

1       **A shared parameter mixture model for longitudinal income data with missing**  
2                               **responses and zero rounding**

3                               Francis K. C. Hui<sup>1\*</sup>, and Howard D. Bondell<sup>2</sup>

4                               *Australian National University and The University of Melbourne*

5                               **Summary**

The analysis of longitudinal income data is often made challenging for several reasons. For example, in a national Australian survey on income over time, a non-negligible proportion of responses are missing, and it is believed the missingness mechanism is non-ignorable. Also, there are a large number of reported zero incomes, some of which may be true zeros (corresponding to individuals who legitimately do not earn an income), while some may be false zeros (corresponding to individuals choosing to round their income to zero). We propose a new shared parameter mixture (SPM) model for analysing semicontinuous longitudinal income data, which addresses the two challenges of income non-response and zero rounding. This is accomplished by jointly modelling an individual's underlying income together with the probability of missingness and rounding to zero, where both probabilities are permitted to vary in a smooth manner with their underlying non-zero income. Applying the SPM model to the Australian income survey reveals that on average, older female individuals and individuals with a long term health condition are considerably less likely to earn an income, while income tended to be highest for male individuals on fixed-term/permanent job contracts between ages 50 to 60. Furthermore there is evidence of both zero rounding, and conditional on the assumed missingness mechanism, individuals with incomes at the higher and lower ends are more likely to not report their income.

6       **Key words:** income; missing not at random; mixed models; non-ignorable; zero inflation

7                               **1. Introduction**

8       Modelling how income distribution changes over time with an individual's physical  
9 and social circumstances is of major interest in economics, with important implications for  
10 understanding the relationship between income and public health issues such as poverty, and

---

\* Author to whom correspondence should be addressed.

<sup>1</sup> Research School of Finance, Actuarial Studies & Statistics, Australian National University, Acton, ACT 2601, Australia

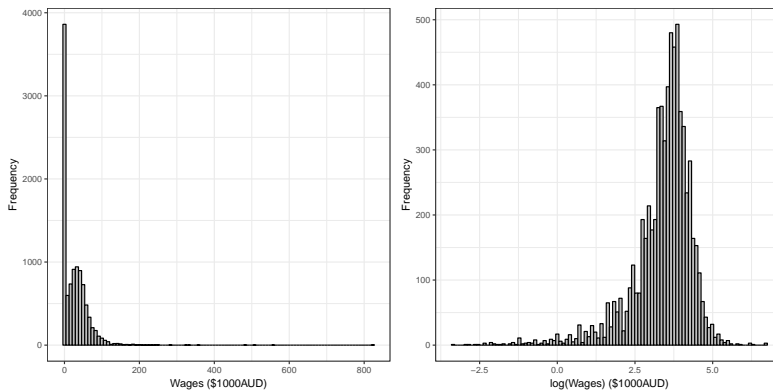
<sup>2</sup> School of Mathematics and Statistics, The University of Melbourne, Melbourne, VIC 3010, Australia  
Email: francis.hui@anu.edu.au

*Acknowledgment.* Both FKCH and HDB were supported Australian Research Council Discovery projects. Thanks to Nicole Watson for useful discussions relating to the HILDA survey. The HILDA Survey was initiated and is funded by the Australian Government Department of Social Services and is managed by the Melbourne Institute of Applied Economic and Social Research. The findings and views reported in this paper, however, are those of the authors.

11 poor physical and mental health (e.g., [Bechtel, Lordan & Rao 2012](#)). This article is motivated  
12 by the Household Income and Labour Dynamics in Australia (HILDA) survey ([Watson &](#)  
13 [Wooden 2012](#)), a national panel survey collected annually in Australia since 2001. One of  
14 the main goals of the HILDA survey is to achieve a greater understanding of how individual  
15 earnings have evolved over time in Australia, along with factors that drive income distribution  
16 and its evolution.

17 As is the case more generally, modelling income data in the HILDA survey is made  
18 challenging for two main reasons. First, a proportion of the income responses (approximately  
19 7.60%) are missing, and it is often more realistic to assume that non-response is missing not  
20 at random ([Rubin 1976](#)) as individuals with very low or very high income tend to refuse  
21 to report their income ([Riphahn & Serfling 2005](#); [Schraepfer 2006](#)). While the proportion  
22 of data missing is not as large as in other national surveys on income (e.g., the British  
23 Household Panel Survey, where it exceeds 15%, [Jenkins 2010](#)), it remains true that with  
24 such non-response we can not ignore the missing data mechanism because it invalidates  
25 many standard methods of inference on longitudinal data ([Little & Rubin 2014](#)). In the  
26 context of longitudinal income data, there has been some research on handling non-ignorable  
27 responses (e.g., the pattern mixture model of [Giusti & Little 2011](#)), although the vast majority  
28 of developments have instead focussed on assuming income is missing at random and thus  
29 ignorable missingness, largely driven by the easier means of analyses, and comparing various  
30 imputation methods (see [Headey & Wooden 2004](#); [Watson & Starick 2011](#), for examples  
31 with the HILDA survey). A second major challenge is that there are a substantial number  
32 of incomes (approximately 32.48%) reported as zero, as exemplified in [Figure 1](#). In other  
33 disciplines, such data comprising of one (or more) spikes at specific values along with a  
34 continuous distribution are known as semicontinuous data ([Bohning & Alfo 2016](#); [Liu](#)  
35 [et al. 2019](#)). While a proportion of these reported zero incomes may correspond to true  
36 zeros e.g., the individual is retired, some may be heaped or rounded zeros arising from an  
37 individual with very low but non-zero income choosing to round their income to zero (see  
38 for instance [Hanisch 2005](#), for known evidence of this phenomenon). If such a rounding  
39 mechanism is not accounted for i.e., all the zeros are assumed to be true zeros, then this can  
40 also invalidate inference on the likelihood of a person having a non-zero response. Recent  
41 work by [Zinn & Wrbach \(2016\)](#) and [Gross & Rendtel \(2016\)](#) among others have proposed  
42 heuristic approaches to first identify values to which incomes tend to be rounded, and then use  
43 kernel-type heaping functions to quantify the probability of an income value being rounded.  
44 However, both articles were solely focused on modelling the rounding process, while we  
45 are interested in modelling potential rounding to zero as part of an overall analysis for how  
46 income varies with an individual's attributes over time.

Figure 1. Histograms of all reported income data (left) and  $\log(\text{income})$  for reported non-zero income data only (right) from the HILDA survey. In the left panel, we can see a substantial proportion of incomes reported equal to zero, while the distribution of non-zero incomes in the right panel is unimodal and presents evidence of rounding at non-zero income values.



47 In this article, we propose a shared parameter mixture (SPM) model for analysing  
 48 semicontinuous longitudinal income data, which explicitly addresses the challenges of non-  
 49 ignorable responses and zero rounding in the HILDA survey. To our knowledge, our article  
 50 is the first in the applied statistics literature that attempts to overcome both challenges  
 51 simultaneously in a single model. This is despite the fact that both missing responses  
 52 and rounding tend to occur at the same time in many income surveys. The SPM model  
 53 is formulated hierarchically by first constructing a latent process relating an individual's  
 54 true mean income to their covariate values. This takes the form of a mixture model, with  
 55 one mixture component being a point mass at zero used to model the probability of an  
 56 individual earning an income, and the second mixture component being a positive continuous  
 57 distribution used to model the mean non-zero income over time, conditional on earning  
 58 an income. For modelling semicontinuous responses comprising a point mass at zero and  
 59 positive continuous data, mixture models are a popular although by no means the only  
 60 approach; see the special issue edited by [Bohning & Alfo \(2016\)](#) as well as the recent review  
 61 article by [Liu et al. \(2019\)](#). More importantly, we believe this mixture model is appropriate  
 62 for studying the true underlying income distribution in the HILDA survey, as it ensures that  
 63 the two processes of whether an individual actually earns an income and then how much  
 64 they earn are modelled explicitly, with potentially differing covariates driving the two mean  
 65 structures characterising these processes; see also [Bacci et al. \(2019\)](#).

66 Given the latent income process, the SPM model then builds in two components which  
 67 account for potential non-ignorable missingness and zero rounding. For the former, we  
 68 propose a shared parameter logistic regression model, such that the probability of missingness

69 depends on the value of their underlying mean non-zero income as well as other covariates.  
70 This approach is similar to shared parameter regression models, where random effects are  
71 used to drive both, and thus induce a dependence between, the responses and the missing  
72 data mechanism (Chapter 19, [Fitzmaurice et al. 2008](#)). However, a key difference is that rather  
73 than sharing the random effects, we share the conditional mean in the non-zero component  
74 of the underlying mixture model, such that the probability of non-response is allowed to vary  
75 in a smooth manner with underlying non-zero income. Apart from being more parsimonious,  
76 this approach also improves interpretability as we can differentiate between direct effects and  
77 indirect effects (acting via income) of covariates on the likelihood of missingness. Indeed,  
78 previous studies have shown that, aside from being more likely for lower and higher income  
79 earners, income non-response may also be directly affected by personal attributes such as age,  
80 which also affect the level of income ([Riphahn & Serfling 2005](#)). As a side note, we point  
81 out that approaches where the linear predictor (or function of it) is included in the model  
82 for missingness are sometimes referred to as ‘joint models’ ([Wulfsohn & Tsiatis 1997](#)).  
83 However in this article, we will use the terminology ‘shared parameter models’ as is done  
84 in [Rizopoulos, Verbeke & Molenberghs \(2008\)](#), [Tsonaka, Verbeke & Lesaffre \(2009\)](#) and  
85 [Creemers et al. \(2011\)](#), among others. Conditional on the responses not being missing, we  
86 then propose a second logistic regression model to account for potential zero rounding. This  
87 component of the SPM model may be also regarded as a shared parameter model, as the  
88 likelihood of rounding depends on the true income and thus indirectly on covariates driving  
89 income.

