

Evaluation of Multiple Imputation Approaches for Handling Missing Covariate Information in a Case-Cohort Study.

Melissa Middleton (✉ melissa.mddleton@mcri.edu.au)

University of Melbourne

Cattram Nguyen

Murdoch Children's Research Institute

Margarita Moreno-Betancur

University of Melbourne

John B Carlin

Murdoch Children's Research Institute

Katherine J Lee

Murdoch Children's Research Institute

Research Article

Keywords: Multiple Imputation, Case-Cohort Study, Simulation Study, Missing Data, Unequal Sampling Probability, Inverse Probability Weighting

Posted Date: October 6th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-922973/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Title:** Evaluation of multiple imputation approaches for handling missing covariate
2 information in a case-cohort study.

3

4 **Authors:** Melissa Middleton^{1,2}, Cattram Nguyen^{1,2}, Margarita Moreno-Betancur^{1,2}, John B
5 Carlin^{1,2}, Katherine J Lee^{1,2}

6 1. Clinical Epidemiology & Biostatistics Unit, Murdoch Children's Research Institute, Melbourne
7 Australia

8 2. Department of Paediatrics, The University of Melbourne, Parkville, Australia

9

10 **Corresponding Author:** Melissa Middleton (melissa.middleton@mcri.edu.au)

11 Murdoch Children's Research Institute, Royal Children's Hospital, 50 Flemington Rd,
12 Parkville, VIC 3052

13 **Abstract:**

14 **Background** In case-cohort studies a random subcohort is selected from the inception cohort
15 and acts as the sample of controls for several outcome investigations. Analysis is conducted
16 using only the cases and the subcohort, with inverse probability weighting (IPW) used to
17 account for the unequal sampling probabilities resulting from the study design. Like all
18 epidemiological studies, case-cohort studies are susceptible to missing data. Multiple
19 imputation (MI) has become increasingly popular for addressing missing data in
20 epidemiological studies. It is currently unclear how best to incorporate the weights from a
21 case-cohort analysis in MI procedures used to address missing covariate data.

22 **Method** A simulation study was conducted with missingness in two covariates, motivated by
23 a case study within the Barwon Infant Study. MI methods considered were: using the
24 outcome, a proxy for weights in the simple case-cohort design considered, as a predictor in
25 the imputation model, with and without exposure and covariate interactions; imputing
26 separately within each weight category; and using a weighted imputation model. These
27 methods were compared to a complete case analysis (CCA) within the context of a standard
28 IPW analysis model estimating either the risk or odds ratio. The strength of associations,
29 missing data mechanism, proportion of observations with incomplete covariate data, and
30 subcohort selection probability varied across the simulation scenarios. Methods were also
31 applied to the case study.

32 **Results** There was similar performance in terms of relative bias and precision with all MI
33 methods across the scenarios considered, with expected improvements compared with the
34 CCA. Slight underestimation of the standard error was seen throughout but the nominal level
35 of coverage (95%) was generally achieved. All MI methods showed a similar increase in

36 precision as the subcohort selection probability increased, irrespective of the scenario. A
37 similar pattern of results was seen in the case study.

38 **Conclusions** How weights were incorporated into the imputation model had minimal effect
39 on the performance of MI; this may be due to case-cohort studies only having two weight
40 categories. In this context, inclusion of the outcome in the imputation model was sufficient to
41 account for the unequal sampling probabilities in the analysis model.

42 **Key Words:** Multiple Imputation, Case-Cohort Study, Simulation Study, Missing Data,
43 Unequal Sampling Probability, Inverse Probability Weighting.

44 **BACKGROUND**

45 Epidemiological studies often collect large amounts of data on many individuals. Some of
46 this information may be costly to analyse, for example biological samples. Furthermore, if
47 there are a limited number of cases, data on all non-cases may provide little additional
48 information to that provided by a subset [1]. In this context, investigators may opt to use a
49 case-cohort study design, in which background covariate data and outcomes are collected on
50 all participants and more costly exposures (e.g. metabolite levels) are collected on a smaller
51 subset. An example of a cohort study adopting the case-cohort design is the Barwon Infant
52 Study (BIS). This is a population-derived cohort study with a focus on non-communicable
53 diseases and the biological processes driving them. Given that a number of investigations
54 within BIS involve exposures collected through costly biomarker and metabolite analysis, for
55 example serum vitamin D levels, the case-cohort design was implemented to minimise cost
56 [2].

57 In the case-cohort design, a subset of the full cohort, hereafter termed the subcohort, is
58 randomly selected from the inception cohort. This subcohort is used as the sample of controls
59 for all subsequent investigations, with exposure data collected from the subcohort and all

60 cases [3]. In such a study, the analysis is conducted on the subcohort and cases only, resulting
61 in an unequal probability of selection into the analysis, with cases having probability of
62 selection equal to 1 and non-case subcohort members having a probability of selection less
63 than 1. This unequal sampling should be accounted for in the analysis so as to avoid bias
64 induced due to the oversampling of cases [4].

65 One way to view the case-cohort design, and to address the unequal sampling, is to treat it as
66 a missing data problem, where the exposure data is ‘missing by design’. Standard practice in
67 the analysis of case-cohort studies is to handle this missing exposure data using inverse
68 probability weighting (IPW) based on the probability of selection into the analysis [1].

69 Additionally, case-cohort studies may be subject to unintended missing data in the covariates,
70 and multiple imputation (MI) may be applied to address this missing data.

71 MI is a two-stage procedure in which missing values are first imputed by drawing plausible
72 sets of values from the posterior distribution of the missing data given the observed data, to
73 form multiple, say $m > 1$, complete datasets. In the second stage, the analysis is conducted on
74 each of these m datasets as though they were fully observed, producing m estimates of the
75 target estimands. An overall estimate for the parameter of interest is produced, along with its
76 variance, using Rubin’s rules [5]. If the imputation model is appropriate, and the assumptions
77 on the missing data mechanism hold, then the resulting estimates are unbiased with standard
78 errors (SE) that reflect not only the variation of the data but also the uncertainty in the
79 missing values [6].

80 When conducting MI, there are two general approaches that can be used to generate the
81 imputed datasets when there is multivariate missingness: joint modelling, most commonly
82 multivariate normal imputation (MVNI), and fully conditional specification (FCS). Under
83 MVNI the missing covariates are assumed to jointly follow a multivariate normal distribution

84 [7]. In contrast, FCS uses a series of univariate imputation models, one for each incomplete
85 covariate, and imputes missing values for each variable by iterating through these models
86 sequentially [8]. To obtain valid inferences, careful consideration must be made when
87 constructing the imputation model such that it incorporates all features of the analysis model,
88 in order to ensure compatibility between the imputation and analysis model [9, 10].

