# A note on the robustness of a full Bayesian method for nonignorable missing data analysis

### Zhiyong Zhang and Lijuan Wang
*University of Notre Dame*

**Abstract.** A full Bayesian method utilizing data augmentation and Gibbs sampling algorithms is presented for analyzing nonignorable missing data. The discussion focuses on a simplified selection model for regression analysis. Regardless of missing mechanisms, it is assumed that missingness only depends on the missing variable itself. Simulation results demonstrate that the simplified selection model can recover regression model parameters under both correctly specified situations and many misspecified situations. The method is also applied to analyzing a training intervention data set with missing data.

## 1 Introduction

Missing data problem is a big challenge in statistical inference even for a well designed study. Little and Rubin (2002) distinguished three kinds of missing data mechanisms—missing completely at random (MCAR), missing at random (MAR), and missing not at random (MNAR). For the MCAR mechanism, every subject (datum) has the same probability to be missing. For example, if a participant's test score is missing because his/her experimenter accidentally forgets to give him/her the test, then the resulting missing datum can be viewed as MCAR. If missingness can be fully predicted by observed data in a model, the missing mechanism is MAR. For example, in a pre-test post-test experiment, all subjects participated in the pre-test. Some subjects missed the post-test because they did not perform well in the pre-test. In this case, missingness during the post-test may be predicted using data from the pre-test. However, if missingness cannot be fully predicted by observed data, missing data are MNAR. For example, in survey research, when asking about salary, those with high salary often choose not to respond. Thus, missing data on income are often considered to be MNAR.

To some extent, MCAR and MAR data are ignorable because problems caused by them can be overcome through sophisticated statistical techniques such as the full information likelihood (FIML) method and multiple imputation (e.g., Little and Rubin, 2002, Schafer, 1997). MNAR data, however, are nonignorable because without extra information or modeling specifications, their influences cannot be

well addressed. To deal with nonignorable missing data, selection models (see recent reviews by Ibrahim et al., 2006 and Ibrahim and Molenberghs, 2009) and the multiple imputation (MI) method with auxiliary variables (e.g., Graham, 2009) can be used. However, both selection models and MI typically require auxiliary variables (information) to account for nonignorable missingness.

Best et al. (1996) proposed a simplified selection model to study cognitive decline in the elderly. The model assumed that missingness in an outcome variable was only related to itself and the model, therefore, did not require auxiliary variables to model the nonignorable missing mechanism. Through empirical data analysis, they demonstrated that model parameters of interest were insensitive to different prior specifications on the parameters predicting missingness. In this article, we examine and extend the model by Best et al. (1996) in several ways with a focus on the robustness of the model in recovering true regression parameter values using full Bayesian methods. First, we derive the full Bayesian posterior distributions for the simplified selection model. Second, we evaluate the performance of the model under different conditions. Third, we apply the model to analyze a set of training intervention data to illustrate the use of the model.

In the remainder of the paper, we will first discuss selection models and the simplified version by Best et al. (1996). Then, we will discuss how to estimate the simplified selection model through a full Bayesian method by obtaining a full set of conditional posterior distributions for model parameters utilizing the data augmentation algorithm. After this, we will evaluate the performance of the model and the estimation method through several simulation studies. Finally, we will present an empirical example to demonstrate how to apply the model and the method in behavioral research.

## 2 Selection models

Consider a multiple regression model with a dependent variable $y$ and a vector of independent variables $\mathbf{x}$. The regression model can be expressed as

$$y_i = \mathbf{x}_i \boldsymbol{\beta} + e_i, \tag{2.1}$$

where $y_i$ and $\mathbf{x}_i$ are observed data for the $i$th person, $\boldsymbol{\beta}$ is a vector of regression coefficients, and $e_i$ is residual that is assumed to be normally distributed such that $e_i \sim N(0, \phi)$ with $\phi$ denoting the variance of the residuals. If both the dependent and independent variables are fully observed, estimates of the regression coefficients can be obtained conveniently.

However, data are often incomplete or not fully observed for a variety of reasons such as non-response and dropout. For example, we consider a scenario that the dependent variable $y$ is partially observed and independent variables $\mathbf{x}$ are fully observed. This scenario is very common in designed experiments and survey research. For example, in an experiment, the controlled factors are often predetermined and thus are known. However, data of the outcome variables may not

be always observed. In survey research, subjects may respond to nonsensitive questions such as demographic variables including age and level of education, but are less likely to answer sensitive questions.

Let $m_i$ be an indicator variable where $m_i = 0$ if $y_i$ is observed and $m_i = 1$ if $y_i$ is missing. Then, the missing probability of $y_i$, in general, can be modeled as

$$\Pr(y_i \text{ is missing}) = \Pr(m_i = 1) = f(\gamma_0 + \gamma_1 y_i + \boldsymbol{\alpha} \mathbf{v}_i), \tag{2.2}$$

where $\mathbf{v}$ is a set of variables that may be related to the missingness of $y$. The variables in $\mathbf{v}$ may or may not be fully observed and could be latent variables. $\mathbf{v}$ may also include part of or all variables of $\mathbf{x}$. The $\gamma_0$, $\gamma_1$, and $\boldsymbol{\alpha}$ are model parameters. The link function $f$ can be any function that maps its input to an output value from 0 to 1. For example, if $f$ is a logistic function, equation (2.2) becomes a logistic regression model.

Equation (2.2) models the missing mechanism of $y$. The missing mechanisms defined by Little and Rubin (2002) can be distinguished according to the parameter values in equation (2.2). If $\gamma_1 = 0$ and $\boldsymbol{\alpha} = 0$, the missing probability is a constant and the missing mechanism is MCAR. If $\gamma_1 = 0$, and $\boldsymbol{\alpha} \neq 0$ and $\mathbf{v}$ are fully observed and included in the regression model as predictors, the missing mechanism is MAR. If $\gamma_1 \neq 0$ or the coefficients $\boldsymbol{\alpha}$ for the partially observed or unobserved variables in $\mathbf{v}$ are not equal to zero, missing data are MNAR.

Selection models (e.g., Heckman, 1976, Little and Rubin, 2002) focus on modeling the joint distribution of observed data and missing indicators as

$$p(y_i, m_i | \beta, \phi, \boldsymbol{\gamma}, \boldsymbol{\alpha}, \mathbf{x}_i, \mathbf{v}_i) = p(y_i | \mathbf{x}_i, \beta, \phi) p(m_i | y_i, \boldsymbol{\gamma}, \mathbf{v}_i, \boldsymbol{\alpha}), \tag{2.3}$$

where $p(\cdot)$ represents the probability density function and $\boldsymbol{\gamma} = (\gamma_0, \gamma_1)'$. A possible choice for $p(m_i | y_i, \boldsymbol{\gamma}, \mathbf{v}_i, \boldsymbol{\alpha})$ is given in equation (2.2). Model parameters in selection models can be estimated by maximizing the following likelihood function (e.g., Little, 1982, Little and Rubin, 2002)

$$L = \int_{\mathbf{v}_i^{\text{miss}}} \int_{y_i^{\text{miss}}} \prod_{i=1}^{n} p(y_i, m_i | \beta, \phi, \gamma, \boldsymbol{\alpha}, \mathbf{x}_i, \mathbf{v}_i) \, dy_i^{\text{miss}} \mathbf{v}_i^{\text{miss}}, \tag{2.4}$$

where $y_i^{\text{miss}}$ denotes a missing datum in $y$ and $\mathbf{v}_i^{\text{miss}}$ denotes unobserved data in $\mathbf{v}$ for individual $i$. To identify selection models, $\mathbf{x}_i$ and $\mathbf{v}_i$ should not completely overlap, or the values of $\lambda$ or $\boldsymbol{\alpha}$ for some covariate(s) are known or preconstrained (e.g., Little, 1985, Olsen, 1980, Tang, Little and Raghunathan, 2003).