90 We estimate and perform inference on the SPM model using Bayesian v Chain  
91 Monte Carlo (MCMC) estimation via JAGS (Just Another Gibbs Sampler, [Plummer et al.](#)  
92 [2003](#)). Simulations demonstrate that the SPM model can successfully capture potential non-  
93 ignorable missingness and zero rounding in the income responses, while models that do  
94 not explicitly account for such processes may produce incorrect inference on the predictors  
95 driving the income distribution and its evolution. Applying the SPM model to the HILDA  
96 survey produced several interesting although not entirely surprising results: 1) older female  
97 individuals and individuals with a long term health condition are comparably less likely to  
98 earn an income; 2) males between ages 50–60 on a permanent or fixed-term job contract tend  
99 to earn the most income. Furthermore, income distribution depends on age in a nonlinear  
100 manner; 3) there is a strong association between the probability of zero rounding and a  
101 person’s true non-zero income. A sensitivity analysis also suggests that many of the parameter  
102 estimates are sensitive to the missing data assumptions. We provide template R code for  
103 fitting the SPM model, as well to reproduce the simulation and application in Supplementary  
104 Material C.

105

## 2. The shared parameter mixture model

106 Consider a longitudinal dataset comprising  $n$  independent individuals, where for  
 107 individual  $i = 1, \dots, n$  we let  $y_{ij, \text{obs}}$  denote their reported income at time point  $j = 1, \dots, n_i$ .  
 108 In the motivating HILDA survey, there are  $n = 1,693$  individuals who have reported incomes  
 109 ranging from  $n_i = 2$  to 10 time points, leading to a total of 10,650 observations. Along with  
 110 the responses, we observe a number of explanatory variables e.g., gender and age. For the  
 111 HILDA survey, all covariate values in the dataset are observed, while some of the  $y_{ij, \text{obs}}$ 's  
 112 are missing. In what follows, we shall use  $\mathbf{x}_{ij}$  and  $\mathbf{z}_{ij}$  to denote covariates to be included as  
 113 fixed effects and random effects, respectively. Unless stated otherwise, both covariate vectors  
 114 are assumed to contain an intercept term as their first elements. Moreover, as in common  
 115 in longitudinal mixed models, the elements of these vectors tend to overlap since random  
 116 slopes are generally also included as fixed slopes (Hui, Miller & Welsh 2017). We formulate  
 117 the SPM model by first describing the latent process governing an individual's true income,  
 118 before discussing the components mapping this latent process to the observed income.

119 We propose a two-component mixture model to quantify the latent process for an  
 120 individual's true income, denoted here as  $y_{ij, \text{true}}$ . Let  $g_{ij}$  be an indicator of whether individual  
 121  $i$  earns an income at time point  $j$ , such that  $g_{ij} = 1$ , if  $y_{ij, \text{true}} > 0$ , and 0 otherwise. Note  $g_{ij}$   
 122 is not fully observed, as we shall discuss later. Then we apply a logistic regression to model  
 123 the probability that an individual earns an income. Conditional on a vector of random effects  
 124  $\mathbf{b}_{i, \text{earn}}$ , we assume the indicators  $g_{ij}$  are independent Bernoulli random variables,

$$f(g_{ij} | \mathbf{b}_{i, \text{earn}}) = \text{Bern}(g_{ij}; \pi_{ij}); \quad \text{logit}(\pi_{ij}) = \mathbf{x}_{ij, \text{earn}}^\top \boldsymbol{\beta}_{\text{earn}} + \mathbf{z}_{ij, \text{earn}}^\top \mathbf{b}_{i, \text{earn}}, \quad (1)$$

125 where  $f(\cdot | \cdot)$  is used to denote conditional distributions,  $\boldsymbol{\beta}_{\text{earn}}$  are the fixed effects coefficients  
 126 and  $\text{Bern}(\cdot; \mu)$  denotes the Bernoulli distribution with probability of success  $\mu$ . The random  
 127 effects are assumed to be drawn from a multivariate normal distribution, such that  $f(\mathbf{b}_{i, \text{earn}}) =$   
 128  $\mathcal{N}(\mathbf{b}_{i, \text{earn}}; \mathbf{0}, \boldsymbol{\Sigma}_{\text{earn}})$  for some unstructured covariance matrix  $\boldsymbol{\Sigma}_{\text{earn}}$ , where  $\mathcal{N}(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$  denotes  
 129 the multivariate normal density with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ . Second,  
 130 conditional on earning an income i.e.,  $g_{ij} = 1$ , we model the individual's true non-zero  
 131 income via a second regression model,

$$f(y_{ij, \text{true}} | g_{ij} = 1, \mathbf{b}_{i, \text{nonzero}}) = \mathcal{F}_+(y_{ij, \text{true}}; \eta_{ij}, \phi);$$

$$\eta_{ij} = \mathbf{x}_{ij, \text{nonzero}}^\top \boldsymbol{\beta}_{\text{nonzero}} + \mathbf{z}_{ij, \text{nonzero}}^\top \mathbf{b}_{i, \text{nonzero}}, \quad (2)$$

132 where  $\mathcal{F}_+(\cdot; \eta, \phi)$  denotes a distribution defined on the positive real line and characterised by  
 133 parameters  $(\eta, \phi)$  as defined below, and  $f(\mathbf{b}_{i, \text{nonzero}}) = \mathcal{N}(\mathbf{b}_{i, \text{nonzero}}; \mathbf{0}, \boldsymbol{\Sigma}_{\text{nonzero}})$ . We consider  
 134 three choices for  $\mathcal{F}_+(\cdot; \eta, \phi)$ , although we acknowledge more complex distributions for

135 income are possible (e.g., [Caldern-Ojeda, Azpitarte & Gmez-Dniz 2016](#)): 1) a gamma  
 136 distribution with mean  $\mu = \exp(\eta)$  and shape parameter  $\phi$ ; 2) an exponential distribution  
 137 with mean  $\mu = \exp(\eta)$ . This is a special case of the gamma distribution with shape parameter  
 138  $\phi = 1$ ; 3) a lognormal distribution with location  $\eta$  and variance  $\phi$  parameters on the log  
 139 scale, such that the mean of the distribution is equal to  $\mu = \exp(\eta + 2^{-1}\phi)$ . In all three  
 140 cases, we note the use of an exponential link function to reflect the notion that wealth is  
 141 driven by a series of multiplicative processes. Let  $\Psi_{\text{earn}} = (\beta_{\text{earn}}, \text{vech}(\Sigma_{\text{earn}}))$  and  $\Psi_{\text{nonzero}} =$   
 142  $(\beta_{\text{nonzero}}, \text{vech}(\Sigma_{\text{nonzero}}), \phi)$  denote the parameters describing the probability of earning an  
 143 income and the distribution of non-zero incomes respectively. Furthermore, by marginalising  
 144 (2) with respect to  $g_{ij}$  and (1), the formulation of the mixture model becomes clear as we see  
 145 that  $f(y_{ij,\text{true}} | \mathbf{b}_{i,\text{true}}, \mathbf{b}_{i,\text{nonzero}}) = (1 - \pi_{ij})\mathbb{I}_0 + \pi_{ij}\mathcal{F}_+(y_{ij,\text{true}}; \eta_{ij}, \phi)$ , where  $\mathbb{I}_0$  is the Dirac  
 146 delta function. The distribution of true income is thus a mixture of a spike at zero and a  
 147 positive continuous distribution, and as reviewed in Section 1 such a model is commonly  
 148 applied for analysing semicontinuous data.

149 We next formulate a model for the missing responses. Let  $r_{ij}$  denote the missing data  
 150 indicators, such that  $r_{ij} = 1$  implies  $y_{ij,\text{obs}}$  is missing at time point  $j$  for individual  $i$ , and  
 151  $r_{ij} = 0$  otherwise. Also, let  $\mu_{ij} = E(y_{ij,\text{true}} | g_{ij} = 1, \mathbf{b}_{i,\text{nonzero}})$  denote the conditional mean  
 152 non-zero income based on (2), e.g.,  $\mu_{ij} = \exp(\eta_{ij} + 2^{-1}\phi)$  if  $y_{ij,\text{true}}$  is assumed to follow a  
 153 lognormal distribution. To allow for income missingness mechanism being potentially non-  
 154 ignorable, we allow the probability of missingness to vary as a function of  $\mu_{ij}$  (conditional on  
 155 the individual earning an income) as well as with a set of fixed and random effect covariates.  
 156 That is, given a vector of random effects  $\mathbf{b}_{i,\text{miss}}$ , we assume the  $r_{ij}$ 's are independent Bernoulli  
 157 random variables,

$$f(r_{ij} | g_{ij}, \mathbf{b}_{i,\text{nonzero}}, \mathbf{b}_{i,\text{miss}}) = \text{Bern}(r_{ij}; \rho_{ij});$$

$$\text{logit}(\rho_{ij}) = \mathbf{x}_{ij,\text{miss}}^\top \boldsymbol{\beta}_{\text{miss}} + \mathbf{z}_{ij,\text{miss}}^\top \mathbf{b}_{i,\text{miss}} + g_{ij} s_{\text{miss}}(\mu_{ij}), \quad (3)$$

158 where  $s_{\text{miss}}(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is some smooth function satisfying  $s_{\text{miss}}(0) = 0$ , and  $f(\mathbf{b}_{i,\text{miss}}) =$   
 159  $\mathcal{N}(\mathbf{b}_{i,\text{miss}}; \mathbf{0}, \Sigma_{\text{miss}})$ . Conditional on earning an income ( $g_{ij} = 1$ ), we see that due to shared  
 160 random effects between  $\mu_{ij}$  and  $\rho_{ij}$ , there is a resulting marginal dependence between the  
 161 missing data indicators and the missing observations themselves. Consider an individual who  
 162 always earn an income, such that  $\mathbf{g}_i = (g_{i1}, \dots, g_{in_i}) = \mathbf{1}$ . If we let  $\mathbf{y}_{i,\text{true}}$  and  $\mathbf{r}_i$  denote the  
 163  $n_i$ -vectors of true income and missing response indicators for individual  $i$ , respectively, then

$$f(\mathbf{r}_i | \mathbf{y}_{i,\text{true}}) = \int \prod_{j=1}^{n_i} f(r_{ij} | \mathbf{b}_{i,\text{nonzero}}, \mathbf{b}_{i,\text{miss}}, \mathbf{y}_{i,\text{true}}) f(\mathbf{b}_{i,\text{nonzero}} | \mathbf{y}_{i,\text{true}}) f(\mathbf{b}_{i,\text{miss}}) d\mathbf{b}_{i,\text{nonzero}} d\mathbf{b}_{i,\text{miss}},$$

164 where

$$f(r_{ij} | \mathbf{b}_{i,\text{nonzero}}, \mathbf{b}_{i,\text{miss}}, \mathbf{y}_{i,\text{true}}) = \text{Bern} \left[ r_{ij}; \text{logit}^{-1} \left\{ \mathbf{x}_{ij,\text{miss}}^\top \boldsymbol{\beta}_{\text{miss}} + \mathbf{z}_{ij,\text{miss}}^\top \mathbf{b}_{i,\text{miss}} + s_{\text{miss}}(\mu_{ij}) \right\} \right].$$

165 Therefore the dependence of the missing indicators on underlying income is made  
 166 explicit through the dependence on the posterior distribution of the random effects,  
 167  $f(\mathbf{b}_{i,\text{nonzero}} | \mathbf{y}_{i,\text{true}})$ . An analogous dependence result can be obtained for individuals who vary  
 168 between earning and not earning an income i.e,  $\mathbf{g}_i$  contains both zeros and ones.