89 In simple terms, to achieve compatibility, the imputation model must at least include all the
90 variables in the analysis model, and in the same form. It may also include additional
91 variables, termed auxiliary variables, which can be used to improve the accuracy of the
92 imputed values and may decrease bias if the auxiliary variables are strong predictors of
93 missingness [11]. In the context of a case-cohort study where the target estimand is the
94 coefficient for the risk ratio (RR) estimated from a log-binomial model with IPW to address
95 unequal probability of selection, two key features should be reflected in the imputation
96 model; 1) the assumed distribution of the outcome, given the exposure and covariates (i.e.
97 log-binomial), and 2) the weights. It is currently unclear how best to incorporate weights into
98 MI in the context of a case-cohort analysis.

99 It has been previously shown that ignoring weights during MI can introduce bias into the
100 point estimates and estimated variance produced through Rubin's rules in the context of an
101 IPW analysis model [12, 13]. Various approaches to incorporate weights into MI have been
102 proposed in the literature. Previous work from Marti and Chavance [14] in a survival analysis
103 of case-cohort data suggests that simply including the outcome, as a proxy for the weights, in
104 the imputation model may be sufficient for incorporating the weights. In the case-cohort
105 setting we are considering, there are only two distinct weights, representing the unequal
106 selection for cases and controls, so the weighting variable is completely collinear with the
107 outcome. An additional approach is to include weights and an interaction between the

108 weights and all of the variables in the analysis model in the imputation model. Carpenter and
109 Kenward [12] illustrated that this corrected for the bias seen when the weights are ignored in
110 the imputation model. One difficulty with this approach is that it may be infeasible if there
111 are several incomplete variables. Another drawback is the increased number of parameters to
112 be estimated during imputation. Another potential approach is to use stratum-specific
113 imputation, in which missing values for cases and non-case subcohort members are imputed
114 separately. While many studies have compared MI approaches in the case-cohort setting [14-
115 17], these are in the context of a time-to-event endpoint and predominantly considered MI
116 only to address the missing exposure due to the design. Keogh [16] considered additional
117 missingness in the covariates, but this was in the context of a survival analysis where weights
118 were dependent on time. While there are many approaches available, it is unclear how these
119 perform in a simple case-cohort context with missing covariate data.

120 The aim of this study was to compare the performance of a range of possible methods for
121 implementing MI to handle missing covariate data in the context of a case-cohort study where
122 the target analysis uses IPW to estimate the i) RR and ii) odds ratio (OR). Whilst the common
123 estimand in case-cohort studies is the RR due to the ability to directly estimate this quantity
124 without the rare-disease assumption [18], we have chosen to additionally consider the target
125 estimand being the coefficient for the OR as this is still a commonly used estimand. The
126 performance of the MI approaches was explored under a range of scenarios through the use
127 of a simulation study closely based on a case study within BIS, and application of these
128 methods to the BIS data. The ultimate goal was to provide guidance on the use of MI for
129 handling covariate missingness in the analysis of case-cohort studies.

130 The paper is structured as follows. We first introduce the motivating example, a case-cohort
131 investigation within BIS, and the target analysis models used for this study. This is followed

132 by a description of the MI approaches to be assessed and details of a simulation study
133 designed to evaluate these approaches based on the BIS case study. We then apply these
134 approaches to the BIS case study. Finally, we conclude with a discussion.

135 **METHODS**

136 **Motivating Example**

137 The motivating example for this study comes from a case-cohort investigation within BIS
138 [19]. A full description of BIS can be found elsewhere [2]. Briefly, it is a population-derived
139 longitudinal birth cohort study of infants recruited during pregnancy (n=1074). The research
140 question focused on the effect of vitamin D insufficiency (VDI) at birth on the risk of food
141 allergy at one-year. Cord blood was collected and stored after birth, and the children were
142 followed up at one-year. During this review, the infant's allergy status to five common food
143 allergens (cow's milk, peanuts, egg, sesame and cashew) was determined through a
144 combination of a skin prick test and a food challenge. Of those who completed the one-year
145 review (n=894, 83%), a random subcohort was selected (n=274), with a probability of
146 approximately 0.31. The exposure, VDI, was defined as 25(OH)D₃ serum metabolite levels
147 below 50nM and was measured from those with a confirmed food allergy at one-year and
148 those who were selected into the subcohort.

149 The planned primary analysis of the case study was to estimate the RR for the target
150 association using IPW in a binomial regression model adjusted for the confounding variables:
151 family history of allergy (any of asthma, hay fever, eczema, or food allergy in an infant's
152 parent or sibling), Caucasian ethnicity of the infant, number of siblings, domestic pet
153 ownership, and formula feeding at 6 and 12 months. The target analysis of this study adjusted
154 for a slightly different set of confounders to the BIS example. A description of these variables
155 and the amount of missing data in each is shown in Table 1.

156 **Table 1:** Detailed description of case study variables used during simulation and their
 157 distribution within the Barwon Infant Study

Variable	Variable Type	Label	n (%) [*] (N=1074)	n (%) missing
Outcome				
Food Allergy at 1 year (present)	Binary; Present/Absent	<i>foodallergy</i>	61 (7.8)	288 (26.8)
Exposure				
Vitamin D Insufficiency at Birth (present)	Binary; Present/Absent	<i>vdi</i>	149 (44.5)	739 (68.8)
Covariates				
Ethnicity (Caucasian)	Binary; Caucasian/Not Caucasian	<i>cauc</i>	772 (72.1)	3 (0.3)
Maternal Vitamin D Supplements Usage (present)	Binary; Present/Absent	<i>antevd</i>	564 (78.8)	358 (33.3)
Family History of Allergy (present)	Binary; Present/Absent	<i>hxfamall</i>	911 (86.1)	16 (1.5)
Number of Siblings	3-Level Categorical	<i>nsib</i>		0 (0.00)
None			453 (42.2)	
One			383 (35.7)	
Two or more			238 (22.2)	
Family Pet Ownership (present)	Binary; Present/Absent	<i>petown</i>	815 (80.5)	62 (5.8)
Formula Feeding at 6 months [#]	3-Level Categorical	<i>formfeed6</i>		189 (17.6)
Exclusively Breast Fed			429 (46.6)	
Exclusively Formula Fed			320 (34.8)	
Mixed Feeding			171 (18.6)	
Formula Feeding at 12 months [#]	3-Level Categorical	<i>formfeed12</i>		154 (14.3)
Exclusively Breast Fed			271 (30.6)	
Exclusively Formula Fed			354 (40.0)	
Mixed Feeding			260 (29.4)	
Auxiliary				
Maternal Age at Birth		<i>mage</i>	32.1 (4.78)	3 (0.3)
Family SEIFA Classification	3-Level Categorical	<i>seifa</i>		20 (1.9)
Low			268 (25.4)	
Middle			204 (19.4)	
High			582 (55.2)	

158 *Mean and standard deviation given for maternal age; percentage given is exclusive of missing data.