To properly apply selection models, variables in $\mathbf{v}$ have to be determined carefully, which often hinders practical adoptions of selection models. Best et al. (1996) discussed a simplified selection model with $\boldsymbol{\alpha} \equiv 0$. Thus, the model assumes that missingness in $y$ is only related to itself and no auxiliary variables need to be used in the model. One of the advantages of the model is the avoidance of selecting auxiliary variables. Although the function $f$ in equation (2.2)

can take many different forms, the logistic function (the logit link) and the cumulative normal distribution function (the probit link) are most widely used (e.g., Ibrahim et al., 2006, Little and Rubin, 2002). Best et al. used the logistic function in equation (2.2). In this study, we use the cumulative normal distribution function because of its convenience in obtaining posterior distributions. Therefore, the missing mechanism in the simplified selection model can be specified as

$$\begin{cases} m_i \sim \text{Bernoulli}(p_i), \\ p_i = \Phi(\gamma_0 + \gamma_1 y_i) = \Phi\left[ (1 \quad y_i) \begin{pmatrix} \gamma_0 \\ \gamma_1 \end{pmatrix} \right] = \Phi(\mathbf{w}_i \boldsymbol{\gamma}), \end{cases} \tag{2.5}$$

where $p_i$ is the probability that $y_i$ is missing and $\Phi$ is the cumulative normal distribution function. $m_i$ is a binary variable following a Bernoulli distribution. Furthermore, $\mathbf{w}_i = (1, y_i)$ and $\boldsymbol{\gamma} = (\gamma_0, \gamma_1)'$.

Note that if the missing mechanism is MCAR or MNAR depending on $y$ only, the missing mechanism is correctly specified. Otherwise, the missing mechanism is misspecified, for example, under the situation of MAR. For convenience, we refer to the model in equations (2.1) and (2.5) together as the simplified selection model and the model in equation (2.1) as the regular model in the remainder of the paper. Because of the focus of Best et al.'s study, the posterior distributions of the selection model were not discussed and conditions under which the simplified selection model would/would not work were not investigated. Thus, in this study, we will present posterior distributions of the simplified selection model and then evaluate its performance under a variety of conditions.

## 3 Full Bayesian estimation method

To estimate the simplified selection model, we will use the Bayesian estimation method based on data augmentation (Albert and Chib, 1993, Tanner and Wong, 1987) and Gibbs sampling (e.g., Casella and George, 1992). We first assume that there is an underlying normal variable $z_i \sim N(\mathbf{w}_i \boldsymbol{\gamma}, 1)$ for each $m_i$. If $z_i > 0$, then $m_i = 1$. Otherwise, $m_i = 0$. In other words, if $z_i > 0$, $y_i$ is unobserved. By augmenting $z_i$ with $m_i$, the joint distribution of $m_i$ and $z_i$ is

$$p(m_i, z_i | y_i, \boldsymbol{\gamma}) = p(m_i | z_i) p(z_i | y_i, \boldsymbol{\gamma}).$$

The distribution of $z_i$, $p(z_i | y_i, \boldsymbol{\gamma})$, is already known as a normal distribution with mean $\mathbf{w}_i \boldsymbol{\gamma}$ and variance 1 and we need to obtain the distribution for $m_i$ conditional on $z_i$. Note that

$$p(m_i = 1 | z_i > 0) = 1, \qquad p(m_i = 1 | z_i \leq 0) = 0,$$
$$p(m_i = 0 | z_i > 0) = 0, \qquad p(m_i = 0 | z_i \leq 0) = 1.$$

Thus, the distribution for $m_i | z_i$ can be expressed as

$$p(m_i | z_i) = \mathcal{I}(m_i = 1)\mathcal{I}(z_i > 0) + \mathcal{I}(m_i = 0)\mathcal{I}(z_i \leq 0),$$

where $\mathcal{I}(A)$ is an indicator function which takes 1 if the expression $A$ is true and otherwise 0.

Furthermore, by augmenting missing data $y_i^{\text{miss}}$ with observed data, $m_i$ and $z_i$, the joint distribution of $y_i$, $m_i$ and $z_i$ is

$$
\begin{aligned}
p(y_i, z_i, m_i | \boldsymbol{\beta}, \boldsymbol{\gamma}, \phi, \mathbf{X}) &= p(y_i | \boldsymbol{\beta}, \phi, \mathbf{X}) p(m_i | z_i) p(z_i | y_i, \boldsymbol{\gamma}) \\
&= \frac{1}{\sqrt{2\pi\phi}} \exp\left[ -\frac{(y_i - \mathbf{x}_i \boldsymbol{\beta})^2}{2\phi} \right] \\
&\quad \times [\mathcal{I}(m_i = 1)\mathcal{I}(z_i > 0) + \mathcal{I}(m_i = 0)\mathcal{I}(z_i \le 0)] \\
&\quad \times \frac{1}{\sqrt{2\pi}} \exp\left[ -\frac{(z_i - \gamma_0 - \gamma_1 y_i)^2}{2} \right].
\end{aligned}
\tag{3.1}
$$

Thus, the likelihood function for the selection model with augmented data can be expressed as

$$
L(\boldsymbol{\beta}, \boldsymbol{\gamma}, \phi | \mathbf{y}, \mathbf{X}, \mathbf{z}, \mathbf{m}) = \prod_{i=1}^{n} p(y_i, z_i, m_i | \boldsymbol{\beta}, \boldsymbol{\gamma}, \phi, \mathbf{X}),
\tag{3.2}
$$

where $\mathbf{y} = (y_1, \ldots, y_n)'$ denotes a vector of data for the dependent variable and $\mathbf{X} = (\mathbf{x}_1', \ldots, \mathbf{x}_n')'$ is the design matrix. Furthermore, $\mathbf{z} = (z_1, \ldots, z_n)'$ and $\mathbf{m} = (m_1, \ldots, m_n)'$. By integrating out missing data $y_i^{\text{miss}}$ and the underlying variable $z_i$, one can obtain the observed data likelihood for the maximum likelihood estimation method. However, the integration is not an easy task. Bayesian estimation procedure, however, can be implemented relatively easily through Gibbs sampling after obtaining the full set of conditional posterior distributions for the selection model.