169 Unlike standard applications of shared parameter models, where the the random effects  
 170  $\mathbf{b}_{i,\text{nonzero}}$  are included directly in the linear predictor (e.g., [Fitzmaurice et al. 2008](#)), the  
 171 probability of missingness in the SPM model is dependent on the mean non-zero income  
 172  $\mu_{ij}$ . The motivation behind this (at least with the current parameterisation) is to distinguish  
 173 between the indirect effects of covariates on missingness occurring through the responses,  
 174 and the direct effects of covariates on the probability of missingness. For example, if an  
 175 individual earns a very large income then they may choose not to report it due to its large  
 176 value, and since income typically has a significant relationship with age (as is the case in the  
 177 HILDA survey; see Section 5), then in such case age has an *indirect* effect on the probability  
 178 of missingness through  $\mu_{ij}$ . On the other hand, the probability of missingness may also be  
 179 directly related to age e.g., as respondents tend to become more familiar with the survey over  
 180 time and thus are more comfortable with reporting their income, in which case age has a direct  
 181 effect as captured by the  $\boldsymbol{\beta}_{\text{miss}}$  and  $\mathbf{b}_{i,\text{miss}}$ . Note also that by allowing the dependence to be on  
 182  $\mu_{ij}$  rather than on  $\mathbf{b}_{i,\text{nonzero}}$ , we require only a small number of parameters to capture the non-  
 183 ignorable missingness in (3), depending on the choice of  $s_{\text{miss}}(\cdot)$ . Given the large number  
 184 of parameters already present in the SPM model, the ability to capture a non-ignorable  
 185 missingness mechanism parsimoniously yet flexibly is an advantage.

186 It is worth discussing in more detail how (3) allows for and transitions between the  
 187 cases of  $g_{ij} = 0$  and  $g_{ij} = 1$  in a smooth manner. If  $g_{ij} = 0$ , then an individual does not earn  
 188 income and the probability of non-response can be considered as being depending upon the  
 189 underlying, zero income. If  $g_{ij} = 1$ , then an individual earns income and the probability of  
 190 non-response is allowed to vary with the underlying income. However, the former may also  
 191 be interpreted as a limiting case of the latter when  $\mu_{ij} \rightarrow 0$ , that is, we suppose an individual  
 192 earns an income but consider how the probability of non-response varies as their income  
 193 tends to zero. To allow this transition to occur in a smooth manner, we constrain  $s_{\text{miss}}(0) = 0$ .  
 194 Moreover, we see that the covariates  $(\mathbf{x}_{ij,\text{miss}}, \mathbf{z}_{ij,\text{miss}})$  in (3) can not only be interpreted as  
 195 direct effects on the probability of missingness (as discussed), but are in fact the only effects  
 196 in the limiting case of  $\mu_{ij} = 0$ .

197 Regarding the choice of  $s_{\text{miss}}(\cdot)$ , while it is possible to consider more data-driven  
 198 smoothing functions, for ease of interpretation (and likely the inability to estimate such  
 199 terms effectively without much missing data) we choose at most to use a quadratic form  
 200  $s_{\text{miss}}(\mu) = \mu\theta_1 + \mu^2\theta_2$ , where  $\theta = (\theta_1, \theta_2)$ . The quadratic form also means the probability of  
 201 missingness may vary nonlinearly as a function of the mean response, even if the mean itself  
 202 only varies linearly in time. For instance, an individual's income may increase monotonically  
 203 over time, but by constructing  $s_{\text{miss}}(\cdot)$  as a quadratic function then the probability of non-  
 204 response is (say) higher at earlier and later time points in the survey, when the individual  
 205 has relatively low and high income, respectively. Also, at a single time point it allows for  
 206 low income and high income earners to have higher propensity of non-response relative  
 207 to average income earners. Let  $\Psi_{\text{miss}} = (\beta_{\text{miss}}, \text{vech}(\Sigma_{\text{miss}}), \theta)$  denote the parameter vector  
 208 corresponding to the component of the SPM model characterising the missing data process. In  
 209 our application to the HILDA survey, we conduct a sensitivity analysis by comparing different  
 210 choices for  $s_{\text{miss}}(\cdot)$ , including the case of  $s_{\text{miss}}(\cdot) = 0$ , and its impact on the parameter  
 211 estimates in the SPM model.

212 We remark that making use of random effects is by no means the only approach  
 213 for modelling the serial correlation between the indicators of missing income; alternative  
 214 approaches such as assuming a Markovian process to account for dependence between  
 215 neighbouring  $r_{ij}$  (and indeed the  $g_{ij}$  for that matter) are possible; see Chapter 18, [Fitzmaurice  
 216 et al. \(2008\)](#). Our main motivation for using random effects to account for the temporal  
 217 correlation is precisely because it provides a relatively straightforward mechanism to  
 218 additionally induce correlation between the missingness indicators and the observed income,  
 219 and indeed also with indicator of zero rounding as we see below; see also [Bacci & Bartolucci  
 220 \(2015\)](#).

221 The last component of the SPM model involves constructing a logistic regression model  
 222 for the probability that an individual rounds their true income to zero, conditional on earning  
 223 an income. Let  $h_{ij} = 1$  if  $y_{ij,\text{true}}$  is rounded to zero and 0 otherwise. As with  $g_{ij}$ , note that  
 224  $h_{ij}$  is not fully observed. We model the  $h_{ij}$ 's as conditionally independent Bernoulli random  
 225 variables,

$$f(h_{ij}|g_{ij} = 1, \mathbf{b}_{i,\text{nonzero}}) = \text{Bern}(h_{ij}; \nu_{ij}); \quad \text{logit}(\nu_{ij}) = s_{\text{round}}(\mu_{ij}), \quad (4)$$

226 such that  $\nu_{ij}$  is the probability of rounding to zero, and  $s_{\text{round}}(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is a smooth function.  
 227 Analogous to (3), we allow the probability of zero rounding to vary with the conditional  
 228 mean non-zero income  $\mu_{ij}$ . Moreover as this is now shared between the models governing  
 229  $y_{ij,\text{true}}$ ,  $r_{ij}$  and  $h_{ij}$ , it induces a correlation between all three observed quantities and causing  
 230 them to co-vary with each other. Note we choose not to include covariates in the linear



231 predictor in (4), as we argue that the propensity of an individual rounding their income to zero  
 232 depends solely on the proximity of that income to zero (analogous to Zinn & Wrbach 2016;  
 233 Gross & Rendtel 2016). Regarding the form of  $s_{\text{round}}(\cdot)$ , we only consider a simple linear  
 234 form  $s_{\text{round}}(\mu) = \tau_0 + \mu\tau_1$  where  $\tau = (\tau_0, \tau_1)$ , with a negative value of  $\tau_1$  meaning that the  
 235 probability of rounding to zero decreases as the true income increases. Let  $\Psi_{\text{round}} = \tau$  denote  
 236 the parameter vector corresponding to the component of the SPM model governing the zero  
 237 rounding process.

238 It is important to acknowledge that we do not incorporate a mechanism to deal with  
 239 rounding or heaping at non-zero income values, even though there is evidence of this in  
 240 our motivating HILDA survey (Figure 1, right panel). Our reasoning behind this is that, in  
 241 the given dataset, the fraction of tied data points at any one *non-zero* income value is at  
 242 most 1.60% of the total number of observations (170 out of 10,650 observations occurring  
 243 at \$40,000 AUD), and as in typical practice, income data, although technically discrete, are  
 244 modelled as approximately continuous. This contrasts to the heaping at zero, which accounts  
 245 for approximately 32.5% of the total sample size. Of course, it is possible to apply data-  
 246 driven and nonparametric methods to both select the rounding points and model the rounding  
 247 process (e.g., Gross & Rendtel 2016). However we believe that, apart from both the estimation  
 248 procedure being considerably more complicated as a result, inference on the non-zero income  
 249 component of the SPM model will not be substantially affected by this.

## 250 2.1. Special cases

251 The SPM model includes two special cases of interest. First, if there is no missing data  
 252 then (3) can be omitted with  $r_{ij} = 1$  for all  $i$  and  $j$ , and the zero rounding process given by  
 253 (2) becomes the sole process affecting  $y_{ij,\text{true}}$ . Second, if there is no zero rounding then (4)  
 254 can be omitted with  $h_{ij} = 1$  for all  $i$  and  $j$ . Note if there were no missing data also then it  
 255 would be possible to treat the two mixture components defined by (1) and (2), separately.  
 256 That is, estimation and inference for the logistic regression model with parameters  $\Psi_{\text{earn}}$   
 257 could be to done separately to the regression model for non-zero income with parameters  
 258  $\Psi_{\text{nonzero}}$ . But with missing responses separate estimation is not possible due to the shared  
 259 missing data mechanism driving both components of the mixture model. More broadly, in  
 260 the full SPM model the effect of zero rounding is to induce a dependence between the two  
 261 mixture components of (1) and (2) at the observation level. That is, observations from the  
 262 non-zero income component can be moved into the zero income component, leading to cases  
 263 of  $y_{ij,\text{obs}} = 0$  even though  $y_{ij,\text{true}} > 0$ .