159 #Formula feeding variables were not included in the simulation study

160 SEIFA: Socioeconomic index for area

161 Target Analysis

162 In this study, we focus on two estimands from two different analysis models. Each model
 163 targets the association between VDI and food allergy at one-year, adjusting for confounders.
 164 The first model estimates the adjusted RR using a Poisson regression model with a log-link

165 and a robust error variance [20] to avoid the known convergence issues of the log-binomial
166 model:

$$\begin{aligned} 167 \quad \log(\Pr(\text{foodallergy} = 1)) &= \theta_0 + \theta_1 \text{vdi} + \theta_2 \text{cauc} + \theta_3 \text{petown} \\ 168 &\quad + \theta_4[\text{nsib} = 1] + \theta_5[\text{nsib} = 2] + \theta_6 \text{antevd} \\ 169 &\quad + \theta_7 \text{hxfamall} \end{aligned} \tag{1}$$

170 The RR of interest is $\exp(\theta_1)$. The second target estimand is the adjusted OR for the
171 exposure-outcome association, estimated via a logistic regression model:

$$\begin{aligned} 172 \quad \text{logit}(\Pr(\text{foodallergy} = 1)) &= \beta_0 + \beta_1 \text{vdi} + \beta_2 \text{cauc} + \beta_3 \text{petown} \\ 173 &\quad + \beta_4[\text{nsib} = 1] + \beta_5[\text{nsib} = 2] + \beta_6 \text{antevd} \\ 174 &\quad + \beta_7 \text{hxfamall} \end{aligned} \tag{2}$$

175 The OR of interest is $\exp(\beta_1)$. Estimation for each model uses IPW, where the weights are
176 estimated using the method outlined by Borgan [21] for stratified sampling of the cohort,
177 noting that the oversampling of the cases is a special case of stratified sampling where
178 stratification depends on the outcome. The weight for i th individual can be defined as $w_i = 1$
179 for cases, and n_0/m_0 for non-cases, where n_0 is the number of non-cases in the full cohort
180 and m_0 is the number of non-cases within the subcohort.

181 **MI Methods**

182 Below we outline the four approaches we considered in the BIS case study and simulation
183 study for incorporating the weights into the imputation model. All MI approaches include the
184 outcome, exposure, covariates, and auxiliary variables in the imputation model except where
185 specified. Where imputation has been applied under FCS, binary variables have been imputed
186 from a logistic model. For MVNI, all variables are imputed from a multivariate normal

187 distribution, conditional on all other variables, with imputed covariates included into the
188 analysis as is (i.e. without rounding).

189 *Weight proxy as a main effect*

190 Under this approach, only the analysis and auxiliary variables listed above were included in
191 the imputation model, with the outcome acting as a proxy for the weights due to the
192 collinearity between the outcome and weights. This approach is implemented under both the
193 FCS and MVNI frameworks.

194 *Weight proxy Interactions*

195 The second approach includes two-way interactions between the outcome (as a proxy for the
196 weights) and all other analysis variables in the imputation model. Within FCS, the
197 interactions were included as predictors, with these derived within each iteration of the
198 imputation [22]. Within MVNI interactions were considered as ‘just another variable’ in the
199 imputation model [22].

200 *Stratum-specific Imputation*

201 In the case-cohort setting, where there are only two weight strata, another option is to impute
202 separately within each weight/outcome stratum. Here, the outcome is not included in the
203 imputation model, but rather the incomplete covariates are imputed using a model including
204 the exposure, other covariates and auxiliary variables, for cases and non-cases separately.

205 *Weighted Imputation Model*

206 A final option is to impute the missing values using a weighted imputation model, where the
207 weights are set those used during analysis. This can only be conducted within the FCS
208 framework.

209 The approaches for handling the missing covariate data are summarised in Table 2.

210 **Table 2:** Description of multiple imputation approaches considered to handle missing
 211 covariate data.

Method*	Accommodation of Weighting in MI	MI Framework	Label
Complete case	No imputation completed. Analysis applied to observations with complete covariate data.	N/A	CCA
Weight only	Imputation models include weights (through the outcome) as a predictor of missingness	FCS MVNI	<i>FCS-WO</i> <i>MVNI-WO</i>
Weight interactions	Interaction between outcome (proxy for weight) and exposure/covariates included in imputation model through passive imputation (FCS) or 'just another variable' (MVNI), in addition to outcome as a predictor.	FCS MVNI	<i>FCS-WX</i> <i>MVNI-WX</i>
Stratum specific imputation	Covariates imputed separately by weight status	FCS MVNI	<i>FCS-SS</i> <i>MVNI-SS</i>
Weighted model	Imputation model weighted with inverse probability of selection, outcome included as a predictor.	FCS	<i>FCS-WM</i>

212 *All methods involve using multiple imputation to address the missing covariates, excluding the complete case analysis,
 213 with a weighted analysis model to address the unequal probabilities and missing exposure.

214 FCS: Fully Conditional Specification

215 MVNI: Multivariate Normal Imputation

216 MI: Multiple Imputation

217 **Simulation Study**

218 A simulation study was conducted to assess the performance of each MI approach for
 219 accommodating the case-cohort weights into the imputation model, across a range of
 220 scenarios. Simulations were conducted using Stata 15.1 [23], and data were generated using
 221 models based on the observed relationships in BIS (except where noted), with a cohort size of
 222 1000 observations.

223 *Complete Data Generation*

224 Complete data, comprising the exposure, five confounders and two auxiliary variables, were
 225 generated sequentially, as shown in Figure 1, using the models listed below. A table showing
 226 the parameter values used during the data generation can be found in Additional file 1.