To use the Bayesian method, we need to specify priors for unknown parameters. We first consider the following semi-conjugate priors and then discuss the use of Jeffreys priors (Gelman et al., 2003, Jeffreys, 1946). For the regression coefficients $\boldsymbol{\beta}$, a multivariate normal prior is used as

$$
p(\boldsymbol{\beta}) = \text{MN}(\boldsymbol{\beta}_0, \Sigma_0) \propto |\Sigma_0|^{-1/2} \exp\left[ -\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)' \Sigma_0^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_0) \right],
\tag{3.3}
$$

where $\boldsymbol{\beta}_0$ and $\Sigma_0$ are predefined hyper-parameters representing the mean vector and covariance matrix of the multivariate normal distribution.[1] For the residual variance parameter $\phi$, an inverse gamma distribution is employed,

$$
p(\phi) = \text{IG}(a_0, b_0) \propto \phi^{-a_0/2 - 1} \exp\left( -\frac{b_0}{2\phi} \right),
\tag{3.4}
$$

---

[1] Note that $\boldsymbol{\beta}_0$, $\Sigma_0$ and other symbols with subscript in a distribution represent the hyper-parameters of the prior or posterior distributions.

where $a_0$ and $b_0$ are assumed to be known shape and scale parameters. Finally, for parameters $\boldsymbol{\gamma}$, a multivariate normal prior is also used,

$$p(\boldsymbol{\gamma}) = \mathrm{MN}(\boldsymbol{\gamma}_0, D_0) \propto |D_0|^{-1/2} \exp\left[-\frac{1}{2}(\boldsymbol{\gamma} - \boldsymbol{\gamma}_0)' D_0^{-1}(\boldsymbol{\gamma} - \boldsymbol{\gamma}_0)\right], \quad (3.5)$$

where $\boldsymbol{\gamma}_0$ and $D_0$ are known mean vector and covariance matrix. These priors are called semi-conjugate priors because the corresponding conditional posteriors are from the same distribution family (Gelman et al., 2003).

With the likelihood function in (3.2) and priors in (3.3), (3.4) and (3.5), the joint posterior distribution of the unknown parameters is readily available. However, the marginal posterior distributions of the parameters are difficult to obtain explicitly. To avoid the difficulty of getting the marginal posterior distributions explicitly, we obtain the conditional distributions of the parameters and then utilize the Gibbs sampling method to generate Markov chains for the parameters and construct the Bayesian parameter estimates through posterior means.

The conditional posterior distribution for $\boldsymbol{\beta}$ is a multivariate normal distribution defined by

$$\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}, \phi \sim \mathrm{MN}(\boldsymbol{\beta}_1, \Sigma_1),$$

where

$$\boldsymbol{\beta}_1 = (\Sigma_0^{-1} + \mathbf{X}'\mathbf{X}\phi^{-1})^{-1}(\Sigma_0^{-1}\boldsymbol{\beta}_0 + \mathbf{X}'\mathbf{y}\phi^{-1})$$

and

$$\Sigma_1 = (\Sigma_0^{-1} + \mathbf{X}'\mathbf{X}\phi^{-1})^{-1},$$

where $\mathbf{X}$ and $\mathbf{y}$ are as defined earlier.

The conditional posterior distribution for $\phi$ is an inverse Gamma distribution given by

$$\phi|\boldsymbol{\beta}, \mathbf{y}, \mathbf{X} \sim \mathrm{IG}(a_1, b_1),$$

where

$$a_1 = \frac{a_0 + n}{2}$$

and

$$b_1 = \frac{b_0 + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2}.$$

The conditional posterior distribution for $\boldsymbol{\gamma}$ is

$$\boldsymbol{\gamma}|\mathbf{z}, \mathbf{y}, \mathbf{W} \sim \mathrm{MN}(\boldsymbol{\gamma}_1, D_1),$$

where

$$\boldsymbol{\gamma}_1 = (D_0^{-1} + \mathbf{W}'\mathbf{W})^{-1}(D_0^{-1}\boldsymbol{\gamma}_0 + \mathbf{W}'\mathbf{z})$$

and

$$D_1 = (D_0^{-1} + \mathbf{W}'\mathbf{W})^{-1},$$

where $\mathbf{W} = (\mathbf{w}_1', \ldots, \mathbf{w}_n')'$.

The conditional posterior distribution for the underlying variable $z_i$ is

$$z_i | \boldsymbol{\gamma}, y_i, m_i \sim \begin{cases} N(\mathbf{w}_i\boldsymbol{\gamma}, 1)I(0, +\infty), & m_i = 1, \\ N(\mathbf{w}_i\boldsymbol{\gamma}, 1)I(-\infty, 0], & m_i = 0. \end{cases}$$

Thus, $z_i$ follows a truncated normal distribution.

Finally, the conditional posterior distribution for missing data $y_i^{\mathrm{miss}}$ is

$$y_i^{\mathrm{miss}} | \boldsymbol{\beta}, \boldsymbol{\gamma}, z_i, \phi, \mathbf{x}_i \sim N(\mu_1, \phi_1),$$

where

$$\mu_1 = \left[ \frac{\mathbf{x}_i\boldsymbol{\beta}}{\phi} + \gamma_1(z_i - \gamma_0) \right] \left( \frac{1}{\phi} + \gamma_1^2 \right)^{-1}$$

and

$$\phi_1 = \left( \frac{1}{\phi} + \gamma_1^2 \right)^{-1}.$$

Note that if missing data are MCAR, one should expect that $\gamma_1 = 0$. Then the posterior for missing data $y_i^{\mathrm{miss}}$ reduces to

$$y_i^{\mathrm{miss}} | \boldsymbol{\beta}, \boldsymbol{\gamma}, z_i, \phi, \mathbf{x}_i \sim N(\mathbf{x}_i\boldsymbol{\beta}, \phi)$$

which can be viewed as the multiple imputation method.

Jeffreys priors (Jeffreys, 1946) can also be used in obtaining the conditional posterior distributions of the model parameters. One form of Jeffreys priors can be

$$p(\phi) \propto 1/\phi$$

and

$$p(\boldsymbol{\beta}, \boldsymbol{\gamma}) \propto 1.$$

Note that these priors are improper and can be viewed as carrying no information (Box and Tiao, 1973). With these Jeffreys priors, the conditional posterior distributions for $z_i$ and $y_i^{\mathrm{miss}}$ remain the same. The conditional posterior distributions for the model parameters become

$$\boldsymbol{\beta} | y_i, \mathbf{x}_i, \phi \sim \mathrm{MN}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}, (\mathbf{X}'\mathbf{X})^{-1}\phi],$$

$$\phi | \boldsymbol{\beta}, y_i, \mathbf{x}_i \sim \mathrm{IG}\left[ \frac{n}{2}, \frac{(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})}{2} \right],$$

$$\boldsymbol{\gamma} | z_i, y_i \sim \mathrm{MN}[(\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\mathbf{Z}, (\mathbf{W}'\mathbf{W})^{-1}],$$

where $\mathbf{Y} = (\mathbf{y}_1', \ldots, \mathbf{y}_n')'$ and $\mathbf{Z} = (\mathbf{z}_1', \ldots, \mathbf{z}_n')'$.

With the conditional posterior distributions, one can implement the following Gibbs sampling procedure.

