

# Bayesian Approaches for Missing Not at Random Outcome Data: The Role of Identifying Restrictions

Antonio R. Linero\* and Michael J. Daniels†

August 31, 2017

## Abstract

Missing data is almost always present in real datasets, and introduces several statistical issues. One fundamental issue is that, in the absence of strong uncheckable assumptions, effects of interest are typically not nonparametrically identified. In this article, we review the generic approach of the use of identifying restrictions from a likelihood-based perspective, and provide points of contact for several recently proposed methods. An emphasis of this review is on restrictions for nonmonotone missingness, a subject that has been treated sparingly in the literature. We also present a general, fully-Bayesian, approach which is widely applicable and capable of handling a variety of identifying restrictions in a uniform manner.

**Key words:** missing data; MNAR; mixture models; multiple imputation; nonignorable missingness; nonparametric Bayes.

## 1 Introduction

Missing data is highly prevalent in real datasets. Within a likelihood-based framework, missing data can best be categorized as either ignorable or nonignorable (Rubin, 1976); the former does not require a model for the missingness process, while the latter does. Nonignorable missingness introduces fundamental identifiability issues because, by virtue of the fact that we did not observe the missing data, we have no data with which to estimate its distribution.

The literature is filled with approaches which resolve identifiability issues by making parametric modeling assumptions (see Section 2 for a review). Following Cox and Donnelly (2011, page 96), however, we believe that if an issue cannot be resolved nonparametrically given an infinite sample then it is “usually dangerous to resolve it parametrically.” While parametric approaches are useful, we argue that they should not indirectly resolve identifiability issues. An alternative approach is to incorporate non-identifiability into the analysis. The full-data distribution can be factored into two components: (1) the observed-data distribution, which is identified by the observed data; and (2) the conditional distribution of the missing data given the observed data, sometimes called the *extrapolation distribution*, which is not identified (Daniels and Hogan 2008, Section 8.2; Little 1995). Different assumptions about the

---

\*Department of Statistics, Florida State University, arlinero@stat.fsu.edu

†Department of Statistics, University of Florida, mdaniels@stat.ufl.edu

33 missing data can be expressed in terms of *identifying restrictions* which allow the analyst to  
34 recover the full-data distribution from the observed data distribution. The most well-known  
35 identifying restriction is the missing at random (MAR) assumption (Rubin, 1976), but many  
36 alternatives exist.

37 The National Research Council (2010) recommends the routine use of *sensitivity analysis*  
38 to assess the impact of assumptions about the missing data on inference. Two approaches to  
39 sensitivity analysis are to first consider many different identifying restrictions (Thijs et al.,  
40 2002) and second (in the spirit of Rotnitzky et al. 1998 and Daniels and Hogan 2008, Chapter  
41 9) to introduce an unidentified *sensitivity parameter*  $\xi$  which represents an interpretable  
42 deviation from a benchmark identifying restriction. The sensitivity parameter  $\xi$  should be  
43 such that (1) there is no information in the data to inform  $\xi$  and (2) upon specification of  $\xi$ ,  
44 the effects of interest are identified.

45 Concerns about parametric assumptions have motivated frequentist semiparametric ap-  
46 proaches (Robins et al., 1995; Scharfstein et al., 1999) which make minimal assumptions about  
47 the full-data distribution. These approaches posit a parametric model for the missing data  
48 mechanism and a semiparametric model for the outcome distribution, and produce estimates  
49 by solving inverse-probability-weighted (IPW) estimating equations. These procedures are  
50 frequently doubly-robust, requiring the analyst to specify one of the two models correctly to  
51 attain consistent estimation (Scharfstein et al., 1999; Rotnitzky et al., 1998; Tsiatis, 2007).  
52 Recently, there have been various likelihood-based approaches proposed which have the flex-  
53 ibility of semiparametric approaches and allow a flexible sensitivity analysis (Wang et al.,  
54 2010; Linero and Daniels, 2015; Linero, 2017). An advantage of the Bayesian approach is  
55 that it allows for uncertainty about the unidentified components of the model to be encoded  
56 in an informative prior, allowing the analyst to incorporate subject-matter expertise formally  
57 into the analysis.

58 This article has three goals. First, we provide a review of model-based approaches to  
59 nonignorable missingness, including parametric approaches which identify the full-data dis-  
60 tribution (see National Research Council, 2010; Ibrahim and Molenberghs, 2009, for addi-  
61 tional reviews of MNAR modeling strategies). Our second goal is to summarize and review  
62 existing identifying restrictions in the literature. A special emphasis is given to recent pro-  
63 posals for nonmonotone missingness, as this subject has received a sparser treatment in the  
64 literature. We highlight several recently proposed identifying restrictions and characterize  
65 them as generalizations of monotone restrictions.

66 Our third goal is to propose a flexible, fully-Bayesian, framework for incomplete outcome  
67 data. First, a flexible Bayesian nonparametric model is chosen for the observed data distri-  
68 bution. Second, we use an identifying restriction to identify the extrapolation distribution.  
69 The framework allows for many different restrictions to be used without needing to change  
70 the model used for the observed data, can accommodate both monotone and nonmonotone  
71 missingness, and allows for the introduction of sensitivity parameters. The proposed ap-  
72 proach might be perceived as a competitor to the IPW approaches which are prevalent in

73 the literature. However, it has several features which IPW approaches do not. First, the  
 74 Bayesian framework allows for expert knowledge to be formally incorporated into the anal-  
 75 ysis by eliciting informative priors on sensitivity parameters. Second, the approach allows  
 76 for simultaneous inference about functionals of the full-data distribution, rather than just a  
 77 specifically chosen functional such as the mean; for example it is possible to make inferences  
 78 about means and quantiles simultaneously. Third, we are not required to fit different models  
 79 depending on the choice of identifying restriction, allowing for a more principled comparison  
 80 of different restrictions.

81 To illustrate the necessity of conducting a principled sensitivity analysis, we analyze data  
 82 from the Breast Cancer Prevention Trial (BCPT). A concern in this study was that the  
 83 treatment tamoxifen might cause depression. We show that the evidence for this hypothesis  
 84 is strongly influenced by the assumptions made about the missingness, and that seemingly  
 85 similar assumptions can yield dramatically different results. This underscores the need for  
 86 statisticians and subject-matter experts to work together in determining which assumptions  
 87 about the missing data are most appropriate for a particular problem.

## 88 1.1 Notation

89 Let  $Y_j^{(i)}$  denote the measurement of variable  $j$  intended to be collected on subject  $i$  for  
 90  $i = 1, \dots, N$ , and let  $Y^{(i)} = (Y_1^{(i)}, \dots, Y_J^{(i)})$ . Let  $R^{(i)} = (R_1^{(i)}, \dots, R_J^{(i)})$  be a vector of  
 91 missingness indicators such that  $R_j^{(i)} = 1$  or  $0$  according to whether  $Y_j^{(i)}$  is observed or not.  
 92 For a given binary vector  $r \in \{0, 1\}^J$ , let  $y_r = (y_j : r_j = 1)$  and  $y_{-r} = (y_j : r_j = 0)$ . The  
 93 observed data on subject  $i$  is then given by  $Y_{R^{(i)}}^{(i)}$ , and the missing data is given by  $Y_{-R^{(i)}}^{(i)}$ .

94 We assume the pairs  $(Y^{(i)}, R^{(i)})$  are iid with density  $p(y, r)$  with respect to some measure;  
 95 implicitly,  $p(y, r)$  may depend on a parameter vector  $\theta$ . We refer to  $p(y, r)$  as the *full-data*  
 96 *distribution*. To lighten notation, we will often work with an iid copy  $(Y, R)$  of  $(Y^{(1)}, R^{(1)})$ .  
 97 For simplicity we omit covariates; in principle all distributions we discuss can be defined  
 98 conditional on fully-observed covariates  $X = x$ .

99 We will abuse notation, for example writing  $p(y)$  for the marginal density of  $Y$  or  $p(r | y)$   
 100 for the probability of  $R = r$  given  $Y = y$ ; it will always be clear from context what density is  
 101 being referred. When specific arguments are required, we will write for example  $p(R_j = 1 |$   
 102  $Y = y)$  for the probability of  $R_j = 1$  given  $Y = y$ .

103 For a fixed  $r$ , let  $p(y, r) = p(y_r, r)p(y_{-r} | y_r, r)$  denote the *extrapolation factorization*  
 104 (Daniels and Hogan, 2008, Section 8.2) of  $p(y, r)$ . This factors  $p(y, r)$  into the product of a  
 105 term which is identified and a term which is unidentified. Note that  $p(y_r, r)$  is the density  
 106 of the observed data  $(Y_R, R)$  while  $p(y_{-r} | y_r, r)$  is the conditional density of the missing  
 107 data  $Y_{-R}$ . We refer to  $p(y_r, r)$  as the *observed-data distribution* and to  $p(y_{-r} | y_r, r)$  as the  
 108 *extrapolation distribution*.

109 Missingness is said to be *monotone* if  $R_j = 0$  implies  $R_{j+1} = 0$ . This commonly occurs in  
 110 longitudinal trials when missingness is due to dropout. Missingness can then be summarized  
 111 by the last time at which a subject is measured  $S^{(i)} = \max\{j : R_j^{(i)} = 1\}$ , which we refer

112 to as the (index of the) *dropout time*. For longitudinal studies it is also useful to let  $\bar{Y}_j^{(i)} =$   
 113  $(Y_1^{(i)}, \dots, Y_j^{(i)})$  denote the history of the response up-to time  $j$ , and let  $\tilde{Y}_j^{(i)} = (Y_{j+1}^{(i)}, \dots, Y_J^{(i)})$   
 114 denote the future of the response strictly after time  $j$ . Thus,  $Y^{(i)} = (\bar{Y}_j^{(i)}, \tilde{Y}_j^{(i)})$ . We similarly  
 115 define  $\bar{R}_j^{(i)}$  and  $\tilde{R}_j^{(i)}$ .

## 116 1.2 Running example: the Breast Cancer Prevention Trial

117 To make the concepts presented concrete, we will focus on applications to the Breast Cancer  
 118 Prevention Trial (BCPT), a clinical trial which assigned women at high-risk of developing  
 119 breast cancer to either a preventative drug, tamoxifen, or to a placebo. One aim of this study  
 120 was to determine if tamoxifen causes depression. The response  $Y_j^{(i)}$  is 1 or 0 according to  
 121 whether subject  $i$  is depressed or not at time  $j$ . Roughly  $N = 5000$  subjects were assigned to  
 122 each of tamoxifen ( $Z = 1$ ) and control ( $Z = 0$ ). Measurements were scheduled to be taken  
 123 at baseline and 3, 6, 12, 18, 24, 30, and 36 months from baseline, for  $J = 8$  intended measure-  
 124 ments. There was a substantial amount of missingness at all time points, and missingness was  
 125 highly nonmonotone. A concern is that depression at time  $j$  might be associated with missing-  
 126 ness at time  $j$ , even after conditioning on other observables, resulting in MNAR missingness.  
 127 Our primary interest is in the intention-to-treat effect  $\psi = E(Y_J | Z = 1) - E(Y_J | Z = 0)$ .