264 With longitudinal income data, neither simplification discussed above may be  
 265 appropriate: there are almost always cases of non-response in income surveys as individuals

266 refuse to report their income (see also the sensitivity analysis in Supplementary Material B.4),  
267 while zero rounding and indeed rounding to other non-zero amounts is a common occurrence  
268 as seen in Figure 1. In Section 4, we show that if there is zero rounding and/or non-ignorable  
269 missingness and this is not accounted for in the analysis, then it can lead to erroneous  
270 inference for one or more components of the SPM model. Furthermore, in Section 5, we  
271 will see that the HILDA survey presents evidence of both features being present.

### 272 3. Estimation

273 The full formulation of the SPM model can be written as follows. For individual  
274  $i = 1, \dots, n$  and time point  $j = 1, \dots, n_i$ , denote the vector of all responses as  
275  $(y_{ij,\text{obs}}, r_{ij}, g_{ij}, h_{ij})$ , i.e., the observed income, the missingness indicator, the indicator of  
276 whether the true income is positive, and the indicator of whether the true income has been  
277 rounded to zero. Apart from  $y_{ij,\text{obs}}$ , both  $g_{ij}$  and  $h_{ij}$  are only partly observed: if  $y_{ij,\text{obs}} > 0$   
278 then we know  $g_{ij} = 1$  and  $h_{ij} = 0$  i.e., the individual truly earns an income and chose not to  
279 round at that time point. If  $y_{ij,\text{obs}} = 0$  then both  $g_{ij}$  and  $h_{ij}$  are missing since we do not know  
280 whether the individual actually did not earn an income, or did earn an income but chose to  
281 round their income to zero. Finally, if  $r_{ij} = 1$  then  $y_{ij,\text{obs}}$  is missing, and we also set  $g_{ij}$  and

282  $h_{ij}$  to be missing. The SPM model is formulated as follows:

$$\begin{aligned}
f(\mathbf{b}_{i,\text{earn}}) &= \mathcal{N}(\mathbf{b}_{i,\text{earn}}; \mathbf{0}, \Sigma_{\text{earn}}) \\
f(\mathbf{b}_{i,\text{nonzero}}) &= \mathcal{N}(\mathbf{b}_{i,\text{nonzero}}; \mathbf{0}, \Sigma_{\text{nonzero}}) \\
f(\mathbf{b}_{i,\text{miss}}) &= \mathcal{N}(\mathbf{b}_{i,\text{miss}}; \mathbf{0}, \Sigma_{\text{miss}}) \\
f(g_{ij} | \mathbf{b}_{i,\text{earn}}) &= \text{Bern}(g_{ij}; \pi_{ij}); \\
&\text{where } \text{logit}(\pi_{ij}) = \mathbf{x}_{ij,\text{earn}}^\top \boldsymbol{\beta}_{\text{earn}} + \mathbf{z}_{ij,\text{earn}}^\top \mathbf{b}_{i,\text{earn}}, \\
f(y_{ij,\text{true}} = 0 | g_{ij} = 0) &= 1 \\
f(y_{ij,\text{true}} | g_{ij} = 1, \mathbf{b}_{i,\text{nonzero}}) &= \mathcal{F}_+(y_{ij,\text{true}}; \eta_{ij}, \phi); \\
&\text{where } \eta_{ij} = \mathbf{x}_{ij,\text{nonzero}}^\top \boldsymbol{\beta}_{\text{nonzero}} + \mathbf{z}_{ij,\text{nonzero}}^\top \mathbf{b}_{i,\text{nonzero}}, \\
f(r_{ij} | g_{ij}, \mathbf{b}_{i,\text{nonzero}}, \mathbf{b}_{i,\text{miss}}) &= \text{Bern}(r_{ij}; \rho_{ij}); \\
&\text{where } \text{logit}(\rho_{ij}) = \mathbf{x}_{ij,\text{miss}}^\top \boldsymbol{\beta}_{\text{miss}} + \mathbf{z}_{ij,\text{miss}}^\top \mathbf{b}_{i,\text{miss}} + g_{ij} s_{\text{miss}}(\mu_{ij}), \\
f(h_{ij} | g_{ij} = 1, \mathbf{b}_{i,\text{nonzero}}) &= \text{Bern}(h_{ij}; \nu_{ij}); \\
&\text{where } \text{logit}(\nu_{ij}) = s_{\text{round}}(\mu_{ij}), \\
f(y_{ij,\text{obs}} = 0 | g_{ij} = 0, r_{ij} = 0) &= f(y_{ij,\text{true}} = 0 | g_{ij} = 0) = 1 \\
f(y_{ij,\text{obs}} = 0 | g_{ij} = 1, r_{ij} = 0, h_{ij} = 1) &= 1 \\
f(y_{ij,\text{obs}} | g_{ij} = 1, r_{ij} = 0, h_{ij} = 0) &= f(y_{ij,\text{true}} | g_{ij} = 1, \mathbf{b}_{i,\text{nonzero}}),
\end{aligned}$$

283 and  $y_{ij,\text{obs}}$  is missing if  $r_{ij} = 1$ . To clarify the final three lines of the SPM model, note we have  
284 the following conditional distributions for reported income: 1)  $f(y_{ij,\text{obs}} = 0 | g_{ij} = 0, r_{ij} =$   
285  $0) = f(y_{ij,\text{true}} = 0 | g_{ij} = 0) = 1$ , which has the interpretation that if an individual does not  
286 earn an income and choose to report it, then we observe a zero income; 2)  $f(y_{ij,\text{obs}} = 0 | g_{ij} =$   
287  $1, r_{ij} = 0, h_{ij} = 1) = 1$ , which has the interpretation that if an individual earns an income,  
288 chooses to report it but rounds to zero, then we observe a zero income; 3)  $f(y_{ij,\text{obs}} | g_{ij} =$   
289  $1, r_{ij} = 0, h_{ij} = 0) = f(y_{ij,\text{true}} | g_{ij} = 1, \mathbf{b}_{i,\text{nonzero}}) = \mathcal{F}_+(y_{ij,\text{true}}; \mu_{ij}, \phi)$ , which means that  
290 if individual earns an income, chooses to report and does not round, then the distribution  
291 of the observed income is the same as the distribution of the true income  $y_{ij,\text{true}}$ .

292 The hierarchical nature of SPM model facilitates using Bayesian estimation via MCMC  
293 sampling, which we implemented using JAGS and the R package `runjags` (Denwood 2016)  
294 to jointly estimate all parameters (and random effects) in the model. As prior distributions  
295 for both our simulation study and application to the HILDA survey, for all the fixed effect  
296 coefficients  $\boldsymbol{\beta}_{\text{earn}}, \boldsymbol{\beta}_{\text{nonzero}}, \boldsymbol{\beta}_{\text{miss}}$  and parameters  $\boldsymbol{\theta}$  and  $\boldsymbol{\tau}$  characterising  $s_{\text{miss}}(\cdot)$  and  $s_{\text{round}}(\cdot)$ ,  
297 respectively, we assigned independent normal priors to each element with mean zero and  
298 variance equal to ten, noting all covariates entered into the SPM model were standardised

299 to have mean zero and variance one. For the nuisance parameter  $\phi$  in the gamma and  
 300 lognormal distributions, we assigned a half-normal prior with scale parameter equal to ten.  
 301 Our motivation for using such weakly informative priors for these parameters was to stabilise  
 302 the MCMC sampling but also reflect the prior belief that these coefficients were unlikely to  
 303 exceed the range  $\pm 10$  (e.g., [Gelman et al. 2008](#)). In both the simulations and application, we  
 304 did test the use of more diffuse prior distributions e.g., setting variance/scale parameters to  
 305 50 instead of 10, and found that results tended to be similar overall but that convergence of  
 306 the MCMC sampling took longer. Finally, for the three random effects covariance matrices  
 307  $\Sigma_{\text{earn}}$ ,  $\Sigma_{\text{nonzero}}$  and  $\Sigma_{\text{miss}}$ , we assigned independent inverse Wishart priors with the degrees of  
 308 freedom set equal to the dimension of the corresponding covariance matrix plus one, and the  
 309 scale matrix set equal to the identity matrix in the case of the simulation study.

310 We ran three MCMC chains, with each chain consisting for 20,000 burn-in iterations  
 311 followed by 50,000 additional samples with a thinning period of 25. We assessed convergence  
 312 based on trace plots and the Gelman-Rubin convergence statistic ([Gelman & Rubin 1992](#))  
 313 and, if required (although this was rarely needed), extended the MCMC sampling until the  
 314 convergence statistic was below 1.1 for all parameters in the model. Without extension, this  
 315 lead to a total of 6,000 MCMC samples for analysis.

#### 316 4. Simulation study

317 We conducted an empirical study to assess the performance of the SPM model, and  
 318 to illustrate how the failure to account for possible zero rounding and/or non-ignorable  
 319 missingness can lead to erroneous inference on various aspects of the longitudinal income  
 320 data. We used the SPM model fitted to the HILDA survey in Section 5 as the basis for  
 321 constructing a true SPM model as follows. First, we included only fixed effect covariates for  
 322 which the 95% highest posterior density interval of the corresponding regression coefficient  
 323 from the fitted SPM model did not contain zero. Next, we considered only a simple  
 324 linear relationship for the fixed effect of age in the simulation, rather than the non-linear  
 325 relationship modelled based on spline functions in the real application. This allowed us to  
 326 more straightforwardly assess the performance of competing models in terms of inference  
 327 on the effect of age. Based on these two modifications, this led to six coefficients in  $\beta_{\text{earn}}$ ,  
 328 five coefficients in  $\beta_{\text{nonzero}}$ , and four coefficients in  $\beta_{\text{miss}}$ . Aside from these two changes, the  
 329 other components of the true SPM model inherited the same form as the model fitted to the  
 330 HILDA survey in Section 5, with the true parameter values taken as the posterior median  
 331 estimate. In particular: 1) the same vector of two random effect terms (intercept and random  
 332 slope for age) was used in  $z_{ij,\text{earn}}$ ,  $z_{ij,\text{nonzero}}$ , and  $z_{ij,\text{miss}}$ ; 2) the non-zero income distribution  
 333  $\mathcal{F}_+(\cdot; \mu, \phi)$  was set to the lognormal distribution; 3) a quadratic form was used for  $s_{\text{miss}}(\cdot)$ ;

334 4) a linear form was used for  $s_{\text{round}}(\cdot)$ . We provide more details and a table of all the true  
 335 parameter values in Supplementary Material A.

