
- Sethuraman, J. (1994). A constructive definition of dirichlet priors. *Statistica Sinica*, *4*, 639–650.
- Shen, W., Tokdar, S. T., Ghosal, S. (2013). Adaptive Bayesian multivariate density estimation with Dirichlet mixtures. *Biometrika*, *100*, 623–640.
- Si, Y., Reiter, J. (2013). Nonparametric Bayesian multiple imputation for incomplete categorical variables in large-scale assessment surveys. *Journal of Educational and Behavioral Statistics*, *38*, 499–521.
- Sinha, D., Ibrahim, J. G., Chen, M. H. (2003). A Bayesian justification of Cox's partial likelihood. *Biometrika*, *90*, 629–641.
- Stekhoven, D. J., Buhlmann, P. (2012). MissForest-non-parametric missing value imputation for mixed-type data. *Bioinformatics*, *28*, 112–118.
- The National Research Council (2010). *The prevention and treatment of missing data in clinical trials*. Washington, DC: National Academic Press.
- Tokdar, S. T., Zhu, Y. M., Ghosh, J. K. (2010). Bayesian density regression with logistic Gaussian process and subspace projection. *Bayesian Analysis*, *5*, 319–344.
- van Buuren, S. (2007). Multiple imputation of discrete and continuous data by fully conditional specification. *Statistical Methods in Medical Research*, *16*, 219–242.
- van Buuren, S. (2012). *Flexible imputation of missing data*. Boca Raton, Florida: Chapman and Hall/CRC.
- Waljee, A. K., Mukherjee, A., Singal, A. G., Zhang, Y., Warren, J., Balis, U., Marrero, J., Zhu, J., Higgins, P. D. (2013). Comparison of imputation methods for missing laboratory data in medicine. *BMJ Open*, *3*, e002847.
- Walker, S. G., Damien, P., Laud, P. W., Smith, A. F. (1999). Bayesian nonparametric inference for random distributions and related functions. *Journal of the Royal Statistical Society: Series B*, *61*, 485–527.
- White, I. R., Royston, P., Wood, A. M. (2011). Multiple imputation using chained equations: Issues and guidance for practice. *Statistics in Medicine*, *30*, 377–399.
- Zhang, X., Boscardin, W. J., Belin, T. R. (2008). Bayesian analysis of multivariate nominal measures using multivariate multinomial probit models. *Computational Statistics and Data Analysis*, *52*, 3697–3708.
- Zhang, Z., Rockette, H. E. (2006). Semiparametric maximum likelihood for missing covariates in parametric regression. *Annals of the Institute of Statistical Mathematics*, *58*, 687–706.