1. Start with initial values $\boldsymbol{\beta}^{(0)}, \phi^{(0)}, \boldsymbol{\gamma}^{(0)}, z_i^{(0)}, y_i^{\text{miss}(0)}$.
2. Assume at the iteration $t$, one has $\boldsymbol{\beta}^{(t)}, \phi^{(t)}, \boldsymbol{\gamma}^{(t)}, z_i^{(t)}, y_i^{\text{miss}(t)}$.
3. At iteration $t + 1$,
   (a) generate $\boldsymbol{\beta}^{(t+1)}$ from $\boldsymbol{\beta}|y_i^{\text{obs}}, y_i^{\text{miss}(t)}, \mathbf{x}_i, \phi^{(t)}$,
   (b) generate $\phi^{(t+1)}$ from $\phi|\boldsymbol{\beta}^{(t+1)}, y_i^{\text{obs}}, y_i^{\text{miss}(t)}, \mathbf{x}_i$,
   (c) generate $z_i^{(t+1)}$ from $z_i|\boldsymbol{\gamma}^{(t)}, y_i^{\text{obs}}, y_i^{\text{miss}(t)}, m_i$ for $i = 1, \ldots, n$,
   (d) generate $\boldsymbol{\gamma}^{(t+1)}$ from $\boldsymbol{\gamma}|z_i^{(t+1)}, y_i^{\text{obs}}, y_i^{\text{miss}(t)}$,
   (e) generate $y_i^{\text{miss}(t+1)}$ from $y_i^{\text{miss}}|\boldsymbol{\beta}^{(t+1)}, \boldsymbol{\gamma}^{(t+1)}, z_i^{(t+1)}, \phi^{(t+1)}, \mathbf{x}_i$ if $y_i$ is missing.

The above Gibbs sampling procedure can be used to generate a Markov chain for each model parameter, underlying variable, and missing datum. After convergence, these Markov chains can be viewed as samples from the joint distribution and marginal distributions of the parameters and thus Bayesian parameter estimates can be constructed (e.g., Casella and George, 1992, Geman and Geman, 1984).

## 4 Simulation studies

In this section, we conduct several simulation studies to investigate the performance of the simplified selection model and the Bayesian estimation method. The focus of the simulation studies is to evaluate whether the simplified selection model with the Bayesian estimation method is robust to choices of priors, choices of link functions, and missing data mechanisms. In the first study, we compare the parameter estimates from the selection model and the regular regression model assuming MNAR and the missing mechanism is known to depend solely on $\mathbf{y}$. In the second simulation, we evaluate whether results from the selection model are influenced by different choices of priors, and in the third simulation, we investigate whether results from the selection model are influenced by different choices of link functions. In Simulation 4, the missing mechanism is related to an external unobserved variable. In Simulation 5, the missingness depends on both the dependent variable $\mathbf{y}$ and independent variables $\mathbf{X}$. In Simulation 6, the missing mechanism is MCAR.

### 4.1 General settings of the simulation studies

In all simulation studies, data are generated from a multiple regression model with two covariates as in

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + e_i = \mathbf{x}_i \boldsymbol{\beta} + e_i. \tag{4.1}$$

The population parameters are $\beta_0 = \beta_1 = \beta_2 = 1$ and $e_i \sim N(0, \phi)$ with variance $\phi = 0.25$. Both covariates $x_{1i}$ and $x_{2i}$ are generated from the standard normal distribution. Complete data are first generated from this model and missing data are

then generated according to each simulation condition. In all simulation studies, sample size is $n = 100$ and the missing data percentage is 40%. Furthermore, for each simulation study, a total of 1,000 data sets with missing data are generated and analyzed.

For each simulation study, four statistics will be reported. The first one is the parameter estimate which is calculated as the average of parameter estimates of 1,000 simulation replications. The second one is the average standard error (ASE) which is the average of standard errors of parameter estimates from the 1,000 replications. The third one is the standard deviation of the 1,000 sets of estimated parameters. Finally, the coverage probability of the 95% credible (confidence) interval of each parameter are also reported.

In the following simulation studies except for the second simulation, the following priors are used without further elaboration. For $\beta$, the trivariate normal prior is used with $\beta_0 = (0, 0, 0)'$ and

$$\Sigma_0 = \begin{pmatrix} 10^6 & 0 & 0 \\ 0 & 10^6 & 0 \\ 0 & 0 & 10^6 \end{pmatrix}.$$

For $\phi$, the inverse gamma prior is used with $a_0 = b_0 = 10^{-3}$. And for $\gamma$, the bivariate normal prior is used with $\gamma_0 = (0, 0)'$ and

$$D_0 = \begin{pmatrix} 10^6 & 0 \\ 0 & 10^6 \end{pmatrix}.$$

These priors can be considered as carrying little information (Congdon, 2003). For the second simulation, the Jeffreys priors are used.

All simulations are conducted using SAS and WinBUGS (Zhang et al., 2008). The convergence of the Markov chains is monitored through the Geweke statistics (Geweke, 1992). The WinBUGS codes for the simplified selection model are provided in the Appendix.

## 4.2 Simulation 1: Missingness depends on y only

4.2.1 *Purpose.* This simulation study investigates whether the regression model parameters ($\beta$ and $\phi$) can be recovered using the simplified selection model and the Bayesian estimation method when the missing mechanism is correctly specified.

4.2.2 *Missing data generation.* In this simulation study, the probability that a datum is missing is assumed to depend on itself. Thus, the missingness is nonignorable. Missing data are generated in the following way. Let $c_\alpha$ denote the $100\alpha$th percentile of $\mathbf{y}$. Then, the probability that $y_i$ is missing is set as

$$\Pr(y_i \text{ is missing}) = \Pr(m_i = 1) = \begin{cases} 0.9r/(1 - \alpha), & \text{if } y_i > c_\alpha, \\ 0.1r/\alpha, & \text{otherwise,} \end{cases}$$

**Table 1** *MNAR data analysis using the simplified selection model and the regular model*

| Parameters | Simplified selection model | | | | Regular model | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Estimates | ASE | SD | Coverage | Estimates | ASE | SD | Coverage |
| $\beta_0$ (Intercept) | 0.998 | 0.082 | 0.081 | 0.947 | 0.855 | 0.075 | 0.076 | 0.477 |
| $\beta_1$ ($x_1$) | 0.994 | 0.075 | 0.077 | 0.945 | 0.922 | 0.074 | 0.076 | 0.808 |
| $\beta_2$ ($x_2$) | 0.997 | 0.075 | 0.076 | 0.932 | 0.924 | 0.074 | 0.077 | 0.823 |
| $\phi$ | 0.261 | 0.056 | 0.054 | 0.945 | 0.236 | 0.045 | 0.044 | 0.922 |
| $\gamma_0$ | −1.488 | 0.365 | 0.503 | – | – | – | – | – |
| $\gamma_1$ ($y$) | 1.023 | 0.218 | 0.323 | – | – | – | – | – |

Note. The results are based on 1,000 simulation replications. ASE: average standard error. SD: standard deviation of parameter estimates. Coverage: coverage probability of the 95% highest posterior density credible interval.

where $r$ is the predefined missing data rate. In this simulation study, we set $\alpha = 60\%$ and $r = 0.4$. Note that the missing probability function is a step function instead of a continuous function. It also indicates that when $y_i$ is larger than the 60th percentile, its missing probability is 0.9. Otherwise, its missing probability is about 0.067.