336 For the simulation study, we considered a subsample of individuals from the HILDA  
 337 survey. This was obtained by randomly selecting  $n = 200$  individuals from the HILDA survey  
 338 who had  $n_i \geq 5$  measurements, and using only the covariates from this subsample. This  
 339 led to a total of 1626 observations in each simulated dataset. We simulated 400 datasets in  
 340 total. With the above simulation design and true parameter values, on average each simulated  
 341 dataset has around 6-8% missing income responses and between 30 to 40% of the observed  
 342 income values were zero; both of these features were consistent with what was seen in the  
 343 actual HILDA dataset.

344 For each simulated dataset, we fitted four competing models: 1) the SPM model in  
 345 Section 3 (SPM model); 2) a model with no zero rounding component (No rounding); 3) a  
 346 SPM model with  $s_{\text{miss}}(\cdot)$  omitted from the missingness component (No smooth); 4) a model  
 347 where the missing data was not modelled at all (No missing). Under this model and with  
 348 Bayesian MCMC sampling using JAGS, the missing data income responses are thus assumed  
 349 to be missing at random and multiple imputation performed based on the relevant predictive  
 350 distribution. This contrasts with Model 3, which also assumes the income responses are  
 351 missing at random but forms an explicit regression model for the probability of missingness as  
 352 a function of the observed covariates. We used the same set of prior distributions and sampling  
 353 scheme for all four models; see Section 3. Note some models do not involve parameters for  
 354 some components.

355 We considered four measures to assess performance of the competing models: 1)  
 356 relative bias, defined as  $400^{-1} \sum_{t=1}^{400} |\Psi_0|^{-1} (\hat{\Psi}^t - \Psi_0)$  where  $\hat{\Psi}^b$  generically denotes  
 357 the posterior median estimate for some parameter  $\Psi$  in simulated dataset  $t$ , and  $\Psi_0$   
 358 denotes the corresponding true value; 2) relative mean squared error (MSE), defined as  
 359  $400^{-1} \sum_{t=1}^{400} \Psi_0^{-2} (\hat{\Psi}^t - \Psi_0)^2$ ; 3) coverage probability, defined as the number of simulated  
 360 datasets out of 400 where the 95% highest posterior density (HPD) intervals contained the  
 361 true parameter; 4) the mean interval score (Gneiting & Raftery 2007). Given a  $100(1 - \alpha)\%$   
 362 HPD interval for a parameter  $\Psi$ , denoted as  $(l_\Psi, u_\Psi)$ , the interval score is calculated as  
 363  $(u_\Psi - l_\Psi) + 2\alpha^{-1} \times \mathbf{1}(l_\Psi < \Psi_0 < u_\Psi)$  where  $\mathbf{1}(\cdot)$  is the indicator function. A lower mean  
 364 interval score corresponds to better overall performance of the constructed HPD intervals i.e.,  
 365 the interval often contains the true parameter and its width is comparably small.

## 366 4.1. Results

367 It was not surprising to observe that Model 1, given it is the true model, performed  
 368 best overall across the four components of the SPM model (Figures 2 and 3): for most

369 parameters it had the smallest relative bias and relative MSE, the coverage probability of the  
 370 HPD intervals were always close to the nominal 95% coverage level, and the mean interval  
 371 score was almost always the lowest of the four models. The most substantial impact of not  
 372 accounting for zero rounding (Model 2; No rounding) was on inference for the probability  
 373 of an individual earning an income, where several of the parameters in  $\Psi_{\text{earn}}$  incurred  
 374 comparably high relative bias, coverage probabilities substantially below the nominal level,  
 375 and the corresponding mean intervals scores being much higher than models that accounted  
 376 for zero rounding (Figure 2, left column). Inference on the probability of non-response was  
 377 also affected, but to a lesser extent (Figure 3, left column). Incorrect inference on modelling  
 378 the probability of an individual earning an income was not surprising given that Model 2  
 379 falsely assumes that all observed zeros in the data correspond to truly zero incomes. Turning  
 380 to Model 3 (No smooth), the failure to model the non-ignorable missingness mechanism  
 381 impacted inference on the probability of missingness the most (Figures 3, left column).  
 382 This is expected given that, by not explicitly separating the direct and indirect effects of  
 383 the covariates, the direct estimates will erroneously incorporate both effects.

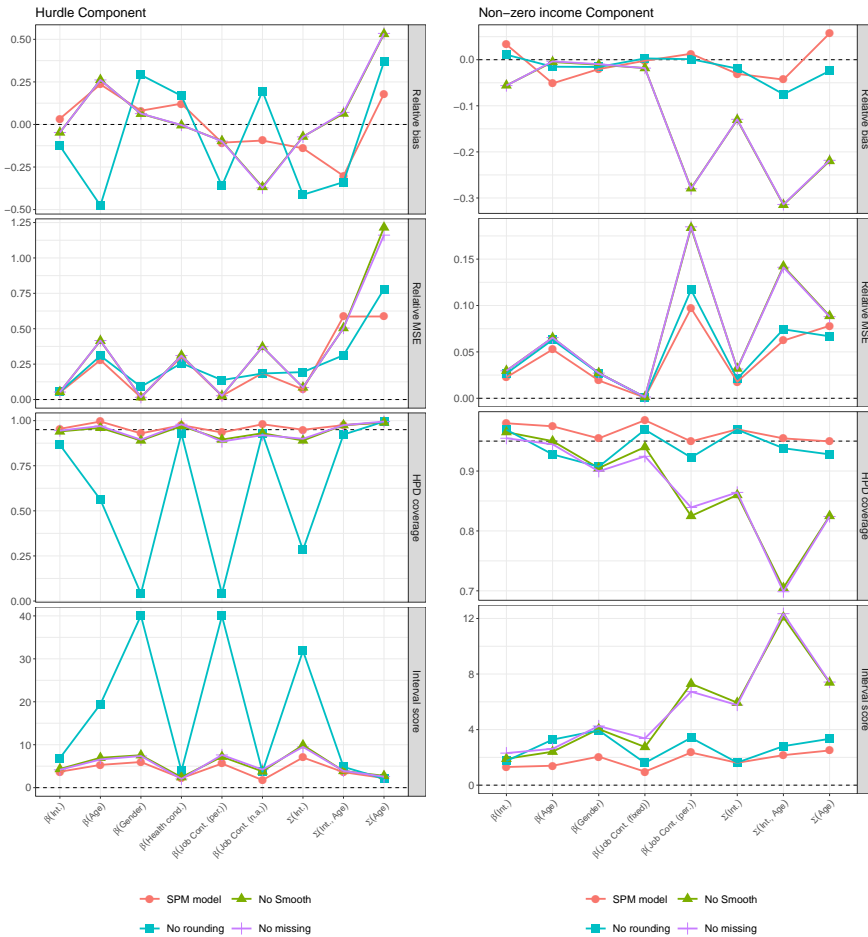
384 Inference on the modelling non-zero income distribution and the probability of zero  
 385 rounding was also negatively impacted to some extent, and this is to expected to be a  
 386 consequence of the failure to account for the shared random effects (and thus the correlation)  
 387 between these component and the missing response component. Finally, given that Model 4  
 388 (No missing) also assumes the income is missing at random but does not explicitly model the  
 389 probability of missingness as a function of the observed covariates, then it was not surprising  
 390 to see the results from this model were rather similar to those of Model 3. In particular, both  
 391 inference on the modelling non-zero income distribution and the probability of zero rounding  
 392 were impacted by this model misspecification.

393

## 5. Application to HILDA dataset

394 We applied the SPM model to a dataset of  $n = 1,693$  individuals from waves 1 to  
 395 10 (corresponding to 2001 to 2010), from the HILDA survey. The data were relatively  
 396 unbalanced, with only 481 individuals having records for all 10 waves, while the fewest  
 397 number of measurements available for an individual was  $n_i = 2$  (306 individuals had this  
 398 cluster size). The average number of measurements for each individual was  $n_i = 7.83$ , with  
 399 the number of observations equal to  $\sum_{i=1}^{1,693} n_i = 10,650$ . As the measure of income and thus  
 400 the response, we used the answer to the question ‘Last financial year, what was your total  
 401 wage and salary income from all jobs before tax or anything else was deducted?’, where the  
 402 financial year runs from July one calendar year to June the following calendar year. We scaled  
 403 the responses to be in terms of thousands of Australian dollars. From these observed income