128 To help illustrate concepts, we will also consider a simplified setting in which  $J = 2$ . We  
 129 refer to this setting as the reduced Breast Cancer Prevention Trial (RBCPT). We assume  
 130 that  $(Y_1, Y_2)$  represent *continuous*, rather than binary, measures of depression level (the  
 131 actual binary responses were created from dichotomizing a quantitative score) to create more  
 132 generality in the development.

## 133 2 Basic MNAR modeling strategies

134 We divide strategies for modeling  $p(y, r)$  into three categories: (1) selection models; (2)  
 135 pattern mixture models; and (3) shared parameter models. In Section 2.4, we describe how  
 136 any of these three approaches can be used to obtain a model for the observed data, without  
 137 modeling the missing data.

### 138 2.1 Selection models

139 The selection modeling approach (Heckman, 1979) is based on the factorization  $p(y, r) =$   
 140  $p(y) \cdot p(r | y)$ . The term  $p(r | y)$  is referred to as the *missing data mechanism*.

141 **Example 1.** Consider the RBCPT. With monotone missingness and  $Y_1$  always observed,  
 142 following Diggle and Kenward (1994), we set

$$\begin{aligned}
 Y &\sim \text{Normal}(\mu, \Sigma), \\
 p(R_2 = 1 | y_1, y_2, R_1 = 1) &= \text{expit}(\phi_0 + \phi_1 y_1 + \phi_2 y_2).
 \end{aligned}
 \tag{1}$$

143 Selection models are attractive for their conceptual simplicity. In the context of the  
 144 BCPT, the selection factorization suggests a causal mechanism in which depression causes  
 145 missingness to occur. As  $p(y)$  is directly available, inference is usually straight-forward.

146 One drawback of parametric selection models is that they may “identify away” the missing  
 147 data problem. Observe that  $\phi_2 = 0$  corresponds to an MAR missing data mechanism in (1).  
 148 One may be tempted to test for MNAR missingness by testing  $\phi_2 = 0$ . As we have stressed,  
 149 testing for MAR cannot be done without recourse to parametric assumptions. As illustrated  
 150 by Kenward (1998), inferences about MAR in this setup are extremely sensitive to parametric  
 151 assumptions. When  $p(y)$  is a Gaussian density,  $(\phi_1, \phi_2)$  function as skewness parameters for  
 152  $p(y_2 | y_1, r)$  and can be estimated from the observed data. Hence, there are no sensitivity  
 153 parameters which can be used as a basis of a sensitivity analysis. In practice, the likelihood of  
 154  $\phi_2$  may be flat enough that it can be used as an approximate sensitivity parameter (Carpenter  
 155 et al., 2002). This problem is mitigated to some extent when semiparametric or nonparametric  
 156 models for  $Y$  are used, although this becomes more difficult as the dimension of the response  
 157 increases. Note also that  $p(y_{-r} | y_r, r)$  is not available in closed form; consequently, it is  
 158 difficult to describe on a conceptual level how missing values are imputed relative to the  
 159 other approaches we describe.

## 160 2.2 Pattern mixture models

161 The pattern mixture approach (Little, 1994, 1993; Hogan and Laird, 1997) is based on the  
 162 factorization  $p(y, r) = p(y | r)p(r)$ . This characterizes  $p(y)$  as a mixture over missingness  
 163 patterns  $\sum_r p(y | r)p(r)$ . The pattern mixture factorization is closely related to the extrap-  
 164 olation factorization, with  $p(y_r, r) = p(y_r | r) \cdot p(r)$ . This makes the pattern mixture approach  
 165 conducive to sensitivity analysis.

166 **Example 2.** Consider the RBCPT and assume monotone missingness with  $Y_1$  always ob-  
 167 served. We set  $\phi = p(R_2 = 1)$ ,  $(Y_1 | R_2 = r) \sim \text{Normal}(\mu^{(r)}, \sigma_1^{(r)})$ , and  $(Y_2 | Y_1 = y_1, R_2 =$   
 168  $r) \sim \text{Normal}(\alpha^{(r)} + \beta^{(r)}y_1, \sigma_2^{(r)})$ . The parameters  $(\alpha^{(0)}, \beta^{(0)}, \sigma_2^{(0)})$  are unidentified. One ap-  
 169 proach to identifying these parameters is to link them to the  $R_2 = 1$  pattern, setting for  
 170 example  $(\beta^{(0)}, \sigma_2^{(0)}) = (\beta^{(1)}, \sigma_2^{(1)})$  and  $\alpha^{(0)} = \alpha^{(1)} + \xi$ . This implies that the influence of  $Y_1$   
 171 on  $Y_2$  and the conditional spread of  $Y_2$  do not depend on  $R_2$ , while the conditional mean of  
 172  $Y_2$  does and is shifted by a fixed amount  $\xi$ . The parameter  $\xi$  is a sensitivity parameter, and  
 173 can be varied as part of a sensitivity analysis.

174 Characteristic of pattern mixture models, the above model allows an interpretable sensi-  
 175 tivity analysis and is transparent in how the it imputes missing values on a conceptual level.  
 176 There are several shortcomings of the pattern mixture approach. Conceptually, it is typically  
 177 not easy to interpret how the response  $Y$  influences the probability of missingness at time  $j$ .  
 178 In the BCPT, a pattern mixture model suggests that those with missing values come from  
 179 a distinct sub-population; an arguably more natural way to capture this intuition is through  
 180 the use of latent class models (Roy, 2003) (though as constructed there, they do not allow

181 sensitivity parameters). Pattern mixture models often possess a large number of unidentified  
 182 parameters that the analyst must specify, with the situation becoming unwieldy in higher  
 183 dimensions. Additionally, sparsity in the observed missing data patterns  $R^{(i)}$  may necessitate  
 184 further modeling of  $p(y | r)$  to share information across times.

### 185 2.3 Shared parameter approaches

186 The shared parameter approach captures dependence between  $Y^{(i)}$  and  $R^{(i)}$  through shared  
 187 random effects (Wu and Carroll, 1988; Henderson et al., 2000), setting  $p(y, r) = \int p(y |$   
 188  $b) p(r | b) G(db)$ . The random effect distribution  $G(\cdot)$  can be specified parametrically, usually  
 189 as a multivariate Gaussian distribution, or nonparametrically.

190 **Example 3.** Consider the BCPT. We set  $(b_1, b_2) \sim \text{Normal}(\mu_b, \Sigma_b)$  and assume that, condi-  
 191 tional on  $b$ , all components of  $(Y, R)$  are mutually independent with logit  $p(Y_j = 1 | b) = Z_j^\top b_1$   
 192 and logit  $p(R_j = 1 | b) = W_j^\top b_2$ . For example, to get a random quadratic trend over time, we  
 193 might set  $Z_j^\top = W_j^\top = (1, t_j, t_j^2)$  where  $t_j$  is the time of measurement  $j$ . This type of shared  
 194 parameter model is referred to as a *correlated random effects model* (Lin et al., 2010).

195 The shared parameter approach provides a highly flexible framework for analyzing non-  
 196 ignorable missingness, and is particularly effective for modeling complex data structures  
 197 (Dunson and Perreault, 2001). Shared parameter models appeal strongly to intuition, sug-  
 198 gesting that  $Y$  and  $R$  have a shared, unobserved, common cause. A drawback of the shared  
 199 parameter approach is that it is difficult to separate  $p(y_r, r)$  from  $p(y_{-r} | y_r, r)$ , making it  
 200 difficult to anchor a sensitivity analysis to an interpretable identifying restriction (see Sec-  
 201 tion 3). Generally, it is not easy to see what assumptions about the missing data mechanism  
 202 are encoded in a shared parameter model.

203 Methods for implementing a sensitivity analysis for shared parameter models have been  
 204 developed by Creemers et al. (2010, 2011). In our example, one might set  $\text{logit}(p(Y_j = 1 |$   
 205  $b, R = r)) = Z_j^\top (b_1^{(i)} + r_j \delta)$  which gives an adjustment to the random effect  $b_1$  at the times  
 206 for which  $r_j = 0$ . One may then set, for example,  $\delta \sim \text{Normal}(\mu_\delta, \Sigma_\delta)$ , with  $\xi = (\mu_\delta, \Sigma_\delta)$  a  
 207 sensitivity parameter. We feel that this is somewhat against the spirit of the shared parameter  
 208 model, as  $Y$  and  $R$  are no longer conditionally independent and the causally suggestive  
 209 motivation is stretched.

### 210 2.4 Observed data modeling

211 The models in Sections 2.1–2.3 have been presented as models for the joint density  $p(y, r)$ .  
 212 An alternative strategy is to model the observed data distribution  $p(y_r, r)$  and leave the  
 213 extrapolation distribution  $p(y_{-r} | y_r, r)$  unspecified. One can then fit a model  
 214 for  $p(y_r, r)$  to the data and complete the model using one of the identifying restrictions  
 215 described in Section 3.

216 Directly modeling  $p(y_r, r)$  can be challenging to do in practice, as it requires a model for  
 217  $Y_r$  for every pattern  $r$ . When missingness is monotone, one approach is to specify models for

218  $p(y_j | S \geq j, \bar{y}_{j-1})$  and  $p(S = j | S \geq j, \bar{y}_j)$ . For examples of this approach, see [Scharfstein](#)  
 219 [et al. \(2014\)](#) and [Wang et al. \(2010\)](#). Other approaches to directly modeling  $p(y_r, r)$  often  
 220 use the pattern mixture approach, specifying models for  $p(y_r | r)$  while leaving  $p(y_{-r} | y_r, r)$   
 221 unspecified. See, for example, [Little \(1994\)](#) and [Thijs et al. \(2002\)](#).

222 A generic approach to modeling the observed data is to specify a *working model* ([Linero](#),  
 223 [2017](#); [Linero and Daniels, 2015](#); [Daniels and Linero, 2015](#)). One then implicitly obtains  
 224 a model for the observed data  $p(y_r, r) = \int p^*(y, r) dy_{-r}$ . In principle,  $p^*(y, r)$  may be a  
 225 selection model, pattern mixture model, or shared parameter model. In [Section 5](#) we will  
 226 apply this approach using a nonparametrically modeled shared parameter to obtain a highly  
 227 flexible model of the observed data.