227 i. Caucasian ethnicity

$$228 \text{cauc} \sim \text{Ber}(p) \quad (3)$$

229 ii. Maternal age at birth, in years, (auxiliary variable)

$$230 \text{mage} = \delta_0 + \delta_1 \text{cauc} + \epsilon \quad (4)$$

231 where $\epsilon \sim N(0, \sigma^2)$

232 iii. Socioeconomic Index for Areas (SEIFA) tertile, (auxiliary variable)

$$233 \log \left\{ \frac{\text{Pr}(\text{seifa} = 1)}{\text{Pr}(\text{seifa} = 0)} \right\} = \zeta_0 + \zeta_1 \text{mage} + \zeta_2 \text{cauc} \quad (5)$$

$$234 \log \left\{ \frac{\text{Pr}(\text{seifa} = 2)}{\text{Pr}(\text{seifa} = 0)} \right\} = \eta_0 + \eta_1 \text{mage} + \eta_2 \text{cauc} \quad (6)$$

235 iv. Family history of allergy

$$236 \text{logit}\{\text{Pr}(\text{hxfamall} = 1)\} = \iota_0 + \iota_1 \text{cauc} \quad (7)$$

237 v. Number of siblings

$$238 \log \left\{ \frac{\text{Pr}(\text{nsib} = 1)}{\text{Pr}(\text{nsib} = 0)} \right\} = \kappa_0 + \kappa_1 \text{mage} + \kappa_2 \text{cauc} + \kappa_3 [\text{seifa} = 1] \\ 239 + \kappa_4 [\text{seifa} = 2] + \kappa_5 \text{hxfamall} \quad (8)$$

$$240 \log \left\{ \frac{\text{Pr}(\text{nsib} = 2)}{\text{Pr}(\text{nsib} = 0)} \right\} = \lambda_0 + \lambda_1 \text{mage} + \lambda_2 \text{cauc} + \lambda_3 [\text{seifa} = 1] \\ 241 + \lambda_4 [\text{seifa} = 2] + \lambda_5 \text{hxfamall} \quad (9)$$

242 vi. Pet ownership and antenatal vitamin D supplement usage

$$243 \text{logit}\{\text{Pr}(\text{petown} = 1)\} = \rho_0 + \rho_1 \text{mage} + \rho_2 \text{cauc} + \rho_3 [\text{seifa} = 1] \\ 244 + \rho_4 [\text{seifa} = 2] + \rho_5 \text{hxfamall} \\ 245 + \rho_6 [\text{nsib} = 1] + \rho_7 [\text{nsib} = 2] \quad (10)$$

$$246 \text{logit}\{\text{Pr}(\text{antevd} = 1)\} = \phi_0 + \phi_1 \text{mage} + \phi_2 \text{cauc} + \phi_3 [\text{seifa} = 1] \\ 247 + \phi_4 [\text{seifa} = 2] + \phi_5 \text{hxfamall} \\ 248 + \phi_6 [\text{nsib} = 1] + \phi_7 [\text{nsib} = 2] \quad (11)$$

249 vii. The exposure, VDI

$$\begin{aligned} 250 \text{ logit}\{\text{Pr}(\text{vdi} = 1)\} &= \psi_0 + \psi_1 \text{mage} + \psi_2 \text{cauc} + \psi_3 [\text{seifa} = 1] \\ 251 &+ \psi_4 [\text{seifa} = 2] + \psi_5 \text{hxfamall} + \psi_6 [\text{nsib} = 1] \\ 252 &+ \psi_7 [\text{nsib} = 2] + \psi_8 \text{petown} + \psi_9 \text{antevd} \end{aligned} \quad (12)$$

253 -----INSERT FIGURE 1 HERE-----

254 Finally, the outcome, food allergy at one-year, was generated from a Bernoulli distribution
255 with a probability determined by either model (1) or model (2) so the target analysis was
256 correctly specified. In these models we set $\theta_1 = \log(RR_{adj}) = \log(1.16)$ and $\beta_1 =$
257 $\log(OR_{adj}) = \log(1.18)$ as estimated from BIS. Given the weak exposure-outcome
258 association in BIS, we also generated food allergy with an enhanced association where we set
259 $\theta_1 = \beta_1 = \log(2.0)$.

260 An additional extreme data generation scenario was considered as a means to stress-test the
261 MI approaches under more extreme conditions. In this scenario, the associations between the
262 continuous auxiliary variable of maternal age and the exposure, missing covariates, and
263 missing indicator variables were strengthened.

264 *Inducing Missingness*

265 Missingness was introduced into two covariates, antenatal vitamin D usage and pet
266 ownership. Missingness was generated such that $p\%$ of overall observations had incomplete
267 covariate information, with $\frac{p}{2}\%$ having missing data in just one covariate and $\frac{p}{2}\%$ having
268 missing data in both, where p was chosen as either 15 or 30. Three missing data mechanisms
269 were considered: an independent missing data mechanism and two dependent missing data
270 mechanisms.

271 Under the independent missing data mechanism, missingness in each covariate was randomly
272 assigned to align with the desired proportions. Under the dependent missingness mechanisms,
273 an indicator for missingness in pet ownership, M_{petown} , was initially generated from a
274 logistic model (13), followed by an indicator for missingness in antenatal vitamin D usage,
275 M_{antevd} (model (14)).

$$276 \quad \text{logit}(\Pr(M_{\text{petown}} = 1)) = \nu_0 + \nu_1 \text{foodallergy} + \nu_2 \text{cauc} + \nu_3 \text{mage} \quad (13)$$

$$277 \quad \text{logit}(\Pr(M_{\text{antevd}} = 1)) = \tau_0 + \tau_1 \text{foodallergy} + \tau_2 \text{cauc} + \tau_3 \text{mage} + \tau_4 M_{\text{petown}} \quad (14)$$

278 The parameters, ν_0 and τ_0 , and τ_4 were iteratively chosen until the desired proportions of
279 missing information were obtained.

280 The two dependent missing mechanisms differed in the strength of association between
281 predictor variables and the missing indicators. The first mechanism used parameter values set
282 to those estimated in BIS (termed Dependent Missingness – Observed, or *DMO*). The second
283 used an enhanced mechanism where the parameters values were doubled (termed Dependent
284 Missingness – Enhanced, or *DME*). The missing data directed acyclic graph (m-DAG)
285 corresponding to the dependent missingness mechanisms is shown in Figure 1. The parameter
286 values used to induce missingness under the dependent missingness mechanisms can be
287 found in Additional file 1.

288 To mimic the case-cohort design, a subcohort was then randomly selected using one of three
289 probabilities of selection (0.20, 0.30, 0.40). The exposure, VDI, was set to missing for
290 participants without the outcome and who had not been selected into the subcohort.

291 Overall, we considered 78 scenarios (2 data generation processes, 2 exposure-outcome
292 associations, 3 missing data mechanisms, 2 incomplete covariate proportions, and 3
293 subcohort selection probabilities, plus another 6 scenarios under extreme conditions).