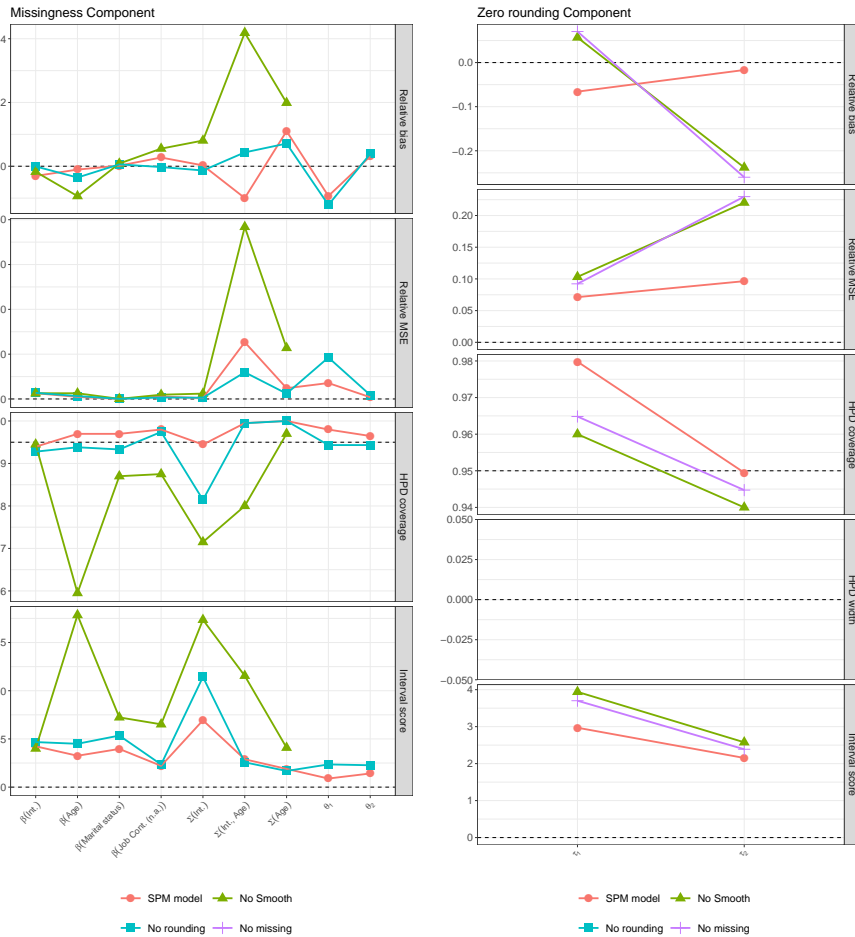
Figure 2. Simulation results for the components of the SPM model relating to the probability of an individual earning an income and the non-zero income distribution (conditional on earning an income). Each row represents one performance measure (from top to bottom: relative bias, relative MSE, coverage probability of 95% HPD intervals, mean interval score). Note some of the four competing models by definition do not contain estimates for certain parameters.



404 values, there were 809 (7.60%) missing values, 3,459 (32.48%) reported zero incomes, and  
 405 6,382 (59.92%) reported non-zero income values.

406 There was evidence of rounding at several non-zero values of income (Figure 1, right  
 407 panel), but as discussed at the end of Section 2 the largest repeated non-zero income value  
 408 only occupied 1.60% of the total number of observations. Of the observed non-zero incomes,  
 409 the mean and median income (per financial year in thousands of AUD) across all individuals  
 410 and time points was \$41.468 and \$36.000 respectively, and values ranged from \$0.035 to  
 411 \$823.241. A histogram of the non-zero income values displayed a strong right skew, while a

Figure 3. Simulation results for the components SPM model relating to the probability of missingness and the probability of zero rounding. Each row represents one performance measure (from top to bottom: relative bias, relative MSE, coverage probability of 95% HPD intervals, mean interval score). Note some of the four competing models by definition do not contain estimates for certain parameters.



412 log transformation produced a more symmetric distribution close to normality (see Figure 1),  
 413 although there was some evidence that it was a little too strong with some outliers on the  
 414 lower end of income values.

415 We considered a set of five covariates for inclusion in the SPM model: gender (a factor  
 416 with levels female/male), presence of a long term health condition (a factor with levels  
 417 no/yes), marital status (a factor with levels single or de facto/divorced/married/separated), job  
 418 contract (a factor with levels casual/fixed-term/permanent/not asked), and age (a continuous  
 419 variable ranging from 17 to 93 years, with the median of 41 years and middle 50% of  
 420 ages occurring between 31 and 54). To clarify, individuals who were categorised as ‘not



asked' for job contract type were those who had not undertaken any paid employment in the seven days prior to them answering the survey. But they may still have been employed at some point within the corresponding financial year. There were no missing values for any of the covariates. We provide statistical and graphical summaries of the five covariates in Supplementary Material B.1.

We conducted an exploratory data analysis to decide which covariates to include in the different components of the SPM model. The details of this are provided in Supplementary Material B.2, and we summarise the results final here. As fixed effect covariates in the vectors  $\mathbf{x}_{ij,\text{earn}}$ ,  $\mathbf{x}_{ij,\text{nonzero}}$ , and  $\mathbf{x}_{ij,\text{miss}}$ , we included an intercept, all four factor covariates i.e., gender (using female as the reference level), presence of long term health condition (using no as the reference level), marital status (using single or de facto as the reference level), and job contract (using casual as the reference level) as fixed effects using dummy variables, and a smooth term for age defined using thin-plate regression splines. As random effect covariates in the vectors  $\mathbf{z}_{ij,\text{earn}}$ ,  $\mathbf{z}_{ij,\text{nonzero}}$ , and  $\mathbf{z}_{ij,\text{miss}}$ , we included a random intercept and slope for age. Note in our exploratory data analysis, the same set of candidate fixed and random effect covariates ended up being included in the components for modelling the probability of earning an income, the distribution of non-zero income, and the probability of non-response.

Finally, one interesting finding from our exploratory analysis was that, ignoring the missing income values, close to half the number of individuals were observed to always have non-zero income (767 out of 1,693; 45.82%), while there were also a sizeable number of individuals who were observed to always have zero income (251 out of 1,693; 14.99%). As a result, we anticipate that the estimated random intercept for the probability of an individual earning an income would be relatively large to accommodate such 'perfect ones' and 'perfect zeros' for the clusters. Analogously, by far the majority of individuals had no missing incomes (1,320 out of 1,693; 78.00%), while there were a small number for which their income was always missing (19 out of 1,693; 1.12%). In turn, we also anticipate that the estimated random intercept for the probability of non-response would also be somewhat large.

## 5.1. Model fitting and results

We began by fitting the SPM model with three choices for the distribution of non-zero incomes: gamma, exponential and lognormal, as discussed below (2), while fixing  $s_{\text{miss}}(\cdot)$  to have a quadratic form; we consider sensitivity analysis of this choice in Supplementary Material B.4. Both the deviance information criterion and penalised expected deviance (Plummer 2008), calculated based only on the observed non-zero incomes, suggested that the lognormal distribution was the most suitable choice for  $\mathcal{F}_+(\cdot; \eta, \phi)$ , and so we opted to use this distribution for the remainder of the analysis; see the results in Supplementary

Table 1. Posterior median estimates and 95% highest posterior density (HPD) intervals for the parameters in the SPM model fitted to the HILDA survey (excluding the fixed effects of age, which was included as a smooth term), assuming a lognormal distribution for non-zero incomes and a quadratic form for  $s_{\text{miss}}(\cdot)$ . Fixed effect coefficients whose highest posterior density intervals do not contain zero are marked with an asterisk.

Parameter	Median (HPD interval)	Parameter	Median (HPD interval)
$\Psi_{\text{earn}}$		$\Psi_{\text{nonzero}}$	
<b>Fixed effects (<math>\beta_{\text{earn}}</math>)</b>		<b>Fixed effects (<math>\beta_{\text{nonzero}}</math>)</b>	
Intercept	5.174 (4.230, 6.229)*	Intercept	2.670 (2.545, 2.774)*
Gender (male)	2.091 (1.643, 2.567)*	Gender (male)	0.736 (0.648, 0.818)*
Health condition (yes)	-0.536 (-0.853, -0.241)*	Health condition (yes)	-0.045 (-0.098, 0.006)
Marital status (divorced)	0.434 (-0.572, 1.442)	Marital status (divorced)	0.207 (-0.051, 0.462)
Marital status (married)	0.135 (-0.336, 0.556)	Marital status (married)	-0.011 (-0.073, 0.046)
Marital status (separated)	-0.364 (-1.115, 0.371)	Marital status (separated)	0.050 (-0.067, 0.172)
Job contract (fixed-term)	1.488 (-0.467, 3.834)	Job contract (fixed-term)	0.400 (0.322, 0.476)*
Job contract (permanent)	1.758 (0.545, 2.983)*	Job contract (permanent)	0.402 (0.348, 0.456)*
Job contract (not asked)	-6.862 (-7.941, -5.968)*	Job contract (not asked)	0.022 (-0.044, 0.082)
<b>Random effects (<math>\Sigma_{\text{earn}}</math>)</b>		<b>Random effects (<math>\Sigma_{\text{nonzero}}</math>)</b>	
Intercept	5.380 (4.201, 6.770)	Intercept	0.624 (0.552, 0.702)
Age	0.522 (0.145, 1.301)	Age	0.456 (0.362, 0.566)
Covariance(Int., Age)	0.897 (0.207, 1.618)	Covariance(Int., Age)	0.330 (0.264, 0.406)
		<b>Dispersion parameter</b>	
		$\phi$	0.288 (0.276, 0.299)
$\Psi_{\text{miss}}$		$\Psi_{\text{round}}(s_{\text{round}}(\cdot))$	
<b>Fixed effects (<math>\beta_{\text{miss}}</math>)</b>		$\tau_0$	-1.696 (-2.062, -1.400)*
Intercept	-4.342 (-4.992, -3.683)*	$\tau_1$	-0.105 (-0.128, -0.074)*
Gender (male)	-0.121 (-0.514, 0.247)		
Health condition (yes)	0.025 (-0.305, 0.352)		
Marital status (divorced)	0.191 (-1.059, 1.309)		
Marital status (married)	0.425 (0.062, 0.827)*		
Marital status (separated)	0.031 (-0.650, 0.728)		
Job contract (fixed-term)	-0.407 (-0.992, 0.178)		
Job contract (permanent)	-0.223 (-0.607, 0.151)		
job contract (not asked)	-0.759 (-1.178, -0.328)*		
<b>Random effects (<math>\Sigma_{\text{miss}}</math>)</b>			
Intercept	2.480 (1.746, 3.337)		
Age	0.356 (0.083, 0.913)		
Covariance(Int., Age)	-0.199 (-0.734, 0.329)		
<b>Smoothing function (<math>s_{\text{miss}}(\cdot)</math>)</b>			
$\theta_1$	-0.009871 (-0.026275, 0.006357)		
$\theta_2$	0.000234 (0.000115, 0.000324)*		

456 Material B.3. However we note that all three candidate distributions produced similar overall  
 457 conclusions for the different components of the SPM model. A table of the posterior median  
 458 estimates and corresponding 95% highest posterior density intervals for the SPM model  
 459 assuming lognormally distributed non-zero incomes and a quadratic form for  $s_{\text{miss}}(\cdot)$  is  
 460 presented in Table 1, while results with the two alternative choices for  $\mathcal{F}_+(\eta, \phi)$  are provided  
 461 in Supplementary Material B.3.