228 A benefit of the working model approach is that it allows models which share information  
 229 across missingness patterns and time, without identifying the extrapolation distribution. This  
 230 allows one to avoid a common pitfall of pattern-mixture models; we can estimate  $p(y_r, r)$   
 231 even when we do not observe some patterns or the amount of data in some patterns is sparse.  
 232 Because the model  $p^*(y, r)$  is used only to obtain a model for  $p(y_r, r)$ , and is not used as  
 233 a basis for inference, we are allowed complete freedom in how to identify the extrapolation  
 234 distribution. Conveniently,  $p^*(y, r)$  can also be used as a basis for Markov chain Monte Carlo  
 235 algorithms.

236 In practice, the working model framework has the drawback of being somewhat difficult  
 237 to implement, in that one must be able to derive the conditional distributions  $p^*(y_r | R = r')$ .  
 238 This places restrictions on which models can be tractably used; in particular, selection models  
 239 and parametric shared parameter models are difficult to use. Fortunately, there are very  
 240 flexible models that are tractable. An additional concern is that, when  $p(y_r, r)$  is modeled  
 241 parametrically,  $p(y, r)$  will usually fall outside of this parametric family. For example, when  
 242 using identifying restrictions, if  $p(y_r | r)$  is modeled with a Gaussian distribution, it will not  
 243 typically be the case that  $p(y | r)$  is Gaussian ([Wang and Daniels, 2011](#)). Consequently,  
 244 the joint model  $p(y, r)$  may not be easily interpretable, although causal effects may still be  
 245 computed using MC integration (see [Section 4](#)).

### 246 3 Identifying restrictions

247 Identifying restrictions provide a useful starting point for identifying the extrapolation dis-  
 248 tribution and conducting a sensitivity analysis. Informally, an identifying restriction is an  
 249 assumption about  $p(y, r)$  which links the observed data distribution  $p(y_r, r)$  to the extrapo-  
 250 lation distribution  $p(y_{-r} | y_r, r)$ .

251 We remark that identifying assumptions differ subtly throughout the literature; for  
 252 example, [Seaman et al. \(2013\)](#) give several non-equivalent definitions of MAR. All restrictions  
 253 we consider will be phrased in the form of conditional independencies, with (for example)  
 254 MAR corresponding to the conditional independence statement  $(Y_{-r} | Y_r, R = r) \stackrel{d}{=} (Y_{-r} | Y_r)$   
 255 for all patterns  $r$ .

256 The goal of specifying an identifying restriction is to nonparametrically identify the pa-  
257 rameters of interest.

258 **Definition 3.1.** Let  $\mathcal{Q}$  denote the set of observed data distributions  $q(y_r, r)$ , and let  $\mathcal{P}$   
259 be some family of full-data distributions  $p(y, r)$ . The family  $\mathcal{P}$  is said to *nonparametrically*  
260 *identify* a parameter  $\psi(p)$  if,

- 261 1. For every  $q \in \mathcal{Q}$ , there exists a  $p \in \mathcal{P}$  such that  $q$  is the associated observed data  
262 density of  $p$ .
- 263 2. For every  $q \in \mathcal{Q}$ , if  $p, p' \in \mathcal{P}$  both marginalize to  $q$ , then  $\psi(p) = \psi(p')$ .

264 The family  $\mathcal{P}$  is said to be *nonparametrically saturated* (Robins, 1997; Vansteelandt et al.,  
265 2006) if, for each  $q \in \mathcal{Q}$ , there exists a unique  $p \in \mathcal{P}$  which marginalizes to  $q$ .

266 In the absence of strong subject-matter knowledge, it is unwise to assume that a par-  
267 ticular identifying restriction holds. Nevertheless, in practice it can be useful to specify a  
268 single identifying restriction as a benchmark assumption, and consider interpretable devi-  
269 ations from that benchmark. For example, one might “anchor” an analysis to MAR and  
270 consider smooth deviations from MAR. Considering several anchors, and deviations from  
271 these anchors, provides insight into how inferences are driven by our assumptions.

272 We differentiate three different types of identifying restrictions. *Joint* restrictions com-  
273 pletely identify  $p(y, r)$ ; that is, they lead to nonparametrically saturated models. *Marginal* re-  
274 strictions do not identify  $p(y, r)$ , but identify the marginals  $p(y_j)$ ; an example is the sequential  
275 explainability assumption (Vansteelandt et al., 2007) discussed later. Marginal restrictions  
276 do not lead to nonparametrically saturated models, but are sufficient to nonparametrically  
277 identify all marginal effects. Marginal restrictions can be useful because (i) they may be more  
278 readily interpretable than joint restrictions, and (ii) they may encode weaker assumptions.  
279 Marginal restrictions are special cases of *partial* restrictions, which are any restrictions which  
280 do not identify  $p(y, r)$ .

### 281 3.1 Identifying restrictions under monotone missingness

282 The missing data problem becomes much simpler when missingness is monotone. In this  
283 case, the missing data pattern can be summarized by the dropout time  $S = \max\{j : R_j = 1\}$ .  
284 Monotonicity occurs naturally when missingness is due to dropout in a longitudinal study.  
285 Techniques for monotone missingness can also be applied if there is a method of ordering the  
286 components of  $Y$  which makes missingness monotone.

287 **Example 4** (NCMV). Consider the BCPT, and assume that missingness is monotone. We  
288 conjecture that the, if a subject drops out at time  $k < j$ , then their missing response at time  
289  $j$  can reasonably be approximated using an equivalent individual who instead drops out at  
290 time  $j$ ; so, we set  $(Y_j \mid \bar{Y}_{j-1}, S = k) \stackrel{d}{=} (Y_j \mid \bar{Y}_{j-1}, S = j)$ . Thijs et al. (2002) refer to this as  
291 the *neighboring case missing value* (NCMV) restriction.



292 **Example 5** (ACMV). Consider again the BCPT with monotone missingness. We conjecture  
 293 that, if a subject drops out at time  $k < j$ , then their response at time  $j$  can reasonably  
 294 be approximated by using an equivalent subject who dropped out *after* time  $j$ ; so, we set  
 295  $(Y_j | \bar{Y}_{j-1}, S = k) \stackrel{d}{=} (Y_j | \bar{Y}_{j-1}, S \geq j)$ . [Little \(1993\)](#) refers to this as the *available case*  
 296 *missing value* (ACMV) restriction.

297 **Example 6** (CCMV). In the BCPT, we decide to use the observations of those who complete  
 298 the study to estimate the conditional distribution of the missing observations; so, we set  
 299  $(Y_j | \bar{Y}_{j-1}, S = k) \stackrel{d}{=} (Y_j | \bar{Y}_{j-1}, S = J)$ ; [Little \(1993\)](#) refers to that as the *complete case*  
 300 *missing value* (CCMV) restriction.

301 The goal of using these restrictions is to provide a starting point for a sensitivity analysis.  
 302 In practice, when missingness is MNAR, none of the conditional independencies asserted  
 303 above is realistic; in fact, ACMV is itself equivalent to MAR ([Molenberghs et al., 1998](#))! In  
 304 the BCPT, if the depression status of an individual at time  $j$  is a strong predictor of  $R_j = 0$   
 305 then one may expect the conditional distribution of  $Y_j$  to be stochastically larger than what  
 306 is implied by ACMV, NCMV, or CCMV.

307 Under monotone missingness, ACMV is equivalent to MAR. This suggests that missing-  
 308 ness at time  $j + 1$  is causally linked only to the past values of  $\bar{Y}_j$ . The NFD restriction  
 309 ([Kenward et al., 2003](#)) generalizes this idea.

310 **Example 7** (NFD). We posit that missingness at time  $j + 1$  is causally due to the past  
 311 and present values of  $Y$ , so that  $p(S = j + 1 | Y) = p(S = j + 1 | \bar{Y}_{j+1})$ , or equivalently  $(Y_{j+1} |$   
 312  $S = k, \bar{Y}_j) \stackrel{d}{=} (Y_{j+1} | S \geq j, \bar{Y}_j)$ . This is referred to as the non-future dependence (NFD)  
 313 assumption.

314 Despite its causal motivation, we note that NFD is not a causal law; for example if  $(Y, R)$   
 315 share an unobserved common cause, NFD will usually be violated. Given that MAR implies  
 316 NFD, but not vice-versa, NFD leads to an under-identified model (and thus is a partial  
 317 restriction); in particular, the distribution  $(Y_j | S = j - 1, \bar{Y}_{j-1})$  is unidentified for  $j > 2$ .  
 318 This is convenient, as it allows the analyst to consider *families* of restrictions, all of which  
 319 satisfy the NFD restriction. For example, [Linero and Daniels \(2015\)](#) centers a sensitivity  
 320 analysis on the MAR assumption by setting  $(Y_j | \bar{Y}_{j-1}, S = j - 1) \stackrel{d}{=} (Y_j + \xi | \bar{Y}_{j-1}, S \geq j)$ ,  
 321 with  $\xi = 0$  corresponding to MAR.

322 The ACMV, NCMV, and CCMV restrictions are all joint restrictions. [Birmingham et al.](#)  
 323 [\(2003\)](#) consider several partial restrictions, including the following marginal restriction which  
 324 is implied by CCMV.

325 **Example 8** (Last-occasion-pattern-mixture). We posit that the conditional distribution of  
 326  $Y_J$  at the end of study, given  $\bar{Y}_j$  and  $S = j$ , can reasonably be approximated by the distribution  
 327 of those who complete the study; hence, we set  $(Y_J | \bar{Y}_j, S = j) \stackrel{d}{=} (Y_J | \bar{Y}_j, S = J)$ .

328 A general tool for extending the restrictions above to the nonmonotone settings is to  
 329 assume that missingness is partially ignorable given  $S$  ([Harel and Schafer, 2009](#)). This sets

330  $p(R = r \mid Y = y, S = s) = p(R = r \mid Y_r = y_r, S = s)$ , and assumes the parameters  
 331 of  $p(r \mid y, s)$  are independent of the parameters of  $p(y, s)$ . Analogously to ignorability,  
 332 partial ignorability ensures that likelihood-based inferences for  $p(y, s)$  do not depend on how  
 333  $p(r \mid y, s)$  is modeled. See Wang et al. (2010) for an application of this assumption to the  
 334 BCPT data.