294 *Evaluation of MI approaches*

295 For each scenario, the MI approaches outlined above were applied and 30 imputed datasets
296 generated, to match the maximum proportion of missing observations. The imputed datasets
297 were analysed using IPW with the corresponding target analysis model. A complete-data
298 analysis (with no missing data in the subcohort) and a complete-case analysis (CCA) were
299 also conducted for comparison. Performance was measured in terms of relative bias, the
300 empirical and model-based SE, and coverage probability of the 95% confidence interval for
301 the target estimand, the effect of VDI on food allergy (θ_1 in model (1) and β_1 in model (2)).
302 In calculating the performance measures, the “true value” was taken to be the value used
303 during data generation. We also report the Monte Carlo standard error (MCSE) for each
304 performance measure. For each scenario we generated 2000 simulations. With 2000
305 simulations, the MCSE for a true coverage of 95% would be 0.49%, and we can expect the
306 estimated coverage probability to fall between 94.0% and 96.0% [24]. To evaluate the
307 methods, 2200 simulations were generated, with the first 2000 simulations to produce results
308 for all methods used for comparison.

309 *Bias in RR estimation*

310 Incompatibility between the imputation and analysis model may arise due to the imputation
311 of missing values from a linear or logistic model when the analysis targets the RR [25]. To
312 explore the bias introduced into the point estimate in this context, MI was conducted on the
313 full cohort with completely observed exposure (i.e. before data were set to missing by design)
314 and analysed without weighting. This was conducted to understand the baseline level bias,

315 prior to the introduction of weighting. Results for this analysis are presented in Additional
316 file 2.

317 **Case Study**

318 Each of the MI methods were also applied to the target analyses using BIS data. For
319 consistency with the simulation study, the analysis was limited to observations with complete
320 outcome and exposure data (n=246). In the case study there were also missing values in the
321 covariates, Caucasian ethnicity (1%) and family history of allergy (1%), and the auxiliary
322 variable SEIFA tertiles (2%), which were imputed alongside pet ownership (1%) and
323 antenatal vitamin D usage (23%). For the FCS approaches, all variables were imputed using a
324 logistic model, except for SEIFA tertile, which was imputed using an ordinal logistic model.
325 Imputed datasets were analysed under each target analysis model with weights of 1 for cases
326 and $(0.31)^{-1}$ for non-case subcohort members. A CCA was also conducted.

327 **RESULTS**

328 Given that the pattern of results were similar across the range of scenarios, we describe the
329 results for the 6 scenarios under extreme conditions (enhanced exposure-outcome association,
330 30% missing covariates under *DME*, and enhanced auxiliary associations). The results for the
331 remaining scenarios are provided in Additional file 3.

332 Across the 2200 simulations, only FCS-WX and FCS-SS had convergence issues (i.e.,
333 successfully completing the analysis without non-convergence of the imputation procedure or
334 numerical issues in the estimation). The rate of non-convergence was greatest for FCS-WX
335 with the smallest subcohort size (probability of selection = 0.2), with 4.0% of simulations
336 under RR estimation and 3.2% under OR estimation failing to converge. Less than 0.2% of

337 simulations failed to converge for FCS-SS under any combination of estimand and subcohort
338 probability of selection.

339 Figure 2 shows the relative bias in the estimate of the association (RR or OR) between VDI
340 and food allergy at one-year for each scenario considered under extreme conditions. The
341 largest bias for the large sample complete-data analysis occurred for the coefficient of the OR
342 and the largest subcohort size at 1.5%, with the MCSE range covering a relative bias of 0%.
343 In all scenarios shown, the CCA resulted in a large relative bias, ranging between 10% and
344 20%. All MI approaches reduced this bias drastically, irrespective of estimand and subcohort
345 size. When the target estimand was the coefficient for the OR and the smallest subcohort
346 probability was used, all MI approaches were approximately unbiased, with MVNI-WX
347 showing the largest relative bias at 1.4% and FCS-WX showing the least at 0.4%. For the
348 remaining scenarios across both estimands, the complete-data analysis and MI approaches
349 showed a positive bias, with the largest occurring with the log(OR) estimand and a
350 probability of selection of 0.3, where the relative bias was centred around 5%. Overall,
351 minimal differences can be seen between the MI approaches when the estimation targeted the
352 RR. When the target estimand was the log(OR), there was a slight decrease in relative bias
353 for FCS-WX, when compared to other MI approaches, and a slight increase in the relative
354 bias for MVNI approaches, when compared to FCS approaches.

355 -----INSERT FIGURE 2 HERE-----

356 The empirical SE and model-based SE are shown in Figure 3. For most scenarios, we can see
357 the SE has been underestimated in the CCA. There was a slight underestimation of the SE
358 when the subcohort selection probability was 0.3 and the target estimand the log(RR),
359 however, the model-based SE appears to fall within the MCSE intervals for the empirical SE.
360 There appears to be no systematic deviation between the empirical SE and the model-based

361 SE for any scenario. An increase in the precision can be seen as the subcohort size increases,
362 and there is an increase in precision for all MI methods compared to the CCA for any given
363 scenario, as expected.

364 -----INSERT FIGURE 3 HERE-----

365 The estimated coverage probabilities are shown in Figure 4. For a nominal coverage of 95%,
366 all MI approaches have a satisfactory coverage with 95% falling within the MC range for all
367 scenarios, with the exception of the smallest subcohort size with OR estimation. Under this
368 scenario, all MI approaches produce over-coverage, as a result of the point estimate being
369 unbiased and the SE overestimated (but with the average model-based SE falling within the
370 MC range). There is no apparent pattern in the coverage probability across the MI methods,
371 with all methods performing similarly. Results from the CCA showed acceptable levels of
372 coverage.

373 -----INSERT FIGURE 4 HERE-----

374 The results from the case study are shown in Figure 5. Results are consistent with the
375 simulation study in that there is little variation in the estimated association across the MI
376 methods. Unlike the simulation study under extreme conditions, the estimated coefficient is
377 similar in the CCA and the MI approaches. There is an expected recovery of information
378 leading to an increase in precision for MI approaches compared to the CCA.

379 -----INSERT FIGURE 5 HERE-----

380 **DISCUSSION**

381 In this study we compared a number of different approaches for accommodating unequal
382 sampling probabilities into MI in the context of a case-cohort study. We found that how the

383 weights were included in the imputation model had minimal effect on the estimated
384 association or performance of MI which, as expected, outperformed CCA. Results were
385 consistent across different levels of missing covariate information, target estimand and
386 subcohort selection probability.