4.2.3 *Results*. The simulated data are analyzed using the simplified selection model in equations (2.1) and (2.5). For the purpose of comparison, we also analyze the data by ignoring the nonignorable missing mechanism, namely, the data are analyzed as MAR through a regular regression model in equation (2.1). The results from the analysis are summarized in Table 1.

When the missing mechanism is MNAR and the data are analyzed ignoring the missing mechanism, the parameter estimates are biased—underestimated in this simulation study. Furthermore, the coverage probabilities are not correct, much smaller than the nominal level 95% especially for the regression coefficients ($\beta$). When the data are analyzed using the simplified selection model, the parameters are well recovered, especially for the regression coefficients, with less than 0.6% bias. The coverage probabilities are also close to 0.95. Finally, the ASE and SD are very close for the regression parameters, which indicates that the standard error estimates of parameters are also accurate. Note that although the missing mechanism is modeled using the cumulative normal distribution function and the missing data are generated according to a step function, the regression model parameter ($\beta$ and $\phi$) estimates are still accurate. This indicates that the nonignorable missing mechanism does not need to be perfectly specified.

### 4.3 Simulation 2: Semi-conjugate priors vs. Jeffreys priors (sensitivity analysis of priors)

4.3.1 *Purpose*. During the discussion of the model estimation method, we discussed two types of priors—the semi-conjugate priors and the Jeffreys priors. This

**Table 2**  *Results from the simplified selection model with different priors*

| Parameters | Semi-conjugate priors | | | | Jeffreys priors | | | |
|---|---|---|---|---|---|---|---|---|
| | Estimates | ASE | SD | Coverage | Estimates | ASE | SD | Coverage |
| $\beta_0$ (Intercept) | 0.998 | 0.082 | 0.081 | 0.947 | 1.004 | 0.084 | 0.082 | 0.954 |
| $\beta_1$ $(x_1)$ | 0.994 | 0.075 | 0.077 | 0.945 | 0.997 | 0.077 | 0.077 | 0.947 |
| $\beta_2$ $(x_2)$ | 0.997 | 0.075 | 0.076 | 0.932 | 1.000 | 0.076 | 0.077 | 0.933 |
| $\phi$ | 0.261 | 0.056 | 0.054 | 0.945 | 0.272 | 0.060 | 0.057 | 0.950 |
| $\gamma_0$ | −1.488 | 0.365 | 0.503 | – | −1.514 | 0.385 | 0.553 | – |
| $\gamma_1$ $(y)$ | 1.023 | 0.218 | 0.323 | – | 1.036 | 0.229 | 0.356 | – |

Note. The same as the previous table.

simulation study is to investigate whether the estimated model parameters ($\beta$ and $\phi$) are influenced by the choice of the two sets of priors.

4.3.2 *Missing data generation*.    The data are generated using the same procedure in Simulation 1.

4.3.3 *Results*.    The results from the analysis using two types of priors are given in Table 2. From the results, the parameter estimates, especially the regression coefficients ($\boldsymbol{\beta}$), are very close from two different types of priors. The results for using the Jeffreys priors can be viewed as the baseline data analysis. Prior information, when available, can be incorporated into data analysis through the semi-conjugate priors.[2]

## 4.4  Simulation 3: Cumulative normal distribution function vs. logistic function

4.4.1 *Purpose*.    In equation (2.2), the $f$ function can be any function that maps a raw value to be within 0 and 1. For derivation convenience, we have used the cumulative normal distribution function in the previous section. With it, the explicit forms of the conditional posterior distributions can be obtained conveniently. However, the other functions such as the logistic function can also be used as in Best et al. (1996). In this simulation study, we investigate whether the choice of the function $f$ such as the cumulative normal distribution function or the logistic function influences regression parameter ($\beta$ and $\phi$) estimates.

4.4.2 *Missing data generation*.    The data are generated using the same procedure in Simulation 1.

---

[2]Best et al. (1996) compared three sets of informative priors and found that parameter estimates were insensitive to the chosen priors in their examples.

**Table 3** *Results from the simplified selection model with different f*

| Parameters | Normal distribution function | | | | Logistic function | | | |
|---|---|---|---|---|---|---|---|---|
| | Estimates | ASE | SD | Coverage | Estimates | ASE | SD | Coverage |
| $\beta_0$ (Intercept) | 0.998 | 0.082 | 0.081 | 0.947 | 0.996 | 0.080 | 0.078 | 0.958 |
| $\beta_1$ $(x_1)$ | 0.994 | 0.075 | 0.077 | 0.945 | 0.994 | 0.074 | 0.075 | 0.947 |
| $\beta_2$ $(x_2)$ | 0.997 | 0.075 | 0.076 | 0.932 | 0.997 | 0.074 | 0.074 | 0.933 |
| $\phi$ | 0.261 | 0.056 | 0.054 | 0.945 | 0.256 | 0.054 | 0.052 | 0.946 |
| $\gamma_0$ | −1.488 | 0.365 | 0.503 | – | −2.721 | 0.686 | 0.663 | – |
| $\gamma_1$ $(y)$ | 1.023 | 0.218 | 0.323 | – | 1.903 | 0.439 | 0.420 | – |

Note. The same as the previous table.

4.4.3 *Results.* The results from the simplified selection model using the cumulative normal distribution and logistic functions are provided in Table 3. The results are almost identical for the regression parameters ($\boldsymbol{\beta}$). For the $\gamma$s, the estimates using the logistic distribution is about 1.83 times of those using the cumulative normal distribution function. This simulation suggests that the selection model does not require the link function to match the missing mechanism exactly.

### 4.5 Simulation 4: Missing data depend on an external and unobserved variable z

4.5.1 *Purpose.* In this simulation, we consider the situation that the missingness of **y** depends on an external variable **z** that is related to **y** but unobserved. The purpose of this simulation is to investigate whether we can recover regression model parameters using the simplified selection model when the missing mechanism is misspecified.

4.5.2 *Missing data generation.* To generate data, $z_i$ is first generated using

$$z_i = ay_i + v_i,$$

where $v_i \sim N(0, 1)$. To investigate whether the correlation between $y$ and $z$ influences parameter estimates, we set $a$ at 1 and 0.2 corresponding to the correlations ($\rho$) between $y$ and $z$ at 0.83 and 0.29, respectively. Let $c_\alpha$ denote the $100\alpha$th percentile of **z**. Then, the probability that $y_i$ is missing is set as

$$\Pr(y_i \text{ is missing}) = \Pr(m_i = 1) = \begin{cases} 0.9r/(1-\alpha), & \text{if } z_i > c_\alpha, \\ 0.1r/\alpha, & \text{otherwise,} \end{cases}$$

where $r$ is the missing data rate. As in the previous simulations, we set $\alpha = 60\%$ and $r = 0.4$.