### 335 3.2 Identifying restrictions for nonmonotone missingness

336 The topic of identifying restrictions under nonmonotone missingness was initiated by Robins  
 337 (1997), who proposed the class of permutation missingness (PM) models. Let  $\bar{O}_j$  denote the  
 338 observed data (including the  $R_j$ 's) up-to-and-including time  $j$ , and  $\tilde{O}_j$  the data observed  
 339 strictly after time  $j$ . The PM restriction assumes

$$(R_j \mid Y, \tilde{R}_j) \stackrel{d}{=} (R_j \mid \bar{Y}_{j-1}, \tilde{O}_j) \quad (2)$$

340 possibly after applying an a-priori known permutation to  $Y$ . In words, (2) states that miss-  
 341 ingness at time  $j$  can depend on the “past” and the “observed future,” but not on the  
 342 present, where the notion of time is determined by the given permutation. For longitu-  
 343 dinal data, one can use (2) without a permutation, or use the reverse permutation to get  
 344  $(R_j \mid Y, \bar{R}_{j-1}) \stackrel{d}{=} (R_j \mid \bar{O}_{j-1}, \tilde{Y}_j)$  which states that missingness depends on the “future” and  
 345 the “observed past.”

346 Our opinion is that PM models are difficult to explain to practitioners. We review several  
 347 alternative assumptions which have been introduced relatively recently.

348 **Example 9** (Sequential explainability). For the BCPT, we believe that the observed depres-  
 349 sion levels prior to time  $j$  are sufficient to predict whether or not a subject will be measured at  
 350 time  $j$ , while the outcome at time  $j$  is not predictive. We therefore impose the sequential ex-  
 351 plainability restriction (Vansteelandt et al., 2007)  $(Y_j \mid \bar{O}_{j-1}, R_j = 0) \stackrel{d}{=} (Y_j \mid \bar{O}_{j-1}, R_j = 1)$ .

352 **Example 10** (NIP). For the BCPT, we believe that, all other observed quantities being  
 353 equal, missingness at time  $j$  is not predictive of depression at time  $j$ . We therefore posit  
 354 the nearest identified pattern (NIP) (Lineró, 2017) restriction,  $(Y_j \mid R = r, Y_r) \stackrel{d}{=} (Y_j \mid R =$   
 355  $r_j^*, Y_r)$ , where  $r_j^*$  is equal to  $r$ , but with  $j$ th component fixed at 1.

356 Both NIP and sequential explainability are marginal restrictions. NIP appears similar to  
 357 NCMV. A more direct analog is the itemwise conditional independence (ICIN) assumption,  
 358 introduced independently by Sadinle and Reiter (2017a) and Shpitser (2016).

359 **Example 11** (ICIN). For the BCPT, we believe that all other quantities (both observed  
 360 and unobserved) being equal, missingness at time  $j$  is not predictive of depression at time  
 361  $j$ . We therefore posit the ICIN restriction  $(Y_j \mid R_j = 0, R_{-j}, Y_{-j}) \stackrel{d}{=} (Y_j \mid R_j = 1, R_{-j}, Y_{-j})$   
 362 where  $R_{-j} = (R_k : k \neq j)$  and  $Y_{-j} = (Y_k : k \neq j)$  denote  $R$  and  $Y$  with the  $j$ th component  
 363 removed.

364 ICIN and NIP differ in that (i) NIP conditions only on the observed components of  $Y$   
 365 and (ii) ICIN is a joint restriction. To the extent that conditioning on additional variables  
 366 makes conditional independence more tenable, ICIN is very attractive. To our knowledge,  
 367 practical algorithms for conducting inference under ICIN are lacking when  $J$  is moderately  
 368 large. Results of [Sadinle and Reiter \(2017a\)](#) imply that ICIN is equivalent to NCMV when  
 369 missingness is monotone. A proof of the following proposition is deferred to the supplementary  
 370 material.

371 **Proposition 3.2.** *ICIN is an extension of NCMV to nonmonotone missingness.*

372 [Tchetgen Tchetgen et al. \(2016\)](#) introduced the pairwise missing at random assumption.  
 373 The name is motivated by the observation that it corresponds to MAR when, for fixed  $r$ , we  
 374 assume  $R \in \{r, \mathbf{1}\}$ , where  $\mathbf{1} = (1, \dots, 1)$ .

375 **Example 12** (PMAR). For the BCPT, we believe that the distribution of the missing values  
 376 of a subject can reasonably be approximated using an equivalent subject who was observed  
 377 at all measurement times. We therefore posit the pairwise missing at random (PMAR)  
 378 restriction,  $(Y_{-r} | R = r, Y_r) \stackrel{d}{=} (Y_{-r} | R = \mathbf{1}, Y_r)$ .

379 Just as ICIN is a joint restriction which generalizes NCMV, PMAR is a joint restriction  
 380 which generalizes CCMV; the following proposition is immediate from the definition.

381 **Proposition 3.3.** *PMAR is an extension of CCMV to nonmonotone missingness.*

### 382 3.3 Sensitivity parameters for identifying restrictions

383 The identifying restrictions in Section 3.1 and Section 3.2 are phrased in terms of conditional  
 384 independence relationships which, as we have noted, are not themselves particularly plausible  
 385 when  $Y_j$  is thought to directly influence  $R_j$ . We consider these assumptions not because  
 386 we believe the conditional independencies they suggest, but rather to use as benchmark  
 387 assumptions. These assumptions can be embedded in a family of restrictions indexed by a  
 388 *sensitivity parameter*  $\xi \in \Xi$  such that (1) there is no information in the data to identify  $\xi$  and  
 389 (2) upon specifying  $\xi$ , the effects of interest are identified. It is essential that the sensitivity  
 390 parameter  $\xi$  be interpretable; our convention will be to associate the benchmark assumption  
 391 with  $\xi = 0$ . The index  $\xi$  can then be thought of as a smooth deviation from our benchmark  
 392 assumption.

393 **Example 13.** For the BCPT, we believe the NIP restriction is unreasonable because depres-  
 394 sion at time  $J$  should increase the risk of missingness, even after accounting for the observed  
 395 data. We instead assume

$$p(y_J | y_r, R = r) = \frac{p(y_J | y_r, R = r_j^*)e^{\gamma y_J}}{E(e^{\gamma Y_J} | Y_r = y_r, R = r_j^*)}.$$

396 Let  $A = \{r, r_j^*\}$ ; using Bayes theorem it can be shown that

$$\log \frac{\text{Odds}(R_J = 0 \mid Y_r, Y_J = 1, R \in A)}{\text{Odds}(R_J = 0, \mid Y_r, Y_J = 0', R \in A)} = \gamma,$$

397 so that  $\gamma$  denotes the effect, on the log-odds scale, of  $Y_J = 1$  on missingness.

398 The exponential tilting strategy is very widely applicable, and we now outline it in a  
 399 general form. Examples of works using this strategy include [Birmingham et al. \(2003\)](#); [Wang](#)  
 400 [et al. \(2010\)](#); [Scharfstein et al. \(2014, 1999\)](#); [Tchetgen Tchetgen et al. \(2016\)](#); [Vansteelandt](#)  
 401 [et al. \(2007\)](#). Consider a restriction of the form

$$(U \mid V = v, W = w) \stackrel{d}{=} (U \mid V = v', W = w), \quad (3)$$

402 where  $U$  is a subset of the missing data,  $W$  is a subset of the complete data distinct from  
 403  $U$ , and  $V$  is a subset of the missing data indicators. The values  $v$  and  $v'$  are such that  $U$   
 404 is missing when  $V = v$ , while  $U$  is observed when  $V = v'$ . For example, under sequential  
 405 explainability, one has  $\{U = Y_j, W = \bar{O}_{j-1}, V = R_j, v = 0, v' = 1\}$  while under PMAR one  
 406 has  $\{U = Y_{-r}, W = Y_r, V = R, v = r, v' = \mathbf{1}\}$ . Let  $f_v(u \mid w)$  and  $f_{v'}(u \mid w)$  denote the  
 407 densities of the distributions in (3). The exponential tilting approach sets

$$f_v(u \mid w) = \frac{f_{v'}(u \mid w) \exp\{t(u, w)\}}{E[\exp\{t(U, w)\} \mid V = v', W = w]}. \quad (4)$$

408 The function  $t(u, w)$  is a function-valued sensitivity parameter. By Bayes theorem,

$$\log \frac{\text{Odds}(V = v \mid U = u, W = w, V \in \{v, v'\})}{\text{Odds}(V = v \mid U = u', W = w, V \in \{v, v'\})} = t(u, w) - t(u', w).$$

409 Hence,  $t(\cdot, w)$  determines the effect of a change in  $U$  on the log-odds of  $V = v$  versus  $V = v'$ .

410 Another option is to consider a transformation-based approach similar to [Daniels and](#)  
 411 [Hogan \(2000\)](#). This is particularly useful when the underlying response is continuous.

412 **Example 14.** Consider the BCPT, but with  $Y$  instead representing a continuous measure of  
 413 depression level. We believe the NIP restriction is unreasonable because we expect depression  
 414 levels to be higher among those missing at time  $j$ , even after conditioning on the observed  
 415 data. We instead assume  $(Y_j \mid Y_r = R = r) \stackrel{d}{=} (Y_j + \xi_j \mid Y_r, R = r_j^*)$ , where  $\xi_j > 0$  represents  
 416 the expected increase in depression level when a subject is missing rather than observed.

417 More generally, starting from (3), one can specify a generic transformation

$$(U \mid V = v, W = w) \stackrel{d}{=} (\mathcal{T}(U, w) \mid V = v', W = w). \quad (5)$$

418 In practice we must specify  $\mathcal{T}(u, w)$  to be interpretable by subject-matter experts. Location  
 419 or location-scale transformations, such as  $\mathcal{T}_j(Y_j) = \xi_{0j} + \xi_{1j}Y_j$ , are popular ([Daniels and](#)

420 Hogan, 2000; Wang and Daniels, 2011; Gaskins et al., 2016) and can be computationally  
 421 advantageous. Non-affine choices for  $\mathcal{T}(\cdot)$  can be used to rescale the data before applying an  
 422 affine transformation.

423 A meaningful sensitivity analysis requires serious engagement with subject-matter ex-  
 424 perts, and as such requires for  $\Xi$  to be low dimensional. A common approach that does  
 425 not formally account for the effect of uncertainty in  $\xi$  is a “tipping point” approach. This  
 426 identifies values, or regions of values, of  $\xi$  which result in substantively different conclusions  
 427 for the effects of interest. If plausible values of  $\xi$  do not include any tipping points, then  
 428 we can have confidence in our substantive conclusions; we note, however, that tipping point  
 429 analyses do not incorporate uncertainty in  $\xi$  when quantifying uncertainty in treatment ef-  
 430 fects. For illustrations of tipping point analyses, see Scharfstein et al. (2014) and Liublinska  
 431 and Rubin (2014). An option which formally incorporates uncertainty in  $\xi$  is to place an  
 432 informative prior on  $\xi$ . As there is no information in the data about  $\xi$ , this prior for  $\xi$  will  
 433 also be the posterior. An advantage of this approach is that it combines all restrictions under  
 434 consideration to achieve a single, final, inference. For examples of this approach, see Daniels  
 435 and Hogan (2008, Chapter 9, Case Study 2), Wang et al. (2010), and Gaskins et al. (2016).