387 While bias was seen in some scenarios, this was minimal (~5%) and consistent across all MI
388 approaches. We conducted a large sample analysis to confirm the data generation process was
389 correct, given the bias observed in the complete-data analysis, which showed minimal bias.

390 We have therefore attributed the positive bias seen in the simulation study to a finite
391 sampling bias, which was observed for large effect sizes in similar studies [16]. The minimal
392 difference across MI methods seen in the current study may be due to the case-cohort setting
393 having only two weight strata that are collinear with the outcome and all MI approaches
394 including the outcome either directly or indirectly (in the case of stratum-specific
395 imputation). The results of this study complement the work by Marti and Chavance [14] who
396 showed that inclusion of the outcome in the imputation model was sufficient to account for
397 the unequal sampling probabilities in the context of a case-cohort survival analysis. In the
398 case study, the MI approaches performed similarly to the CCA and we believe this is due to
399 the observed weak associations in BIS.

400 Our simulation study was complicated by potential bias due to the incompatibility between
401 the imputation and analysis model when the target analysis estimated a RR through a Poisson
402 regression model. The same would be true if the RR was estimated using a binomial model,
403 as in the case study. We assessed this explicitly through imputation of the full cohort prior to
404 subcohort selection, with results shown in Additional file 2. Minimal bias was seen due to
405 this incompatibility. This may be because we only considered missing values in the
406 covariates, which have been generated from a logistic model. Studies that have shown bias

407 due to this incompatibility had considered missing values in both the outcome and exposure
408 [25].

409 One strength of the current study was that it was based on a real case study, with data
410 generated under a causal structure depicted by m-DAGs informed by subject matter
411 knowledge. This simulation study also examined a range of scenarios; however, it is
412 important to note that not all possible scenarios can be considered, and these results may not
413 extend to scenarios with missingness dependent on unobserved data or with unintended
414 missingness in the exposure or outcome. This study also has a number of limitations. The
415 simulations were conducted under controlled conditions such that the analysis model was
416 correctly specified, and the missing data mechanism was known. Under the specified missing
417 data mechanisms, the estimand was known to be recoverable and therefore MI expected to
418 perform well [26]. The missing data mechanism is generally unknown in a real data setting
419 and results may not be generalisable.

420 Another limitation of this study is that only covariates have been considered incomplete.
421 Often there can be missingness in the outcome (e.g. subcohort members drop-out prior to
422 one-year follow-up and outcome collection) and/or unintended missingness in the exposure
423 (e.g. cord blood not stored for infants selected into the subcohort or with food allergy). This
424 study has also only considered a combination of MI and IPW. There are other analytic
425 approaches that could have been used, for example using weighting to account for the
426 missing covariates as well as the missing data by design, or imputing the exposure in those
427 not in the subcohort and conducting a full cohort analysis. These approaches have been
428 explored in a time-to-event setting [14, 16, 17] but little is known on the appropriateness for
429 the case-cohort setting with a binary outcome. Furthermore, no study to date has considered

430 the scenario of additional exposure missing by chance within the subcohort. The limitations
431 mentioned here offer an avenue for future work.

432 **CONCLUSIONS**

433 When performing MI in the context of case-cohort studies, how unequal sampling
434 probabilities were accounted for in the imputation model made minimal difference in the
435 analysis. In this setting, inclusion of the outcome in the imputation model, which is already
436 standard practice, was a sufficient approach to account for the unequal sampling probabilities
437 incorporated in the analysis model.

438

439 **REFERENCES**

- 440 1. Prentice RL. A case-cohort design for epidemiologic cohort studies and disease
441 prevention trials. *Biometrika*. 1986;73(1):1-11.
- 442 2. Vuillermin P, Saffery R, Allen KJ, Carlin JB, Tang ML, Ranganathan S, et al.
443 Cohort Profile: The Barwon Infant Study. *Int J Epidemiol*. 2015;44(4):1148-60.
- 444 3. Lumley T. *Complex Surveys: A Guide to Analysis using R*. Hoboken, NJ: Wiley;
445 2010.
- 446 4. Cologne J, Preston DL, Imai K, Misumi M, Yoshida K, Hayashi T, et al.
447 Conventional case-cohort design and analysis for studies of interaction. *International*
448 *Journal of Epidemiology*. 2012;41(4).
- 449 5. Rubin D. Inference and missing data. *Biometrika*. 1976;63(3):581-92.
- 450 6. Rubin DB. *Multiple imputation for nonresponse in surveys*: John Wiley & Sons;
451 1987.
- 452 7. Schafer JL. *Analysis of incomplete multivariate data*: Chapman and Hall/CRC;
453 1997.
- 454 8. van Buuren S. Multiple imputation of discrete and continuous data by fully
455 conditional specification. *Statistical methods in medical research*. 2007;16(3):219-
456 42.
- 457 9. Bartlett J, Seaman SR, White IR, Carpenter JR. Multiple imputation of covariates by
458 fully conditional specification: accommodating the substantive model. *Statistical*
459 *methods in medical research*. 2015;24(4):462-87.
- 460 10. Meng X-L. Multiple-imputation inferences with uncongenial sources of input.
461 *Statistical Science*. 1994;9(4):538-58.
- 462 11. Lee K, Simpson J. Introduction to multiple imputation for dealing with missing data.
463 *Respirology*. 2014;19(2):162-7.
- 464 12. Carpenter J, Kenward M. *Multiple imputation and its application*: John Wiley &
465 Sons; 2012.
- 466 13. Kim JK, Brick JM, Fuller WA, Kalton G. On the Bias of the Multiple-Imputation
467 Variance Estimator in Survey Sampling. *Journal of the Royal Statistical Society*
468 *Series B (Statistical Methodology)*. 2006;68(3):509-21.