**Table 4** *Results from the simulation study that missing mechanism is related to an external variable* **z**

| Parameters | $a = 1$ and $\rho = 0.83$ | | | | $a = 0.2$ and $\rho = 0.29$ | | | |
|---|---|---|---|---|---|---|---|---|
| | Estimates | ASE | SD | Coverage | Estimates | ASE | SD | Coverage |
| $\beta_0$ (Intercept) | 1.009 | 0.078 | 0.078 | 0.945 | 1.003 | 0.067 | 0.067 | 0.952 |
| $\beta_1$ ($x_1$) | 1.001 | 0.072 | 0.072 | 0.955 | 0.998 | 0.066 | 0.067 | 0.942 |
| $\beta_2$ ($x_2$) | 1.000 | 0.072 | 0.071 | 0.946 | 1.003 | 0.066 | 0.069 | 0.944 |
| $\phi$ | 0.263 | 0.055 | 0.052 | 0.960 | 0.259 | 0.049 | 0.047 | 0.952 |
| $\gamma_0$ | −1.061 | 0.258 | 0.252 | – | −0.557 | 0.174 | 0.127 | – |
| $\gamma_1$ ($y$) | 0.680 | 0.150 | 0.162 | – | 0.191 | 0.099 | 0.103 | – |

Note. The same as the previous table.

4.5.3 *Results.* Results from this simulation are summarized in Table 4. Biases in parameter estimates ($\boldsymbol{\beta}$) are very small and the coverage probabilities are close to 95% regardless of the size of the correlation between the external variable and the outcome variable. Thus, the regression model parameters can still be recovered well using the simplified selection model even when the missing mechanism is misspecified in the current situation.

## 4.6 Simulation 5: Missing data depend on both y and X (a mixture of MNAR and MAR)

4.6.1 *Purpose.* This simulation study investigates whether the simplified selection model is robust to the missing mechanism situation that the missing probability of $y$ depends on both **y** and **X**. The missing mechanism can be viewed as a mixture of MNAR and MAR. Thus, the missing mechanism is also misspecified in the simplified selection model in this simulation study.

4.6.2 *Missing data generation.* To generate missing data, we first generate a variable **z** such that

$$z_i = 0.5y_i + x_{1i} - 0.5x_{2i}.$$

Then the missing probability is set by

$$\Pr(y_i \text{ is missing}) = \Pr(m_i = 1) = \begin{cases} 0.9r/(1-\alpha), & \text{if } z_i > c_\alpha, \\ 0.1r/\alpha, & \text{otherwise,} \end{cases}$$

where $c_\alpha$ is the $100\alpha$th percentile of **z**. In this simulation, we have $\alpha = 60\%$ and $r = 0.4$.

4.6.3 *Results.* Results from this simulation study are given in Table 5. The biases of the regression coefficient parameters ($\boldsymbol{\beta}$) are slightly larger than the previous simulations. For example, for the intercept $\beta_0$, the relative bias is about 3.8%.

**Table 5**  *Results from the simulation study that missing probability is related to both* **y** *and* **X**

| Parameters | Estimates | ASE | SD | Coverage |
|---|---|---|---|---|
| $\beta_0$ (Intercept) | 1.038 | 0.084 | 0.082 | 0.934 |
| $\beta_1$ ($x_1$) | 1.026 | 0.080 | 0.077 | 0.954 |
| $\beta_2$ ($x_2$) | 0.998 | 0.073 | 0.072 | 0.949 |
| $\phi$ | 0.271 | 0.059 | 0.055 | 0.953 |
| $\gamma_0$ | $-1.438$ | 0.364 | 0.476 | – |
| $\gamma_1$ ($y$) | 0.948 | 0.209 | 0.289 | – |

Note. The same as the previous table.

However, the bias is still small (less than 5%). Furthermore, the coverage probabilities are close to 95%. Thus, the simplified selection model appears to work well under the current condition.

### 4.7  Simulation 6: MCAR data analysis using the selection model

4.7.1 *Purpose.*  This simulation study investigates whether the simplified selection model can be applied when the missing mechanism is MCAR. Note that for MCAR, the simplified selection model is correctly specified and one would expect that $\hat{\gamma}_1 = 0$.

4.7.2 *Missing data generation.*  To generate missing data, the missing probability is set as

$$\Pr(y_i \text{ is missing}) = \Pr(m_i = 1) = r,$$

where $r$ is a constant and is set as 0.4 in this simulation.

4.7.3 *Results.*  Results from this simulation are given in Table 6. First, the estimate of $\gamma_1$ is very close to 0. Thus, the selection model correctly identifies that the missing mechanism does not depend on **y**. Second, the regression coefficient ($\beta$) estimates are very close to the true values with the maximum relative bias about 0.3%. Third, the estimate of $\phi$ is not as accurate as the regression coefficients estimates but close to the true value 0.25. Finally, the coverage probabilities are close to the nominal value 0.95. Clearly, if the missing mechanism is MCAR, the simplified selection model can still be applied.

## 5  An empirical example

In this section, we illustrate the application of the simplified selection model through the analysis of a subset of data from the Advanced Cognitive Training

**Table 6**    *Analyze MCAR data using the simplified selection model*

| Parameters | Estimates | ASE | SD | Coverage |
|------------|-----------|-----|-----|----------|
| $\beta_0$ (Intercept) | 1.002 | 0.068 | 0.066 | 0.958 |
| $\beta_1$ ($x_1$) | 1.002 | 0.067 | 0.066 | 0.944 |
| $\beta_2$ ($x_2$) | 1.003 | 0.067 | 0.067 | 0.937 |
| $\phi$ | 0.259 | 0.051 | 0.049 | 0.951 |
| $\gamma_0$ | $-0.258$ | 0.159 | 0.096 | – |
| $\gamma_1$ ($y$) | $-0.002$ | 0.093 | 0.093 | 0.955 |

Note. The same as the previous table.

for Independent and Vital Elderly (ACTIVE) study. The ACTIVE study is a randomized and controlled study designed to determine whether cognitive training interventions can affect cognitively based measures of daily functioning (Jobe et al., 2001, Tennstedt, 2001). For the purpose of illustration, the analysis here focuses on whether booster training on memory can improve everyday problem solving ability (EPT) of the elderly.

The sample size for this data analysis is $N = 703$ with about 53% (372) participants selected randomly to receive the booster training on memory. Before and after the booster training, everyday problem solving ability test was administered to the participants. The time interval between the two tests was about one year. In this data set, the change scores ($\Delta$EPT) as the difference between test scores before and after training were available for 76.5% (583) participants (about 23.5% participants had missing data). It is hypothesized that the training group has a larger $\Delta$EPT than the control group. To control possible confounding factors, we also included demographic variables, age and education level, in our data analysis. There are no missing data on the demographic variables.

Table 7 presents the summary statistics of each group in the data analysis. From Table 7, we can see that the control group had more missing data on EPT than the training group. There were no significant differences in age and education between the two groups. Both groups on average had negative change on EPT and the control group seemed to have more everyday functioning decline. Based on the two-sample $t$-test on $\Delta$EPT with list-wise deletion, the average difference on $\Delta$EPT between two groups was not significant at the significance level of 0.05.