### 436 3.4 A pattern mixture modeling example

437 We now show how one might combine the identifying restrictions described above with a  
 438 model for the observed data for the RBCPT (using the original depression score). We specify  
 439 a pattern mixture model

$$p(R_1 = i, R_2 = j) = \phi_{ij}, \quad [Y_1, Y_2 \mid R = (1, 1)] \sim \text{Normal}(\mu^{(1,1)}, \Sigma^{(1,1)}),$$

$$[Y_1 \mid R = (1, 0)] \sim \text{Normal}(\mu_1^{(1,0)}, \sigma_1^{(1,0)}), \quad [Y_2 \mid R = (0, 1)] \sim \text{Normal}(\mu_2^{(0,1)}, \sigma_2^{(0,1)}).$$

440 All parameters above can be estimated from the observed data using standard techniques;  
 441 for example, we have  $\hat{\mu}^{(1,1)} = \frac{1}{N^{(1,1)}} \sum_{i: R_1^{(i)} = R_2^{(i)} = 1} (Y_1^{(i)}, Y_2^{(i)})^\top$ . For convenience, we write

$$(Y_1 \mid Y_2, R_1 = 1, R_2 = 1) \sim \text{Normal}(\alpha + \beta Y_2, \tau^2),$$

442 where  $(\alpha, \beta, \tau)$  is a function of  $(\mu^{(1,1)}, \Sigma^{(1,1)})$ . Suppose that interest is in the parameter  
 443  $\zeta = E(Y_1)$ . We demonstrate how  $\zeta$  is identified under the PMAR, sequential explainability,  
 444 and NIP assumptions. First, by iterated expectation,

$$\zeta = \sum_{i=0}^1 \sum_{j=0}^1 \phi_{ij} E(Y_1 \mid R_1 = i, R_2 = j).$$

445 Observe that  $E(Y_1 \mid R_1 = 1, R_2 = 1) = \mu_1^{(1,1)}$  and  $E(Y_1 \mid R_1 = 1, R_2 = 0) = \mu_1^{(1,0)}$ . This  
 446 leaves  $E(Y_1 \mid R_1 = 0, R_2 = 0)$  and  $E(Y_1 \mid R_1 = 0, R_2 = 1)$  to be identified.

447 Consider first the PMAR assumption. This implies  $E(Y_1 \mid R_1 = 0, R_2 = 0) = E(Y_1 \mid$

448  $R_1 = 1, R_2 = 1) = \mu_1^{(1,1)}$ . Using iterated expectation, PMAR also implies

$$\begin{aligned} E(Y_1 \mid R_1 = 0, R_2 = 1) &= E\{E(Y_1 \mid Y_2, R_1 = 0, R_2 = 1) \mid R_1 = 0, R_2 = 1\} \\ &= E\{E(Y_1 \mid Y_2, R_1 = 1, R_2 = 1) \mid R_1 = 0, R_2 = 1\} \\ &= E(\alpha + \beta Y_2 \mid R_1 = 0, R_2 = 1) = \alpha + \beta \mu_2^{(0,1)}. \end{aligned}$$

449 This gives

$$\zeta_{\text{PMAR}} = \phi_{00}\mu_1^{(1,1)} + \phi_{10}\mu_1^{(1,0)} + \phi_{01}(\alpha + \beta\mu_2^{(0,1)}) + \phi_{11}\mu_1^{(1,1)}.$$

450 Next, we consider NIP. The derivations under NIP are exactly the same as those under  
 451 PMAR, with the exception that  $E(Y_1 \mid R_1 = 0, R_2 = 0) = E(Y_1 \mid R_1 = 1, R_2 = 0) = \mu_1^{(1,0)}$ .  
 452 Therefore, under NIP we have

$$\zeta_{\text{NIP}} = \zeta_{\text{PMAR}} + \phi_{00}(\mu_1^{(1,0)} - \mu_1^{(1,1)}).$$

453 Hence,  $\zeta_{\text{NIP}}$  will be larger than  $\zeta_{\text{PMAR}}$  when  $\mu_1^{(1,0)} > \mu_1^{(1,1)}$ , and vice versa. Lastly, we  
 454 consider sequential explainability. At time  $j = 1$  there is no observed history, so sequential  
 455 explainability implies the marginal independence  $Y_1 \perp R_1$ . Consequently,

$$\zeta_{\text{SE}} = E(Y_1 \mid R_1 = 1) = \frac{\phi_{10}}{\phi_{10} + \phi_{11}}\mu_1^{(1,0)} + \frac{\phi_{11}}{\phi_{10} + \phi_{11}}\mu_1^{(1,1)}.$$

456 Sequential explainability differs fundamentally from NIP and PMAR as, due to its sequential  
 457 nature, it does not use the distribution of  $(Y_2, R_2)$  to identify  $\zeta$ .

458 We now incorporate sensitivity parameters under sequential explainability. Note that if  
 459  $(Y_1 \mid R_1 = 0) \stackrel{d}{=} (Y_1 + \xi \mid R_1 = 1)$ , then  $\xi = 0$  is consistent with sequential explainability.  
 460 Under this assumption we have

$$\begin{aligned} \zeta(\xi) &= p(R_1 = 1)E(Y_1 \mid R_1 = 1) + p(R_1 = 0)E(Y_1 \mid R_1 = 0) \\ &= p(R_1 = 1)\zeta_{\text{SE}} + p(R_1 = 0)(\zeta_{\text{SE}} + \xi) \\ &= \zeta_{\text{SE}} + (\phi_{00} + \phi_{01})\xi. \end{aligned}$$

461 Sensitivity analysis may now proceed either by eliciting an informative prior on  $\xi$ , or by  
 462 identifying values of  $\xi$  which lead to substantively different inferences.

## 463 4 Inference and computation

464 We discuss two approaches to computation. First, we describe a fully-Bayesian approach,  
 465 which can be computationally demanding. Second, we describe multiple imputation, which  
 466 is a computationally simpler approximation. Let  $\theta$  denote the parameters of the model of  
 467  $p(y, r)$ ,  $\pi(\theta)$  a prior for  $\theta$ , and  $\mathcal{O} = (Y_{R^{(1)}}^{(1)}, R^{(1)}, \dots, Y_{R^{(N)}}^{(N)}, R^{(N)})$  the observed data. We first

---

**Algorithm 1** Monte Carlo integration for sequential explainability

---

```
1: procedure GCOMP( $\theta, T, j$ )  $\triangleright$  Approximates  $\mu_j$  by simulating  $T$  samples from  $p_\theta(y)$ 
2:   for  $t = 1, \dots, T$  do
3:     Sample  $(Y_{R^{(t)}}, R^{(s)}) \sim p_\theta(y_r, r)$ .
4:     if  $R_j^{(t)} = 0$  then
5:       Sample  $Y_j^{(t)} \sim p_\theta(y_j \mid \bar{o}_{j-1}, R_j^{(t)} = 1)$ 
6:     end if
7:   end for
8:   Set  $\mu_j = T^{-1} \sum_{s=1}^T Y_j^{(t)}$ .
9:   return  $\mu_j$ 
10: end procedure
```

---

468 obtain samples of  $\theta$  from its posterior distribution  $\pi(\theta \mid \mathcal{O}) \propto \prod_{i=1}^N p(Y_{R^{(i)}}, R^{(i)}) \pi(\theta)$ , usually  
469 by Markov chain Monte Carlo. When the working model framework described in Section 2.4  
470 is used, samples of  $\theta$  can be obtained by fitting the working model by data augmentation,  
471 taking advantage of the fact that  $p_\theta(y_r, r) = \int p_\theta^*(y, r) dy_{-r}$ . A benefit of this approach is  
472 that sampling  $\theta$  can often be accomplished using general-purpose software for fitting Bayesian  
473 models. Software packages, such as JAGS and WinBUGS, allow for fast fitting of custom models,  
474 and accommodate missing values.

475 Fully-Bayesian inference then proceeds by computing effects of interest directly from the  
476 sampled  $\theta$ 's and the chosen identifying restriction. Multiple imputation, by contrast, uses the  
477 sampled  $\theta$ 's to impute completed datasets some number of times using the identifying restric-  
478 tion. Practically, these approaches are operationally quite similar. We begin by describing  
479 fully-Bayesian inference, and describe the changes required to perform multiple imputation.

#### 480 4.1 Fully-Bayesian inference

481 Given sampled values  $\theta \sim \pi(\theta \mid \mathcal{O})$ , fully-Bayesian inference requires computing the desired  
482 effects. These will typically not be available in closed form, but can be computed by Monte  
483 Carlo integration. For illustrative purposes, we present the algorithm for sequential explain-  
484 ability in Algorithm 1. In the supplementary material we provide Monte Carlo integration  
485 algorithms for PMAR and NIP as well. While we do not pursue this here, Monte Carlo inte-  
486 gration can also be implemented using IPW methods (see Robins, 1997; Birmingham et al.,  
487 2003; Shpitser, 2016, for such schemes). The number of Monte Carlo samples should be large  
488 relative to the sample size; in Section 5, we use 100 times the sample size. This appeal to  
489 Monte Carlo to estimate causal effects was initially proposed by Robins (1986) to implement  
490 G-computation. While computationally intensive, post-processing of the MCMC output is  
491 parallelizable and our experience is that the Monte Carlo integration is not a computational  
492 bottleneck. We can avoid repeating these computations for each  $\xi$  by using an informative  
493 prior, providing another advantage to the fully-Bayesian approach.