469 14. Marti H, Chavance M. Multiple imputation analysis of case-cohort studies. *Statistics in Medicine*. 2011;30:1595-607.
470
471 15. Breslow N, Lumley T, Ballantyne CM, Chambless LE, Kulich M. Using the whole
472 cohort in the analysis of case-cohort data. *American journal of epidemiology*.
473 2009;169(11):1398-405.
474 16. Keogh RH, Seaman SR, Bartlett J, Wood AM. Multiple imputation of missing data
475 in nested case-control and case-cohort studies. *Biometrics*. 2018;74:1438-49.
476 17. Keogh RH, White IR. Using Full-cohort data in nested case-control and case-cohort
477 studies by multiple imputation. *Statistics in Medicine*. 2013;32:4021-43.
478 18. Sato T. Risk ratio estimation in case-cohort studies. *Environ Health Perspect*.
479 1994;102 Suppl 8:53-6.
480 19. Molloy J, Koplin JJ, Allen KJ, Tang MLK, Collier F, Carlin JB, et al. Vitamin D
481 insufficiency in the first 6 months of infancy and challenge-proven IgE-mediated
482 food allergy at 1 year of age: a case-cohort study. *Allergy*. 2017;72:1222-31.
483 20. Zhou G. A modified poisson regression approach to prospective studies with binary
484 data. *American Journal of Epidemiology*. 2004;159(7):702-6.
485 21. Borgan O, Langholz B, Samuelsen SO, Goldstein L, Pogoda J. Exposure stratified
486 case-cohort designs. *Lifetime data analysis*. 2000;6:39-58.
487 22. von Hippel PT. How to impute interactions, squares, and other transformed
488 variables. *Sociological Methodology*. 2009;39(1).
489 23. StataCorp. *Stata Statistical Software: Release 15*. In: StataCorp, editor. College
490 Station, TX: StataCorp LLC; 2017.
491 24. Morris TP, White IR, Crowther MJ. Using simulation studies to evaluate statistical
492 methods. *Statistics in Medicine*. 2019:1-29.
493 25. Sullivan TR, Lee K, J., Ryan P, Salter AB. Multiple imputation for handling missing
494 outcome data when estimating the relative risk. *BMC Medical Research Methodology*.
495 2017;17(134).
496 26. Moreno-Betancur M, Lee KJ, Leacy FB, White IR, Simpson JA, Carlin J. Canonical
497 causal diagrams to guide the treatment of missing data in epidemiologic studies.
498 *American Journal of Epidemiology*. 2018;187(12).

499 **ABBREVIATIONS**

500	BIS	Barwon Infant Study
501	CCA	Complete-Case Analysis
502	DME	Dependent Missing – Enhanced
503	DMO	Dependent Missing – Observed
504	FCS	Fully Conditional Specification
505	IPW	Inverse Probability Weighting
506	m-DAG	Missingness Directed Acyclic Graph

507	MCSE	Monte Carlo Standard Error
508	MI	Multiple Imputation
509	MVNI	Multivariate Normal Imputation
510	OR	Odds Ratio
511	RR	Risk Ratio
512	SE	Standard Error
513	SEIFA	Socioeconomic Index for Areas
514	VDI	Vitamin D Insufficiency

515

516 **DECLARATIONS**

517 **Ethics approval and consent to participate**

518 The case study used data from the Barwon Infant Study, which has ethics approval from the
 519 Barwon Health Human Research and Ethics Committee (HREC 10/24). Participating parents
 520 provided informed consent and research methods followed national and international
 521 guidelines.

522 **Consent for publication**

523 Not applicable

524 **Availability of data and materials**

525 The datasets used and/or analysed during the current study are available from the
 526 corresponding author on reasonable request.

527 **Competing interests**

528 The authors declare that they have no competing interests.

529 **Funding**

530 This work was supported by the Australian National Health and Medical Research Council
 531 (Postgraduate Scholarship 1190921 to MM, career development fellowship 1127984 to KJL,
 532 and project grant 1166023). MMB is the recipient of an Australian Research Council

533 Discovery Early Career Researcher Award (project number DE190101326) funded by the
534 Australian Government. MM is funded by an Australian Government Research Training
535 Program Scholarship. Research at the Murdoch Children’s Research Institute is supported by
536 the Victorian Government’s Operational Infrastructure Support Program. The funding bodies
537 do not have any role in the collection, analysis, interpretation or writing of the study.

538 **Authors’ contributions**

539 MM, CN, MMB, JBC and KJL conceived the project and designed the study. MM designed
540 the simulation study and conducted the analysis, with input from co-authors, and drafted the
541 manuscript. KJL, CN, MMB and JBC provided critical input to the manuscript. All of the co-
542 authors read and approved the final version of this paper.

543 **Acknowledgements**

544 The authors would like to thank the Melbourne Missing Data group and members of the
545 Victorian Centre for Biostatistics for providing feedback in designing and interpreting the
546 simulation study. We would also like to thank the BIS investigator group for providing access
547 to the case-study data for illustrative purposes in this work.

548

549 **Figure 1** Missing data directed acyclic graph (m-DAG) depicting the assumed causal
550 structure between simulated variables and missingness indicators under the dependent
551 missing mechanisms. For the independent missing mechanism, the dashed lines are absent.
552 For simplicity, associations between baseline covariates have not been shown.

553 **Figure 2** Relative bias in the coefficient under the extreme scenarios with 30% missing
554 covariate information. Error bars represent 1.96xMonte Carlo standard errors.

555 **Figure 3:** Empirical standard error and model based standard error under the extreme
556 scenarios with 30% missing covariate information. Error bars represent 1.96xMonte Carlo
557 standard errors.

558 **Figure 4:** Coverage probability across 2000 simulations under the extreme scenarios with
559 30% missing covariate information. Error bars represent 1.96xMonte Carlo standard errors.

560 **Figure 5:** Estimated parameter value with 95% confidence interval in case study dataset