Two models were fitted to the data. The first one is a regular regression model assuming that missingness on $\Delta$EPT is ignorable. Thus, the model can be written as

$$\Delta\text{EPT}_i = \beta_0 + \beta_1 \text{Training}_i + \beta_2 \text{Age}_i + \beta_3 \text{Education}_i + e_i.$$

The second one is the simplified selection model assuming that missingness of $\Delta$EPT is related to change in EPT before and after training. Therefore, the model

**Table 7**  *Summary statistics for the ACTIVE sub-sample*

| Variable | Training group | Control group | Difference |
|---|---|---|---|
| Sample size | 372 | 331 | |
| Missing rate | 17.2% | 30.5% | −13.3% |
| Age (years) | 73.25 (5.78) | 73.84 (6.28) | −0.59 |
| Education (years) | 13.71 (2.60) | 13.45 (2.87) | 0.26 |
| ΔEPT | −0.11 (3.27) | −0.28 (3.13) | 0.17 |

Note: Values in the parentheses are standard deviations.

can be specified as

$$\Delta \text{EPT}_i = \beta_0 + \beta_1 \text{Training}_i + \beta_2 \text{Age}_i + \beta_3 \text{Education}_i + e_i,$$

$$m_i \sim \text{Bernoulli}(p_i),$$

$$p_i = \Phi(\gamma_0 + \gamma_1 \Delta \text{EPT}_i).$$

In the two models, Training is a binary variable with 1 denoting that a participant is from the booster training group and 0 denoting the control group. The missingness indicator variable $m$ was created with 1 indicating that the score of ΔEPT is missing. Both models were estimated through the Bayesian method as discussed earlier. For priors, each regression coefficient ($\beta$) was given a normal distribution with mean 0 and variance $10^6$ and the variance of the residuals $\phi$ was given an inverse gamma distribution with both scale and shape parameters equal to $10^{-3}$. The results from the two models are summarized in Table 8.

First, the Geweke statistics show that the Markov chain for each model parameter converged to its marginal distribution because all Geweke statistics are in the range of −1 to 1. Second, for all model parameters, the ratios between the Monte Carlo error and the standard deviation are smaller than 5%. This indicates that the parameter estimates were accurate. Thus, we can make our inference based on the results in Table 8.

When missing data are assumed to be ignorable, results from the regular regression model show that booster training did not improve everyday problem solving ability of the elderly after controlling effects of age and education level. Age did not predict change in everyday problem solving ability, which is not consistent with aging literature (e.g., Finkel et al., 2003, Hedden and Gabrieli, 2004). However, analyzing the data as MNAR using the simplified selection model reveals a different picture. Training did have a positive effect in the change of everyday problem solving ability. Overall, participants in the training group on average had 0.739 (s.d.: 0.316; HPD: 0.110–1.351) more positive change than those in the control group after controlling effects of age and education level. In addition, age is negatively related to change in everyday problem solving ability (older adults had

**Table 8** *Bayesian parameter estimates from the regular regression model and the simplified selection model*

|                 |              | Estimate | s.d.  | MC/s.d. | HPD    |        | Geweke |
|-----------------|--------------|----------|-------|---------|--------|--------|--------|
| Regular model   | Intercept    | 1.816    | 1.925 | 0.036   | −2.011 | 5.506  | −0.124 |
|                 | Training     | 0.117    | 0.281 | 0.006   | −0.429 | 0.671  | −0.740 |
|                 | Age          | −0.041   | 0.024 | 0.035   | −0.087 | 0.006  | −0.018 |
|                 | Education    | 0.070    | 0.051 | 0.018   | −0.031 | 0.171  | 0.805  |
|                 | $\phi$       | 10.280   | 0.631 | 0.003   | 9.057  | 11.520 | 0.649  |
| Selection model | Intercept    | 1.384    | 2.083 | 0.040   | −2.843 | 5.273  | 0.483  |
|                 | Training     | 0.739    | 0.316 | 0.011   | 0.110  | 1.351  | −0.819 |
|                 | Age          | −0.057   | 0.026 | 0.039   | −0.105 | −0.005 | −0.426 |
|                 | Education    | 0.090    | 0.054 | 0.020   | −0.018 | 0.194  | −0.277 |
|                 | $\phi$       | 13.010   | 1.288 | 0.021   | 10.460 | 15.500 | −0.961 |
|                 | $\gamma_0$   | −1.312   | 0.230 | 0.028   | −1.766 | −0.866 | 0.994  |
|                 | $\gamma_1$ ($\Delta$EPT) | −0.256 | 0.058 | 0.028 | −0.366 | −0.142 | 0.865 |

Note. s.d.: standard deviation, can also be viewed as standard error from a frequentist's perspective. MC/s.d.: the ratio between Monte Carlo error and standard deviation of a parameter. HPD: highest posterior density credible interval.

more negative change than younger adults), which is consistent with previous findings (e.g., see the review by Hedden and Gabrieli, 2004).

With different assumptions on missing mechanisms and consequently different models, our data analysis led to different conclusions. Comparison between empirical results and aging literature indicates that it is more likely that missing data were MNAR in this example. From the estimate of $\gamma_1$ in the simplified selection model, $\Delta$EPT was negatively related to the missingness of itself. This means that participants with a lower $\Delta$EPT (less positive change or more negative change) were more likely to have missing data on post-test. In other words, if a participant expected that he/she would not gain much on EPT through training, he/she was more likely to miss the post-test. Note that we have found that for the training group, participants on average had 0.739 more positive change in their EPT scores than the control group after controlling the effects of age and education from the simplified selection model. Thus, participants in the training group would have a lower probability to have missing data than participants in the control group. Actually, from the empirical missing data rates of the ACTIVE study in Table 7, about 30.5% percent of participants have missing data in the control group and about 17.2% of participants have missing data in the booster training group. Therefore, the control group had a higher missing data rate, which means that relatively more lower $\Delta$EPTs were missed and relatively more higher $\Delta$EPTs were included in the analysis than the training group in terms of proportions. This explains why results from the two models are different.

# 6 Conclusion and discussion

To ease the choice of appropriate variables in explaining missing mechanisms in selection models, we examined and extended the simplified selection model used by Best et al. (1996) in which missingness depends solely on the missing variable itself. We first derived the full conditional posterior distributions for the simplified selection model using the data augmentation algorithm. Then, we conducted six simulation studies to evaluate the performance of the simplified selection model under a variety of conditions. Finally, we demonstrated the application of the simplified selection model through an example using real data from the ACTIVE study.

Our simulation results portrayed important features of the simplified selection model. Simulation study 1 showed that when MNAR data were analyzed as MAR data, parameter estimates were incorrect. When the simplified selection model was applied, parameter estimates were accurate. Simulation studies 2 and 3 demonstrated that the simplified selection model was insensitive to some choices of priors and link functions. For example, using either semi-conjugate priors or Jeffreys priors, model parameters were recovered equally well. Furthermore, although missing data were generated through a step function, both the cumulative normal distribution and logistic functions can be used to obtain correct regression parameter estimates.

MNAR data may be resulted from different situations other than the simple situation that missingness depends on variables themselves of interest. For example, missingness could be related to an auxiliary and unobserved variable. Simulation study 4 investigated such a scenario and found that the simplified selection model was still able to recover regression model parameters very well. Simulation study 5 looked into the situation where missingness in $\mathbf{y}$ depends on both $\mathbf{y}$ and $\mathbf{X}$. The results showed that the simplified selection model again performed well in this situation. Overall, simulation results showed that the simplified selection model is robust to the misspecification situations designed in our simulation studies and the simplified selection model can recover regression model parameters under both correctly specified situations and many misspecified situations.