---

**Algorithm 2** Multiple imputation algorithm for sequential explainability

---

```
1: procedure MI( $M, \mathcal{O}, j$ ) ▷ Inference for  $\mu_j$  using multiple imputation
2:   for  $m = 1, \dots, M$  do
3:     Sample  $\theta \sim \pi(\theta \mid \mathcal{O})$ 
4:     for  $i = 1, \dots, N$  do
5:       if  $R_j^{(i)} = 0$  then
6:         Sample  $Y_j^{(i)} \sim p_\theta(y_j \mid \bar{O}_{j-1}^{(i)}, R_j^{(i)} = 1)$ 
7:       end if
8:     end for
9:     Compute  $\hat{\mu}_j^{(m)} = \frac{1}{N} \sum_{i=1}^N Y_j^{(i)}$ .
10:  end for
11:  Compute  $\hat{\mu}_j$  and  $\hat{\sigma}_{\mu,j}^2$  using the rules for combining inferences under MI.
12: end procedure
```

---

494 **4.2 Multiple imputation**

495 Multiple imputation (MI) proceeds by specifying two, potentially different, models. First,  
496 we use the sampled values  $\theta \sim \pi(\theta \mid \mathcal{O})$  to impute the missing data from  $p_\theta(y_{-r} \mid y_r, r)$   
497 some number  $M > 1$  times. The model  $p_\theta(y, r)$  is referred to as the *imputation model*.  
498 Next, an *analysis model* is specified to compute a point estimate  $\hat{\psi}^{(m)}$  and standard error  
499  $\hat{\sigma}_\psi^{(m)}$  from each of the  $m = 1, \dots, M$  completed datasets  $\mathcal{C}^{(m)}$ . Rubin’s rules (see [Harel and](#)  
500 [Zhou, 2007](#), for a review) are then used to produce a point estimate  $\hat{\psi}$  and standard error  
501  $\hat{\sigma}_\psi$ . The imputation is referred to as *congenial* ([Meng, 1994](#)) when  $\hat{\psi}^{(m)} \approx E(\psi \mid \mathcal{C}^{(m)})$   
502 and  $\hat{\sigma}_\psi^{(m)2} = \text{Var}(\psi \mid \mathcal{C}^{(m)})$ , in which case MI-based inference approximates fully-Bayesian  
503 inference. MI inference may be valid in the absence of congeniality, particularly when the  
504 analysis model is a sub-model of the imputation model. For further discussion of this issue,  
505 see [Rubin \(1996\)](#). For textbook level treatments of multiple imputation, see [Rubin \(1987\)](#) or  
506 [Carpenter and Kenward \(2012\)](#). For an exploration of impact of uncongeniality, see [Daniels](#)  
507 [and Luo \(2017\)](#).

508 The imputation step for MI is operationally similar to the Monte Carlo integration used in  
509 Section 4.1, as it requires simulating from the same conditional distributions. Unlike Monte  
510 Carlo integration, MI only requires imputation of the missing data. Additionally, imputations  
511 can be used with different analysis models. MI is much more practical for large datasets than  
512 fully-Bayesian inference, at the cost of using an approximation. An MI-based algorithm for  
513 estimating  $\mu_j = E(Y_j)$  under sequential ignorability is given in Algorithm 2

514 Extreme caution is required in using MI with partial restrictions in terms of what analysis  
515 models can be used. A minimal condition for MI to be valid is that the analysis model is a  
516 submodel of the imputation model. Hence, when a partial restriction is used, the analysis  
517 model should not identify any part of the joint distribution which is unidentified by the  
518 imputation model. For example, if a marginal restriction identifies the marginals  $p(y_j)$  but  
519 not the joints  $p(y_j, y_k)$ , then the analysis model may also identify  $p(y_j)$  but must not identify



520  $p(y_j, y_k)$ .

521 We remark that there are other approaches to sensitivity analysis which are applied  
522 with multiple imputation. One approach is the so-called “ $\delta$ -adjustment” (Leacy et al., 2017;  
523 Van Buuren, 2012, Section 3.9.1) in which imputations are adjusted, say, by a location shift  $\delta$ .  
524 This approach is ad-hoc and somewhat lacking in transparency regarding what assumptions  
525 it encodes about the missing data, but is highly appealing due to its simplicity. Graphical  
526 methods for conducting a tipping-point analysis are given by Liublinska and Rubin (2014).

## 527 5 Application to the Breast Cancer Prevention Trial data

528 We apply the working model approach described in Section 2.4, using an infinite product-  
529 multinomial mixture (Dunson and Xing, 2009) which is implicitly stratified by treatment,

$$p^*(y, r) = \sum_{k=1}^{\infty} \pi_k \left\{ \prod_{j=1}^J \gamma_{kj}^{r_j} (1 - \gamma_{kj})^{1-r_j} \right\} \left\{ \prod_{j=1}^J \beta_{kj}^{y_j} (1 - \beta_{kj})^{1-y_j} \right\}. \quad (6)$$

530 In the context of missing data, Si and Reiter (2013) applied this model to conduct multi-  
531 ple imputation in large-scale survey data under MAR. For longitudinal responses, various  
532 improvements are possible. One shortcoming of this model is that it does not incorporate  
533 temporal structure; additionally, a model with dependence within the mixture components  
534 would likely perform better (Murray and Reiter, 2016).

535 We give  $\{\pi_k\}_{k=1}^{\infty}$  the stick-breaking prior associated with the Dirichlet process (Sethu-  
536 raman, 1994),  $\pi_k = V_k \prod_{j=1}^{k-1} (1 - V_j)$ ,  $V_k \stackrel{\text{iid}}{\sim} \text{Beta}(1, \alpha)$ . We approximate this by setting  
537  $V_K = 1$  so that  $\pi_k = 0$  for  $k > K$ . For the BCPT data, we set  $K = 50$  and  $\alpha = 1$ . We  
538 view this truncation as a computational concession, leading to an approximation of infer-  
539 ence under  $K = \infty$ ; as pointed out by a referee, one may instead view the truncated model  
540 as a model in its own right, which is parametric rather than nonparametric. We model  
541  $\gamma_{kj} \stackrel{\text{indep}}{\sim} \text{Beta}\{\rho_{\gamma j} a_{\gamma j}, (1 - \rho_{\gamma j}) a_{\gamma j}\}$  and  $\beta_{kj} \stackrel{\text{indep}}{\sim} \text{Beta}\{\rho_{\beta j} a_{\beta j}, (1 - \rho_{\beta j}) a_{\beta j}\}$ . For  $\rho_{\gamma j}$  and  $\beta_{\gamma j}$   
542 we specify independent Uniform(0, 1) priors. Finally, for  $a_{\gamma j}$  and  $a_{\beta j}$  we use a uniform shrink-  
543 age prior, with density  $f_{\sigma}(a) = \sigma/(\sigma + a)^2$  with scale  $\sigma = 15$ . Larger values of  $\sigma$  encourage  
544 heavier shrinkage of the  $\beta_{kj}$ ’s and  $\gamma_{kj}$ ’s towards their means. See Daniels (1999) and Wang  
545 et al. (2010) for motivation and details for the choice of this uniform shrinkage prior.

546 We use MCMC to draw samples of  $\theta = (\pi, \gamma, \beta)$  from the posterior; details are provided  
547 in the supplementary material. We will focus our inference on the effect  $\psi = p(Y_J = 1 \mid Z =$   
548  $1) - p(Y_J = 1 \mid Z = 0)$ , where recall that  $Z = 1$  corresponds tamoxifen and  $Z = 0$  corresponds  
549 to the control. We consider four assumptions which identify  $\psi$ ; the conditional distributions  
550 and algorithms needed are given in the supplementary material. First, we consider MAR  
551 by fitting the  $Y$ -marginal of (6) under ignorability. We also consider PMAR, sequential

552 explainability, and the assumption

$$\begin{aligned}
 & p(Y_J = 1 \mid R = r, Y_r = y_r) \\
 &= \frac{p(Y_J = 1 \mid R = \mathbf{1}, Y_r = y_r)e^\xi}{p(Y_J = 1 \mid R = \mathbf{1}, Y_r = y_r)e^\xi + p(Y_J = 0 \mid R = \mathbf{1}, Y_r = y_r)}.
 \end{aligned} \tag{7}$$

553 Assumption (7) is a nonmonotone, exponentially-tilted, variant of the last-occasion restriction  
 554 of [Birmingham et al. \(2003\)](#). We refer to it as the tilted-last-occasion restriction. In addition  
 555 to the interpretation of the exponential tilting strategy in [Section 3.3](#), the parameter  $\xi$  can  
 556 be interpreted as a location-shift on the logit-scale,

$$p(Y_J = 1 \mid R = r, Y_r = y_r) = \text{expit}[\xi + \text{logit}\{p(Y_J = 1 \mid R = \mathbf{1}, Y_r = y_r)\}],$$

557 where  $\xi$  represents the log-odds ratio of  $[Y_J = 1]$  relative to equivalent individuals with  
 558  $[R = r]$  and  $[R = \mathbf{1}]$ . We posit independent priors for  $\xi$  for each treatment; this has the  
 559 effect of making the posterior variance of  $\psi$  large relative to dependent priors. Alternatively,  
 560 one might take  $\xi$  constant across treatments to encode the belief that the effect of depression  
 561 on missingness does not interact with treatment. To account for the fact that depression  
 562 is expected to be positively correlated with missingness, we set  $\xi \sim \text{Uniform}(0, B)$ . We set  
 563  $B = 0.8$ , corresponding to the belief that it is unlikely that the odds ratio of depression exceeds  
 564  $e^{0.8} \approx 2.2$ . The above specification is made for illustrative purposes and is highly stylized.  
 565 For a more realistic specification which seriously engages with subject-matter expertise, see  
 566 [Wang et al. \(2010\)](#), who elicited informative priors from four subject-matter experts about  
 567 analogous sensitivity parameters  $\xi$ ; none posited values of  $\xi$  larger than 0.8.

568 As a sanity check on the model, it is useful to verify that the posterior gives inferences  
 569 which are consistent with the empirical distribution of the observed data. Let  $\mu_{\text{obs},j} = E(Y_j \mid$   
 570  $R_j = 1)$ . In [Figure 1](#), we compare the inferences based on the posterior distribution of  
 571 the  $\mu_{\text{obs},j}$ 's to the inferences that would be obtained from the standard model-free estimates  
 572  $\hat{\mu}_{\text{obs},j} = \sum_{i=1}^N Y_j^{(i)} R_j^{(i)} / \sum_{i=1}^N R_j^{(i)}$ . We see that the posterior means are essentially identical  
 573 to the the  $\hat{\mu}_{\text{obs},j}$ 's, and the posterior credible intervals agree with the model-free intervals.

574 We report inferences for  $\psi$  obtained using the fully-Bayesian approach in [Figure 2](#); results  
 575 using multiple imputation with a nonparametric analysis model for  $p(y_J)$  are similar, and  
 576 are given in the supplementary material, along with exact numerical results. The most  
 577 striking feature is that inferences obtained under sequential explainability are very different  
 578 from inferences obtained under either PMAR or MAR. First, the magnitude of the effect of  
 579 tamoxifen on depression is much larger under sequential explainability; second, the posterior  
 580 uncertainty is large. This is surprising, as one would expect the additional uncertainty in  $\xi$   
 581 to cause the tilted-last-occasion model to have the most posterior uncertainty.