462 To visualise different parts of the SPM model, we constructed a set of four graphs. The  
 463 first two were plots of marginal mean income versus age, where marginal income is defined

464 as

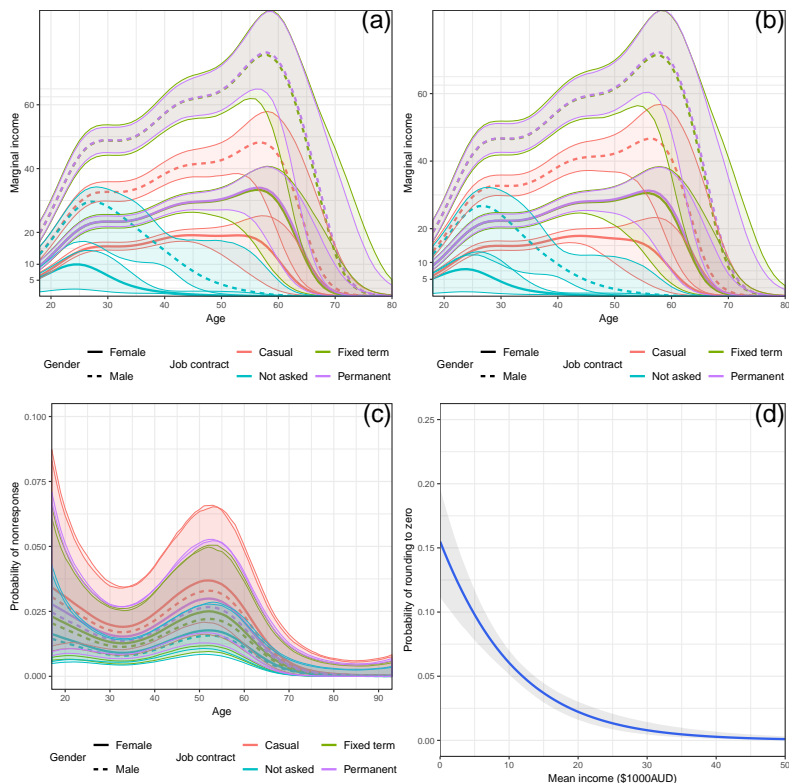
$$\begin{aligned}
 E_{Y_{\text{true}}}(y_{ij,\text{true}} | \mathbf{b}_{i,\text{earn}} = \mathbf{0}, \mathbf{b}_{i,\text{nonzero}} = \mathbf{0}) &= E_G \{ E_{Y_{\text{true}}|G}(y_{ij,\text{true}} | g_{ij}, \mathbf{b}_{i,\text{earn}} = \mathbf{0}, \mathbf{b}_{i,\text{nonzero}} = \mathbf{0}) \} \\
 &= \frac{\exp(\mathbf{x}_{ij,\text{earn}}^\top \boldsymbol{\beta}_{\text{earn}})}{1 + \exp(\mathbf{x}_{ij,\text{earn}}^\top \boldsymbol{\beta}_{\text{earn}})} \exp\left(\mathbf{x}_{ij,\text{nonzero}}^\top \boldsymbol{\beta}_{\text{nonzero}} + \frac{1}{2}\phi\right).
 \end{aligned}$$

465 Note this is equal to the mean of the mixture model that results from marginalising the non-  
 466 zero income regression model in (2) with respect to the model for probability of earning  
 467 an income in (1), and then setting the random effects to zero. Figures 4a,b depict this for  
 468 individuals with absence and presence of long term health condition, respectively. We also  
 469 plot the effect of age on the average probability of missing income,  $E(r_{ij} | g_{ij}, \mathbf{b}_{i,\text{nonzero}} =$   
 470  $\mathbf{0}, \mathbf{b}_{i,\text{miss}} = \mathbf{0})$ , in Figure 4c. Finally, a plot of how the probability of rounding to zero  $\nu_{ij}$   
 471 depends on the underlying mean income,  $\mu_{ij}$ , is provided in Figure 4d.

472 On average, male individuals who had a permanent job contract (in the past seven days  
 473 prior to taking the survey) and no long term health condition were substantially more likely  
 474 to earn an income (Table 1). Individuals who had not undertaken any paid employment in  
 475 the past seven days prior to taking the survey were substantially less likely to earn income,  
 476 keeping in mind that income was based on the entire financial year. Conditioning on earning  
 477 an income, male individuals on a fixed-term or permanent job contract were also on average  
 478 the highest income earners. There was no clear evidence that marital status had any effect  
 479 on the propensity to earn income or on the amount of income earned. Finally, both the fixed  
 480 and random intercept in the model for the probability of earning income were quite large in  
 481 magnitude: This is in line with what we anticipated from the exploratory data analysis, which  
 482 found close to half of the individuals observed to have non-zero income, while another 15%  
 483 of individuals were observed to always report zero income.

484 From Figures 4a,b, we see that on average marginal income tended to increase in a non-  
 485 linear manner from age 17, peaking around age 57 before dropping dramatically afterwards to  
 486 zero. This finding is consistent with previous studies which show income tends to maximise  
 487 for individuals between the ages 50 to 60. Moreover, the dramatic drop in income after age  
 488 60 reflects the probability of earning income itself decreasing rapidly around this age range  
 489 i.e., retirement age. We caution against over interpreting the model particularly beyond the  
 490 age of 70, given the data here are comparably sparse; see Supplementary Material B.1. The  
 491 exception to this trend was individuals who had no job contract (in the past seven days prior  
 492 undertaking the survey), where marginal income peaks around age 25 before declining to  
 493 zero. The marginal income plots also clearly showed that for a given age, males tend to have  
 494 higher incomes than females, while individuals on fixed-term or permanent job contracts have  
 495 higher incomes than those on casual or no job contracts. Finally, marginal income tended to

Figure 4. Plots illustrating various components of the SPM model fitted to the HILDA survey: (a) marginal income as a function of age, gender, and job contract, for individuals with no long term health condition; (b) marginal income as a function of age, gender, and job contract, for individuals with a long term health condition; (c) probability of missing response as a function of age, gender, and job contract; and (d) probability of zero rounding as a function of mean income  $\mu_{ij}$ . In all plots, 95% pointwise highest posterior density intervals are indicated by shading of the corresponding color.



496 be slightly lower for individuals with a long term health condition, although this effect was  
 497 not particularly evident when comparing Figures 4a,b: this was consistent with the fixed effect  
 498 for health condition being significant (but not large compared to the effects of gender, say)  
 499 for  $\Psi_{\text{earn}}$  but not for  $\Psi_{\text{nonzero}}$ .

500 Turning to the model for the probability of missing response, and under the missing data  
 501 mechanism of the SPM model, we observe a large negative estimated fixed intercept while the  
 502 random intercept also suggested considerable heterogeneity between individuals (Table 1).  
 503 This was in line with the exploratory analyses which found over 75% of individuals had no  
 504 missing incomes at all. Being legally married and not undertaking any paid employment  
 505 in the seven days prior to undertaking the survey were the only two fixed effect terms  
 506 found to have substantial direct effects on the probability of non-response. The estimates

507 of the quadratic smoothing function also provide some evidence that low and high income  
508 earners are more likely to not report their income. We explore the sensitivity of the results  
509 to the specific missingness assumptions in Supplementary Material B.4. It should be noted  
510 that the magnitudes of both estimated coefficients parameterising the quadratic smoothing  
511 function were very small. This suggests that the actual effect of the mean response  $\mu_{ij}$  on the  
512 probability of non-response was dominated by the direct effects of covariates. Plotting the  
513 overall probability of non-response against age, we see a non-linear effect where probability  
514 of missing incomes starts out relatively high at a young age before decreasing and then  
515 increasing again, reaching a peak of around 0.05 around ages 53, before decreasing towards  
516 zero (Figure 4). This interesting non-linear effect was somewhat consistent with what we saw  
517 in our exploratory analysis; see Supplementary Material B.1. More importantly, it should be  
518 kept in mind that this non-linear effect is only producing very small changes in the actual  
519 probability of missing response e.g., from 0.01 to 0.05. Finally, there was strong evidence  
520 of low income earners tended to round and reporting their income as zero (Table 1), with  
521 the probability of zero rounding rapidly decreasing towards zero and becoming negligible  
522 beyond \$20,000 AUD (Figure 4d). Note however the actual probability of zero rounding was  
523 at most around 15%.

524 In Supplementary Material B.4, we present results from a sensitivity analysis used to  
525 assess the impact of different choices of  $s_{\text{miss}}(\cdot)$ , and consequently of different assumptions  
526 about the missing data process, on parameter estimates in the SPM model. We adopted an  
527 imputation based approach (motivated by [Carpenter, Kenward & White 2007](#), among others)  
528 to conduct this assessment. Results show that many of the parameter estimates are indeed  
529 sensitive to the assumption of incomes being non-ignorable, with linear and quadratic choices  
530 for  $s_{\text{miss}}(\cdot)$  tending to produce more similar results to each other compared to assuming  
531  $s_{\text{miss}}(\cdot) = 0$ . We emphasise that our goal for performing the sensitivity analysis is not to  
532 provide evidence in favour of one missingness mechanism over another (since this choice is  
533 not explicitly possible to begin with, [Molenberghs et al. 2008](#)). Rather, it is simply to show  
534 that the HILDA survey is sensitivity to the choice of the missing data assumption.