Our ACTIVE data analysis provided a substantive example on how to apply the simplified model to analyze real data. From previous research, we knew that with the increase of age, there was a accelerated decline in cognitive ability (e.g., Finkel et al., 2003, Hedden and Gabrieli, 2004). However, the regular data analysis assuming ignorable missing data showed that age was not related to the change of EPT. This signaled that the missing mechanism here may not be ignorable. On the other hand, when the simplified selection model was applied, the negative relation between age and the change in EPT showed up. Furthermore, other results from the simplified selection model seemed also reasonable. It demonstrated that if a participant did not expect much help from booster training, she/he was more likely to miss a test. It should be noted that this is not a formal test for distinguishing

missing mechanisms. Thus, we suggest whenever a selection model is used to analyze missing data, data analysis based on a corresponding regular model should also be conducted and reported for comparison.

This study has its limitations. Admittedly, the regression model discussed in the current study is a relatively simple model. However, this model is very widely used practically. Especially, with such a model, we can disentangle the complexity of missing data analysis in a transparent way. For example, conditional posterior distributions of missing data are readily available and they clearly show the difference and connection between a regular model and a selection model. The method and strategy used in this study can be readily generalized to more complex models, for example, growth curve models and growth mixture models. By presenting the details of a simpler model, it is our hope that more future research on nonignorable missing data can be conducted for regression models and other sophisticated models.

Because of the focus of the current study, we did not formally discuss model fit and model selection when applying selection models. Potentially, the model fit can be evaluated through posterior predictive checking and the model selection can be conducted using Bayes factors (Gelman et al., 2003). However, the validity and reliability of posterior predictive checking and Bayes factors for model fit and model selection involving nonignorable missing data need careful evaluation.

## Appendix: WinBUGS codes for the selection model

```
## Model
model{
  for (i in 1:N){
    mu[i]<-b[1]+b[2]*x1[i]+b[3]*x2[i]
    y[i]~dnorm(mu[i], pre.phi)

    z[i]~dnorm(muz[i], 1)I(L[i],U[i])
    muz[i]<-b[4]+b[5]*y[i]

    L[i]<- -(1-m[i])*10000
    U[i]<- m[i]*10000
  }
  for (i in 1:5){
    b[i]~dnorm(0, 1.0E-6)
    Para[i]<-b[i]
  }
  pre.phi~dgamma(.001,.001)
  Para[6]<-1/pre.phi
```

```
}
## Starting values
list(b=c(0,0,0,0,0), pre.phi=1)
## Data are omitted for the sake of saving space
```

## Acknowledgments

## References

Albert, J. H. and Chib, S. (1993). Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association* **88**, 669–679. MR1224394

Best, N. G., Spiegelhalter, D. J., Thomas, A. A. and Brayne, C. E. (1996). Bayesian analysis of realistically complex models. *Journal of Royal Statistical Society, Ser. A* **159**, 323–342.

Box, G. E. P. and Tiao, G. C. (1973). *Bayesian Inference in Statistical Analysis*. New York, NY: Wiley. MR0418321

Casella, G. and George, E. I. (1992). Explaining the Gibbs sampler. *The American Statistican* **46**, 167–174. MR1183069

Congdon, P. (2003). *Applied Bayesian Modelling*. Chichester, UK: Wiley. MR1990543

Finkel, D., Reynolds, C. A., McArdle, J. J., Gatz, M. and Pedersen, N. L. (2003). Latent growth curve analyses of accelerating decline in cognitive abilities in late adulthood. *Developmental Psychology* **39**, 535–550.

Gelman, A., Carlin, J. B., Stern, H. S. and Rubin, D. B. (2003). *Bayesian Data Analysis*. London, UK: Chapman & Hall/CRC. MR1385925

Geman, S. S. and Geman, D. D. (1984). Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **6**, 721–741.

Geweke, J. J. (1992). Evaluating the accuracy of sampling-based approaches to calculating posterior moments. In *Bayesian Statistics* (J. M. Bernado, J. O. Berger, A. P. Dawid and A. F. M. Smith, eds.) **4**, 169–193. Oxford, UK: Clarendon Press. MR1380276

Graham, J. W. (2009). Missing data analysis: Making it work in the real world. *Annual Review of Psychology* **60**, 549–576.

Heckman, J. J. (1976). The common structure of statistical models of truncation, sample selection and limited dependent variables and a simple estimator for such models. *Annals of Economic and Social Measurement* **5**, 475–492.

Hedden, T. and Gabrieli, J. D. E. (2004). Insights into the ageing mind: A view from cognitive neuroscience. *Nature Reviews Neuroscience* **5**, 87–96.

Ibrahim, J. G. and Molenberghs, G. (2009). Missing data methods in longitudinal studies: A review. *Test* **18**, 1–43. MR2495958

Ibrahim, J. G., Chen, M.-H., Lipsitz, S. R. and Herring, A. H. (2006). Missing-data methods for generalized linear models: A comparative review. *Journal of the American Statistical Association* **100**, 332–346. MR2166072

Jeffreys, H. (1946). An invariant form for the prior probability in estimation problems. *Proceedings of the Royal Society of London. Series A, Mathematical and Physical Sciences* **186**, 453–461. MR0017504

Jobe, J. B., Smith, D. M., Ball, K., Tennstedt, S. L., Marsiske, M., Willis, S. L., Rebok, G. W., Morris, J. N., Helmers, K. F., Leveck, M. D. and Kleinman, K. (2001). Active: A cognitive intervention trial to promote independence in older adults. *Controlled Clinical Trials* **22**, 453–479.

Little, R. J. A. (1982). Models for nonresponse in sample surveys. *Journal of American Statistical Association* **77**, 237–250. MR0664675

Little, R. J. A. (1985). A note about model for selectivity bias. *Econometrica* **53**, 1469–1474.

Little, R. J. A. and Rubin, D. B. (2002). *Statistical Analysis with Missing Data*, 2nd ed. New York: Wiley-Interscience. MR1925014

Olsen, R. J. (1980). A least squares correction for selectivity bias. *Econometrica* **48**, 1815–1820. MR0617823

Schafer, J. L. (1997). *Analysis of Incomplete Multivariate Data*. London, UK: Chapman & Hall/CRC. MR1692799

Tang, G., Little, R. J. A. and Raghunathan, T. E. (2003). Analysis of multivariate missing data with nonignorable nonresponse. *Biometrika* **90** 747–764. MR2024755

Tanner, M. A. and Wong, W. H. (1987). The calculation of posterior distributions by data augmentation. *Journal of the American Statistical Association* **82** 528–540. MR0898357

Tennstedt, S. (2001). ACTIVE (advanced cognitive training for independent and vital elderly), 1999–2001 [United States] [Computer file]. ICPSR04248-v1. New England Research Institute [producer], Watertown, MA, 2001. Inter-university Consortium for Political and Social Research [distributor], Ann Arbor, MI, 2005-10-11.

Zhang, Z., McArdle, J. J., Wang, L. and Hamagami, F. (2008). A SAS interface for Bayesian analysis with WinBUGS. *Structural Equation Modeling* **15**, 705–728. MR2530373

Department of Psycholgy
University of Notre Dame
118 Haggar Hall
Notre Dame, Indiana 46545
USA
E-mail: zzhang4@nd.edu