582 The additional posterior uncertainty can be explained from the fact that most of the  
 583 missingness in the data was monotone. As a result, there is little information about  $p(y_J \mid$   
 584  $\bar{o}_{J-1}, R_J = 1)$  for most missingness patterns. On the other hand, there are many fully

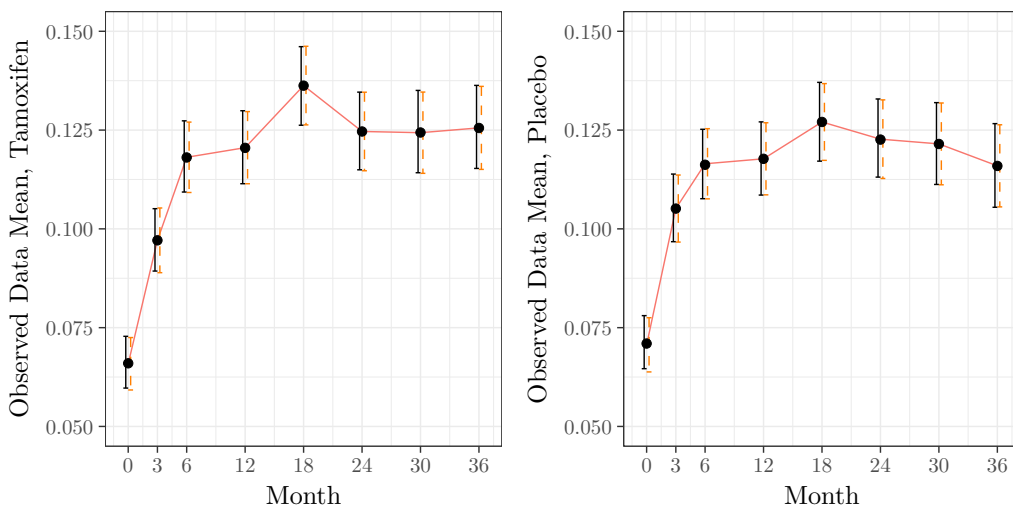


Figure 1: Observed data means over time for the tamoxifen and placebo arms of the study. Dots correspond to the posterior mean using the prior outlined in this section. The line corresponds to the empirical mean of the observed data for each time point. Solid error bars give the 95% credible interval for the observed data mean; dashed error bars given the usual 95% confidence interval based on asymptotic normality of the observed-data means.

585 observed individuals, so there is ample data to estimate  $p(Y_J | Y_r, R = \mathbf{1})$  for all patterns.

586 The fact that sequential explainability produces a larger effect size and leads to substan-  
 587 tively different conclusions is concerning, and necessitates an explanation. Further investigat-  
 588 ion revealed that, among those who were observed at the end of study, but who missed at  
 589 least one visit (roughly 650 individuals per treatment), the difference in depression levels was  
 590 a massive 6%. Moreover, this difference was highly significant, with Fisher’s exact test giving  
 591 a  $P$ -value of 0.002. Under sequential explainability, those who were not observed at the end  
 592 of the study are associated to this group, whereas under PMAR and the tilted-last-occasion  
 593 model these individuals are associated to fully observed individuals. As there was no evidence  
 594 of a difference in depression levels for fully observed subjects ( $P$ -value  $> 0.5$  using Fisher’s  
 595 exact test) the estimate of  $\psi$  is much smaller.

596 Whether PMAR or sequential explainability is more appropriate depends on subject mat-  
 597 ter considerations, as well as the causes of missingness. Regardless, the sensitivity analysis  
 598 led us to find a treatment effect in a sub-population (those who were observed at the end  
 599 of the study, but missed at least one prior visit) which is perhaps itself of interest. Hence,  
 600 in addition to determining the robustness of our inferences, a sensitivity analysis can give  
 601 substantive insight into the relationship between the missingness and the response.

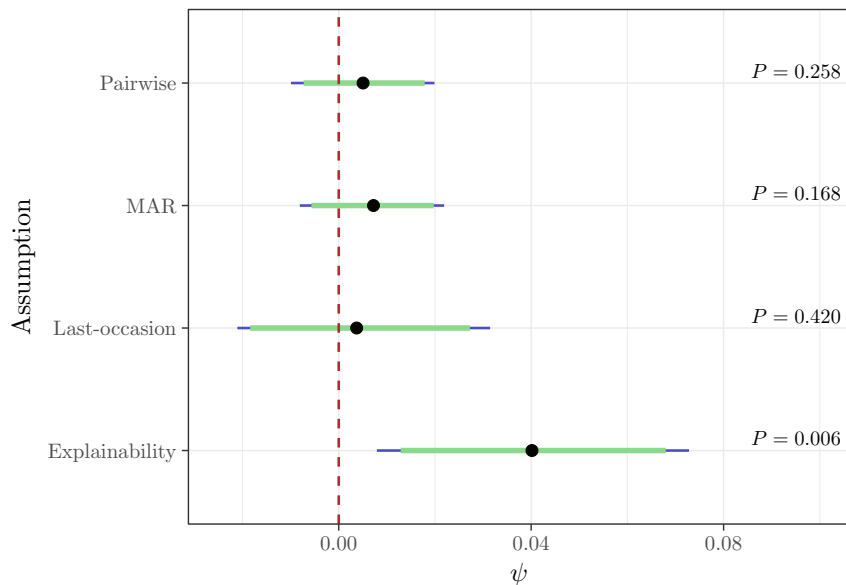


Figure 2: Posterior credible intervals for  $\psi$  under different assumptions. Dots give the posterior mean, green bars give two-sided 90% credible intervals, blue bars give two-sided 95% credible intervals. On the right, the posterior probability  $P = \Pr(\psi < 0)$  is given for each assumption.

## 602 6 Discussion and Open problems

603 In this paper we reviewed identifying restrictions with a focus on recent proposals for non-  
 604 monotone missingness. We also combined a flexible modeling approach for the observed data  
 605 with a variety of identifying restrictions to analyze data from the BCPT.

606 Several interesting avenues of research exist. Auxiliary covariates are often used to impute  
 607 missing outcomes under an assumption that MAR holds only conditional on these additional  
 608 covariates. This is sometimes called A-MAR missingness. The inclusion of such covariates  
 609 can create parameter interpretation problems (Daniels et al., 2014) for categorical outcomes.  
 610 A proposal similar to that introduced here for continuous outcomes and auxiliary covariates  
 611 can be found in Zhou et al. (working paper).

612 This paper has focused on missing outcome data. Handling missing covariate data is also  
 613 of general importance (see, e.g., Ibrahim et al., 1999; Xu et al., 2016; Murray and Reiter,  
 614 2016). One approach to addressing this would be to specify joint Bayesian nonparametric  
 615 models, along with identifying restrictions for the combined vector of outcomes and covariates.  
 616 An interesting problem here is how to specify a parsimonious set of sensitivity parameters  
 617 which will correspond to conditional distributions of both missing outcomes and missing co-  
 618 variates. Multiple identifying restrictions could be used for such analyses, similar to what  
 619 was used in Linero and Daniels (2015) for different types of dropouts (see also Sadinle and  
 620 Reiter, 2017b). Nonignorable missingness for more complex data-structures, such as longitu-

dinal images or networks, remains an underdeveloped area. Much of what has been proposed here could also be used for causal inference. Kim et al. (2017) and Roy et al. (2016) propose Bayesian nonparametric approaches similar to ours in the context of causal mediation and marginal structural models respectively. We are also intrigued by the ICIN restriction as an anchoring assumption, and believe practical methods for performing inference under ICIN would be valuable.

There are few software implementations for conducting sensitivity analysis using identifying restrictions, especially when missingness is nonmonotone. The primary challenge lies in imputing the missing data from the appropriate conditional distributions, as this requires model-specific software. Our R implementation of the multinomial mixture model is provided in the supplementary material. Beyond this, we mention several tools available for sensitivity analysis. Bunouf et al. (2015) provide SAS and R code for implementing pattern-mixture models under monotone missingness and a Gaussian assumption. Scharfstein et al. (2017) provide the R/SAS package SAMON for implementing semiparametric models under monotone missingness. Outside of our proposed framework, proc MI in SAS now supports  $\delta$ -adjustments using the MNAR option, and the SOLAS software package implements the tipping-point strategy of Liublinska and Rubin (2014).

## Acknowledgments

This work was partially supported by NIH grant R01CA183854. The BCPT data was collected under NIH grants U10-CA37377, U10-CA69974, R01AI078835, P30MH086043, and R01HL79457.

## References

- Birmingham, J., Rotnitzky, A., and Fitzmaurice, G. M. (2003). Pattern-mixture and selection models for analysing longitudinal data with monotone missing patterns. *Journal of the Royal Statistical Society, Series B.*, 65:275–297.
- Bunouf, P., Molenberghs, G., Grouin, J.-M., and Thijs, H. (2015). A SAS program combining R functionalities to implement pattern-mixture models. *Journal of Statistical Software*, 68(8).
- Carpenter, J. and Kenward, M. (2012). *Multiple Imputation and its Application*. John Wiley & Sons.
- Carpenter, J., Pocock, S., and Johan Lamm, C. (2002). Coping with missing data in clinical trials: A model-based approach applied to asthma trials. *Statistics in Medicine*, 21(8):1043–1066.

- 654 Cox, D. and Donnelly, C. A. (2011). *Principles of Applied Statistics*. Cambridge University  
655 Press, first edition.
- 656 Creemers, A., Hens, N., Aerts, M., Molenberghs, G., Verbeke, G., and Kenward, M. G.  
657 (2010). A sensitivity analysis for shared-parameter models for incomplete longitudinal  
658 outcomes. *Biometrical Journal*, 52(1):111–125.
- 659 Creemers, A., Hens, N., Aerts, M., Molenberghs, G., Verbeke, G., and Kenward, M. G. (2011).  
660 Generalized shared-parameter models and missingness at random. *Statistical modelling*,  
661 11(4):279–310.
- 662 Daniels, M. J. (1999). A prior for the variance in hierarchical models. *Canadian Journal of*  
663 *Statistics*, 27(3):567–578.
- 664 Daniels, M. J. and Hogan, J. W. (2000). Reparameterizing the pattern mixture model for  
665 sensitivity analyses under informative dropout. *Biometrics*, 56(4):1241–1248.
- 666 Daniels, M. J. and Hogan, J. W. (2008). *Missing Data In Longitudinal Studies*. Chapman  
667 and Hall/CRC, first edition.
- 668 Daniels, M. J. and Linero, A. R. (2015). Bayesian nonparametrics for missing data in longitu-  
669 dinal clinical trials. In *Nonparametric Bayesian Inference in Biostatistics*, pages 423–446.  
670 Springer.
- 671 Daniels, M. J. and Luo, X. (2017). A note on “congeniality” for missing data in the presence  
672 of auxiliary covariates. Technical report.
- 673 Daniels, M. J., Wang, C., and Marcus, B. H. (2014). Fully Bayesian inference under ignorable  
674 missingness in the presence of auxiliary covariates. *Biometrics*, 70(1):62–72.
- 675 Diggle, P. and Kenward, M. G. (1994). Informative drop-out in longitudinal data analysis.  
676 *Applied Statistics*, 43:49–73.
- 677 Dunson, D. B. and Perreault, S. D. (2001). Factor analytic models of clustered multivariate  
678 data with informative censoring. *Biometrics*, 57(1):302–308.
- 679 Dunson, D. B. and Xing, C. (2009). Nonparametric Bayes modeling of multivariate categorical  
680 data. *Journal of the American Statistical Association*, 104(487):1042–1051.
- 681 Gaskins, J., Daniels, M., and Marcus, B. (2016). Bayesian methods for nonignorable dropout  
682 in joint models in smoking cessation studies. *Journal of the American Statistical Associa-*  
683 *tion*, 111(516):1454–1465.
- 684 Harel, O. and Schafer, J. L. (2009). Partial and latent ignorability in missing-data problems.  
685 *Biometrika*, 96:37–50.