535

## 6. Discussion

536 We proposed and applied a shared parameter mixture model to analyse national  
537 Australian income distributions and evolution over time. The SPM model simultaneously  
538 handles the complications of non-ignorable missingness and zero rounding by first modelling  
539 the latent, true income process through a two-component mixture model combining a  
540 point mass at zero describing the probability of an individual earning an income with a  
541 positive continuous distribution describing how non-zero income varies. A shared parameter

542 component is then used to model the probability of non-response, and where the probability  
543 of missingness can vary with the latent non-zero income in a nonlinear manner. Finally, a  
544 zero rounding component is included such that the probability of an individual reporting  
545 zero income, when they actually earn an income, depends solely on the proximity of the  
546 latent income to zero. Applying the SPM model to the HILDA survey produces several key  
547 results related to longitudinal income: 1) gender, the presence of a long term health condition,  
548 job contract type, and age were important factors in determining the likelihood of earning an  
549 income; 2) gender, age, and job contract type have important effects in driving the distribution  
550 of non-zero incomes over time; 3) job contract type and marital status had direct effects  
551 on the probability of non-response, while the parameter estimate in the SPM model were  
552 sensitive to the missingness mechanism assumption; 4) there is a strong association between  
553 the probability of zero rounding and a person's true non-zero income.

554 One of the underlying motivations behind proposing the SPM model was to deal with  
555 potential income-dependent non-responses in income. In many datasets on income however,  
556 it is quite common to provide imputed values based on assuming the responses are missing  
557 at random for practical reasons (as is the case in the HILDA survey e.g., [Watson & Starick](#)  
558 [2011](#)), even though there is reasonable sociological evidence of missing income responses  
559 being non-ignorable (see the references in Section 1 and also [Bollinger et al. 2014](#)). While  
560 not conclusive, our results from analysing the HILDA survey do show that inference is  
561 sensitive to the assumption surrounding the missing data mechanism on income, and so  
562 future analyses involving such data from the survey should more seriously acknowledge  
563 and/or address this issue. Also, while the use of Bayesian MCMC estimation was strongly  
564 motivated by the hierarchical nature of the SPM model, given the increasingly large volume  
565 of many longitudinal datasets then future research should explore more efficient estimation  
566 approaches (e.g., [Kim & Wand 2018](#); [Niku et al. 2019](#)), as well other more efficient Bayesian  
567 estimation approaches such as those implemented in Stan ([Carpenter et al. 2017](#)) and nimble  
568 ([de Valpine et al. 2017](#)).

569

### References

- 570 BACCI, S. & BARTOLUCCI, F. (2015). A multidimensional finite mixture structural equation model for  
571 nonignorable missing responses to test items. *Structural Equation Modeling: A Multidisciplinary*  
572 *Journal* **22**, 352–365.
- 573 BACCI, S., BARTOLUCCI, F., BETTIN, G. & PIGINI, C. (2019). A latent class growth model for migrants  
574 remittances: an application to the German Socio-Economic Panel. *Journal of the Royal Statistical*  
575 *Society: Series A (Statistics in Society)* **182**, 1607–1632.
- 576 BECHTEL, L., LORDAN, G. & RAO, D. (2012). Income inequality and mental health – empirical evidence  
577 from Australia. *Health economics* **21**, 4–17.

- 578 BOHNING, D. & ALFO, M. (2016). Editorial: Special issue on models for continuous data with a spike at  
579 zero. *Biometrical Journal* **58**, 255–258.
- 580 BOLLINGER, C.R., HIRSCH, B.T., HOKAYEM, C.M. & ZILIAK, J.P. (2014). Trouble in the Tails: Earnings  
581 Non-Response and Response Bias across the Distribution. In *Joint Statistical Meetings, Boston, August*.
- 582 CALDERN-OJEDA, E., AZPITARTE, F. & GMEZ-DNIZ, E. (2016). Modelling income data using two  
583 extensions of the exponential distribution. *Physica A: Statistical Mechanics and its Applications* **461**,  
584 756–766.
- 585 CARPENTER, B., GELMAN, A., HOFFMAN, M.D., LEE, D., GOODRICH, B., BETANCOURT, M.,  
586 BRUBAKER, M., GUO, J., LI, P. & RIDDELL, A. (2017). Stan: A probabilistic programming language.  
587 *Journal of statistical software* **76**.
- 588 CARPENTER, J.R., KENWARD, M.G. & WHITE, I.R. (2007). Sensitivity analysis after multiple imputation  
589 under missing at random: a weighting approach. *Statistical Methods in Medical Research* **16**, 259–275.
- 590 CREEMERS, A., HENS, N., AERTS, M., MOLENBERGHS, G., VERBEKE, G. & KENWARD, M.G. (2011).  
591 Generalized shared-parameter models and missingness at random. *Statistical modelling* **11**, 279–310.
- 592 DE VALPINE, P., TUREK, D., PACIOREK, C.J., ANDERSON-BERGMAN, C., LANG, D.T. & BODIK, R.  
593 (2017). Programming with models: writing statistical algorithms for general model structures with  
594 NIMBLE. *Journal of Computational and Graphical Statistics* **26**, 403–413.
- 595 DENWOOD, M.J. (2016). *runjags: An R Package Providing Interface Utilities, Model Templates, Parallel*  
596 *Computing Methods and Additional Distributions for MCMC Models in JAGS*. *Journal of Statistical*  
597 *Software* **71**, 1–25.
- 598 FITZMAURICE, G., DAVIDIAN, M., VERBEKE, G. & MOLENBERGHS, G. (2008). *Longitudinal data*  
599 *analysis*. Boca Raton, Florida: Chapman and Hall/CRC.
- 600 GELMAN, A., JAKULIN, A., PITTAU, M.G. & SU, Y. (2008). A weakly informative default prior distribution  
601 for logistic and other regression models. *The Annals of Applied Statistics* **2**, 1360–1383.
- 602 GELMAN, A. & RUBIN, D.B. (1992). Inference from iterative simulation using multiple sequences.  
603 *Statistical Science* , 457–472.
- 604 GIUSTI, C. & LITTLE, R.J. (2011). An analysis of nonignorable nonresponse to income in a survey with a  
605 rotating panel design. *Journal of Official Statistics* **27**, 211–229.
- 606 GNEITING, T. & RAFTERY, A.E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal*  
607 *of the American Statistical Association* **102**, 359–378.
- 608 GROSS, M. & RENDTEL, U. (2016). Kernel density estimation for heaped data. *Journal of Survey Statistics*  
609 *and Methodology* **4**, 339–361.
- 610 HANISCH, J.U. (2005). Rounded responses to income questions. *Allgemeines Statistisches Archiv* **89**, 39–  
611 48.
- 612 HEADEY, B. & WOODEN, M. (2004). The effects of wealth and income on subjective well-being and ill-  
613 being. *Economic record* **80**, 24–33.
- 614 HUI, F.K.C., MLLER, S. & WELSH, A.H. (2017). Hierarchical selection of fixed and random effects in  
615 generalized linear mixed models. *Statistica Sinica* **27**, 501–518.
- 616 JENKINS, S.P. (2010). The British household panel survey and its income data. *ISER Working Paper Series*  
617 **33**.
- 618 KIM, A.S. & WAND, M.P. (2018). On expectation propagation for generalised, linear and mixed models.  
619 *Australian & New Zealand Journal of Statistics* **60**, 75–102.
- 620 LITTLE, R.J. & RUBIN, D.B. (2014). *Statistical analysis with missing data*. Hoboken, New Jersey: John  
621 Wiley & Sons.
- 622 LIU, L., SHIH, Y.C.T., STRAWDERMAN, R.L., ZHANG, D., JOHNSON, B.A. & CHAI, H. (2019).  
623 Statistical analysis of zero-inflated nonnegative continuous data: A review. *Statistical Science* **34**, 253–  
624 279.
- 625 MOLENBERGHS, G., BEUNCKENS, C., SOTTO, C. & KENWARD, M.G. (2008). Every missingness not at  
626 random model has a missingness at random counterpart with equal fit. *Journal of the Royal Statistical*

- 627 *Society: Series B (Statistical Methodology)* **70**, 371–388.
- 628 NIKU, J., BROOKS, W., HERLIANSYAH, R., HUI, F.K.C., TASKINEN, S. & WARTON, D.I. (2019).  
629 Efficient estimation of generalized linear latent variable models. *PloS one* **14**, e0216129.
- 630 PLUMMER, M. (2008). Penalized loss functions for Bayesian model comparison. *Biostatistics* **9**, 523–539.
- 631 PLUMMER, M. et al. (2003). JAGS: A program for analysis of Bayesian graphical models using Gibbs  
632 sampling. In *Proceedings of the 3rd International Workshop on Distributed Statistical Computing (DSC*  
633 *2003)*. March. p. 20–22.
- 634 RIPHAHN, R.T. & SERFLING, O. (2005). Item non-response on income and wealth questions. *Empirical*  
635 *Economics* **30**, 521–538.
- 636 RIZOPOULOS, D., VERBEKE, G. & MOLENBERGHS, G. (2008). Shared parameter models under random  
637 effects misspecification. *Biometrika* **95**, 63–74.
- 638 RUBIN, D.B. (1976). Inference and missing data. *Biometrika* **63**, 581–592.
- 639 SCHRAEPLER, J.P. (2006). Explaining income nonresponse – A case study by means of the British Household  
640 Panel Study (BHPS). *Quality & Quantity* **40**, 1013–1036.
- 641 TSONAKA, R., VERBEKE, G. & LESAFFRE, E. (2009). A semi-parametric shared parameter model to handle  
642 nonmonotone nonignorable missingness. *Biometrics* **65**, 81–87.
- 643 WATSON, N. & STARICK, R. (2011). Evaluation of Alternative Income Imputation Methods for a  
644 Longitudinal Survey. *Journal of Official Statistics* **27**, 693–715.
- 645 WATSON, N. & WOODEN, M.P. (2012). The HILDA survey: a case study in the design and development of  
646 a successful household panel survey. *Longitudinal and Life Course Studies* **3**, 369–381.
- 647 WULFSOHN, M.S. & TSIATIS, A.A. (1997). A joint model for survival and longitudinal data measured with  
648 error. *Biometrics* **53**, 330–339.
- 649 ZINN, S. & WRBACH, A. (2016). A statistical approach to address the problem of heaping in self-reported  
650 income data. *Journal of Applied Statistics* **43**, 682–703.