- 686 Harel, O. and Zhou, X.-H. (2007). Multiple imputation: Review of theory, implementation  
687 and software. *Statistics in Medicine*, 26(16):3057–3077.
- 688 Heckman, J. J. (1979). Sample selection bias as a specification error. *Econometrica*, 47:153–  
689 161.
- 690 Henderson, R., Diggle, P. J., and Dobson, A. (2000). Joint modelling of longitudinal mea-  
691 surements and event time data. *Biostatistics (Oxford)*, 1:465–480.
- 692 Hogan, J. W. and Laird, N. M. (1997). Mixture models for the joint distribution of repeated  
693 measures and event times. *Statistics in Medicine*, 16:239–257.
- 694 Ibrahim, J. G., Lipsitz, S. R., and Chen, M.-H. (1999). Missing covariates in generalized  
695 linear models when the missing data mechanism is non-ignorable. *Journal of the Royal*  
696 *Statistical Society: Series B (Statistical Methodology)*, 61(1):173–190.
- 697 Ibrahim, J. G. and Molenberghs, G. (2009). Missing data methods in longitudinal studies: a  
698 review. *Test*, 18(1):1–43.
- 699 Kenward, M. G. (1998). Selection models for repeated measurements with non-random  
700 dropout: an illustration of sensitivity. *Statistics in Medicine*, 17(23):2723–2732.
- 701 Kenward, M. G., Molenberghs, G., and Thijs, H. (2003). Pattern-mixture models with proper  
702 time dependence. *Biometrika*, 90:53–71.
- 703 Kim, C., Daniels, M. J., Marcus, B. H., and Roy, J. A. (2017). A framework for Bayesian  
704 nonparametric inference for causal effects of mediation. *Biometrics*, 73:401–409.
- 705 Leacy, F. P., Floyd, S., Yates, T. A., and White, I. R. (2017). Analyses of sensitivity to the  
706 missing-at-random assumption using multiple imputation with delta adjustment: Appli-  
707 cation to a tuberculosis/hiv prevalence survey with incomplete hiv-status data. *American*  
708 *journal of epidemiology*, 185(4):304–315.
- 709 Lin, H., Liu, D., and Zhou, X.-H. (2010). A correlated random-effects model for normal  
710 longitudinal data with nonignorable missingness. *Statistics in medicine*, 29(2):236–247.
- 711 Linero, A. R. (2017). Bayesian nonparametric analysis of longitudinal studies in the presence  
712 of informative missingness. *Biometrika*, 104(2):327–341.
- 713 Linero, A. R. and Daniels, M. J. (2015). A flexible Bayesian approach to monotone missing  
714 data in longitudinal studies with informative dropout with application to a schizophrenia  
715 clinical trial. *Journal of the American Statistical Association*, 110(1):45–55.
- 716 Little, R. J. (1995). Modeling the drop-out mechanism in repeated-measures studies. *Journal*  
717 *of the American Statistical Association*, 90(431):1112–1121.

- 718 Little, R. J. A. (1993). Pattern-mixture models for multivariate incomplete data. *Journal of*  
719 *the American Statistical Association*, 88:125–134.
- 720 Little, R. J. A. (1994). A class of pattern-mixture models for normal incomplete data.  
721 *Biometrika*, 81:471–483.
- 722 Liublinska, V. and Rubin, D. B. (2014). Sensitivity analysis for a partially missing binary  
723 outcome in a two-arm randomized clinical trial. *Statistics in medicine*, 33(24):4170–4185.
- 724 Meng, X.-L. (1994). Multiple-imputation inferences with uncongenial sources of input. *Sta-*  
725 *tistical Science*, pages 538–558.
- 726 Molenberghs, G., Michiels, B., Kenward, M. G., and Diggle, P. J. (1998). Monotone missing  
727 data and pattern-mixture models. *Statistica Neerlandica*, 52:153–161.
- 728 Murray, J. S. and Reiter, J. P. (2016). Multiple imputation of missing categorical and contin-  
729 uous values via bayesian mixture models with local dependence. *Journal of the American*  
730 *Statistical Association*, 111(516):1466–1479.
- 731 National Research Council (2010). *The Prevention and Treatment of Missing Data in Clinical*  
732 *Trials*. The National Academies Press.
- 733 Robins, J. (1986). A new approach to causal inference in mortality studies with a sustained  
734 exposure periodapplication to control of the healthy worker survivor effect. *Mathematical*  
735 *Modelling*, 7(9-12):1393–1512.
- 736 Robins, J. M. (1997). Non-response models for the analysis of non-monotone non-ignorable  
737 missing data. *Statistics in medicine*, 16(1):21–37.
- 738 Robins, J. M., Rotnitzky, A., and Zhao, L. P. (1995). Analysis of semiparametric regression  
739 models for repeated outcomes in the presence of missing data. *Journal of the American*  
740 *Statistical Association*, 90(429):106–121.
- 741 Rotnitzky, A., Robins, J. M., and Scharfstein, D. O. (1998). Semiparametric regression for  
742 repeated outcomes with non-ignorable non-response. *Journal of the American Statistical*  
743 *Association*, 93:1321–1339.
- 744 Roy, J. (2003). Modeling longitudinal data with nonignorable dropouts using a latent dropout  
745 class model. *Biometrics*, 59:441–456.
- 746 Roy, J., Lum, K. J., and Daniels, M. J. (2016). A Bayesian nonparametric approach to  
747 marginal structural models for point treatments and a continuous or survival outcome.  
748 *Biostatistics*, 18(1):32–47.
- 749 Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63:581–592.
- 750 Rubin, D. B. (1987). *Multiple Imputation for Nonresponse in Surveys*. Wiley.



- 751 Rubin, D. B. (1996). Multiple imputation after 18+ years. *Journal of the American statistical*  
752 *Association*, 91(434):473–489.
- 753 Sadinle, M. and Reiter, J. P. (2017a). Itemwise conditionally independent nonresponse mod-  
754 eling for incomplete multivariate data. *Biometrika*, 104(1):207–220.
- 755 Sadinle, M. and Reiter, J. R. (2017b). Sequential identification of nonignorable missing data  
756 mechanisms. *Statistica Sinica*. To appear.
- 757 Scharfstein, D., Mcdermott, A., Diaz, I., Carone, M., Lunardon, N., and Turkoz, I. (2017).  
758 Global sensitivity analysis for repeated measures studies with informative dropout: A semi-  
759 parametric approach. *Biometrics*. To appear.
- 760 Scharfstein, D., McDermott, A., Olson, W., and Wiegand, F. (2014). Global sensitivity anal-  
761 ysis for repeated measures studies with informative dropout: A fully parametric approach.  
762 *Statistics in Biopharmaceutical Research*, 6(4):338–348.
- 763 Scharfstein, D. O., Rotnitzky, A., and Robins, J. M. (1999). Adjusting for nonignorable  
764 dropout using semiparametric nonresponse models. *Journal of the American Statistical*  
765 *Association*, 94.
- 766 Seaman, S., Galati, J., Jackson, D., and Carlin, J. (2013). What is meant by missing at  
767 random? *Statistical Science*, 28(2):257–268.
- 768 Sethuraman, J. (1994). A constructive definition of Dirichlet priors. *Statistica Sinica*, 4:639–  
769 650.
- 770 Shpitser, I. (2016). Consistent estimation of functions of data missing non-monotonically and  
771 not at random. In *Advances in Neural Information Processing Systems*, pages 3144–3152.
- 772 Si, Y. and Reiter, J. P. (2013). Nonparametric Bayesian multiple imputation for incom-  
773 plete categorical variables in large-scale assessment surveys. *Journal of Educational and*  
774 *Behavioral Statistics*, 38(5):499–521.
- 775 Tchetgen Tchetgen, E. J., Wang, L., and Sun, B. (2016). Discrete choice models for  
776 nonmonotone nonignorable missing data: Identification and inference. *arXiv preprint*  
777 *arXiv:1607.02631*.
- 778 Thijs, H., Molenberghs, G., Michiels, B., Verbeke, G., and Curran, D. (2002). Strategies to  
779 fit pattern-mixture models. *Biostatistics*, 3:245–265.
- 780 Tsiatis, A. A. (2007). *Semiparametric Theory and Missing Data*. Springer.
- 781 Van Buuren, S. (2012). Flexible imputation of missing data.

- 782 Vansteelandt, S., Goetghebeur, E., Kenward, M. G., and Molenberghs, G. (2006). Ignorance  
783 and uncertainty regions as inferential tools in a sensitivity analysis. *Statistica Sinica*,  
784 16:953–979.
- 785 Vansteelandt, S., Rotnitzky, A., and Robins, J. M. (2007). Estimation of regression models for  
786 the mean of repeated outcomes under nonignorable nonmonotone nonresponse. *Biometrika*,  
787 94(4):841–860.
- 788 Wang, C. and Daniels, M. J. (2011). A note on MAR, identifying restrictions, model com-  
789 parison, and sensitivity analysis in pattern mixture models with and without covariates for  
790 incomplete data. *Biometrics*, 67:810–818.
- 791 Wang, C., Danies, M. J., Scharfstein, D. O., and Land, S. (2010). A Bayesian shrinkage model  
792 for incomplete longitudinal binary data with application to the breast cancer prevention  
793 trial. *Journal of the American Statistical Association*, 105:1333–1346.
- 794 Wu, M. C. and Carroll, R. J. (1988). Estimation and comparison of changes in the presence  
795 of informative right censoring by modeling the censoring process. *Biometrics*, 45:175–188.
- 796 Xu, D., Daniels, M. J., and Winterstein, A. G. (2016). Sequential BART for imputation of  
797 missing covariates. *Biostatistics*, 17(3):589–602.