Figures

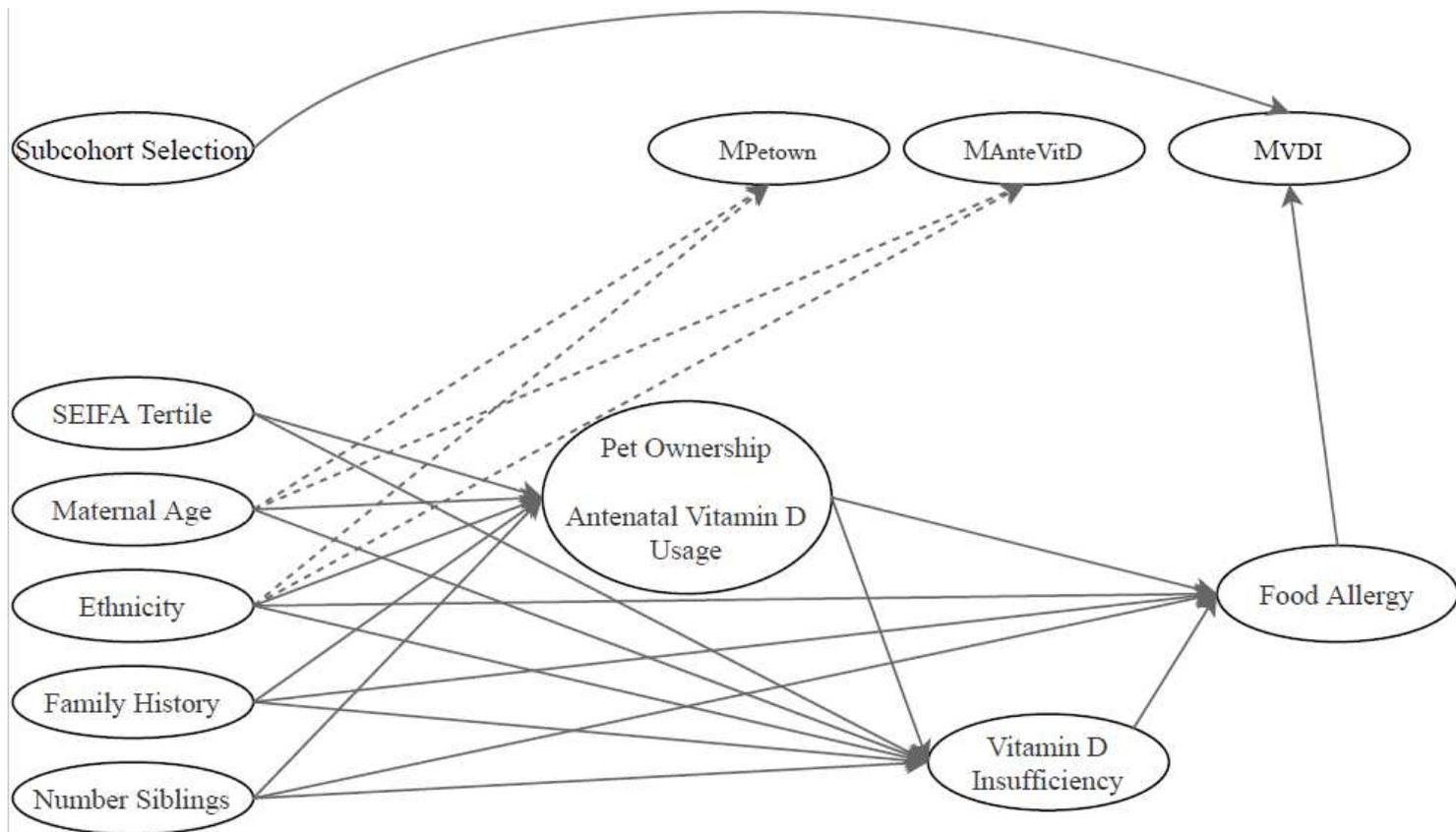


Figure 1

Missing data directed acyclic graph (m-DAG) depicting the assumed causal structure between simulated variables and missingness indicators under the dependent missing mechanisms. For the independent missing mechanism, the dashed lines are absent. For simplicity, associations between baseline covariates have not been shown.

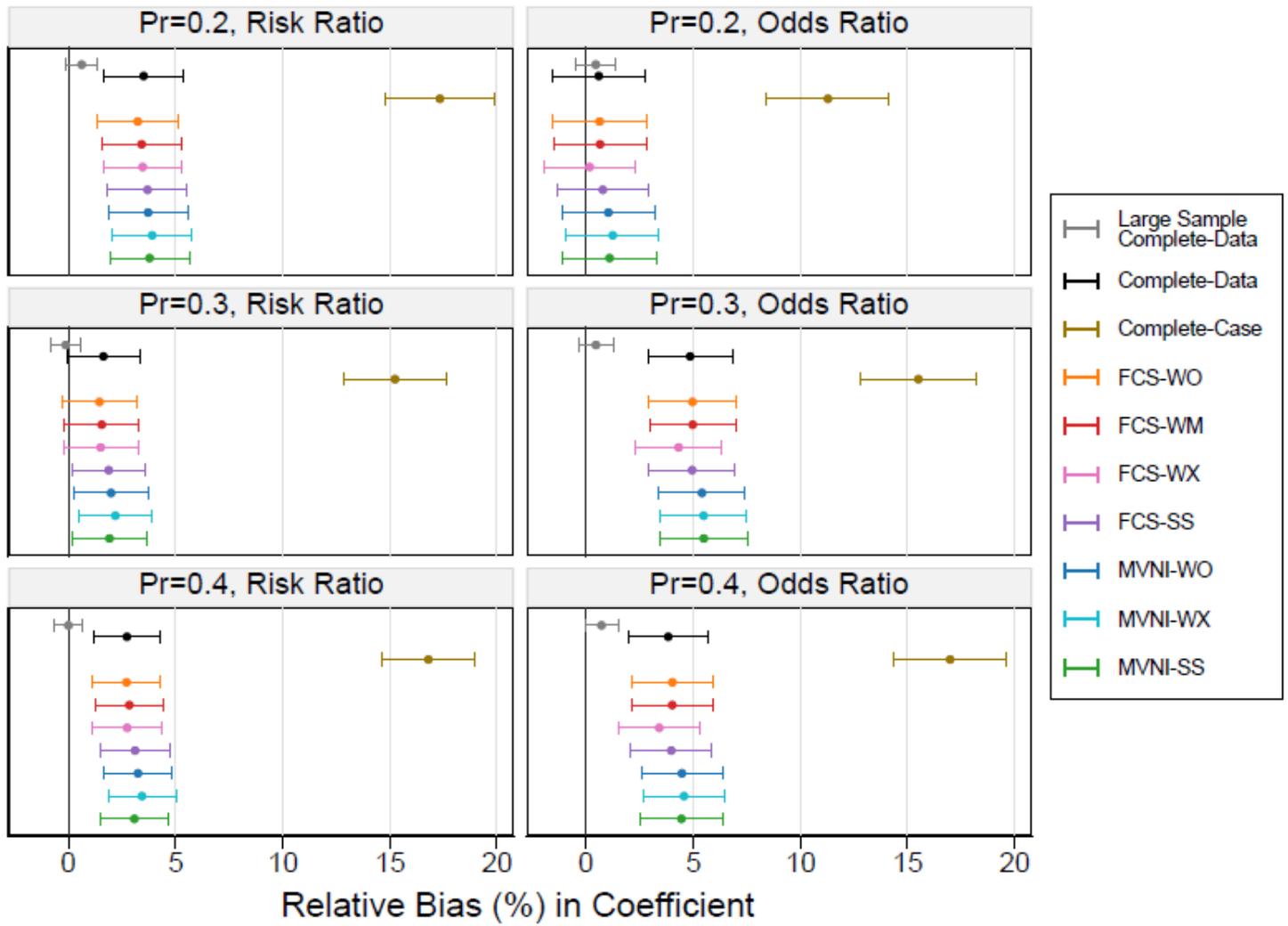


Figure 2

Relative bias in the coefficient under the extreme scenarios with 30% missing covariate information. Error bars represent 1.96xMonte Carlo standard errors.

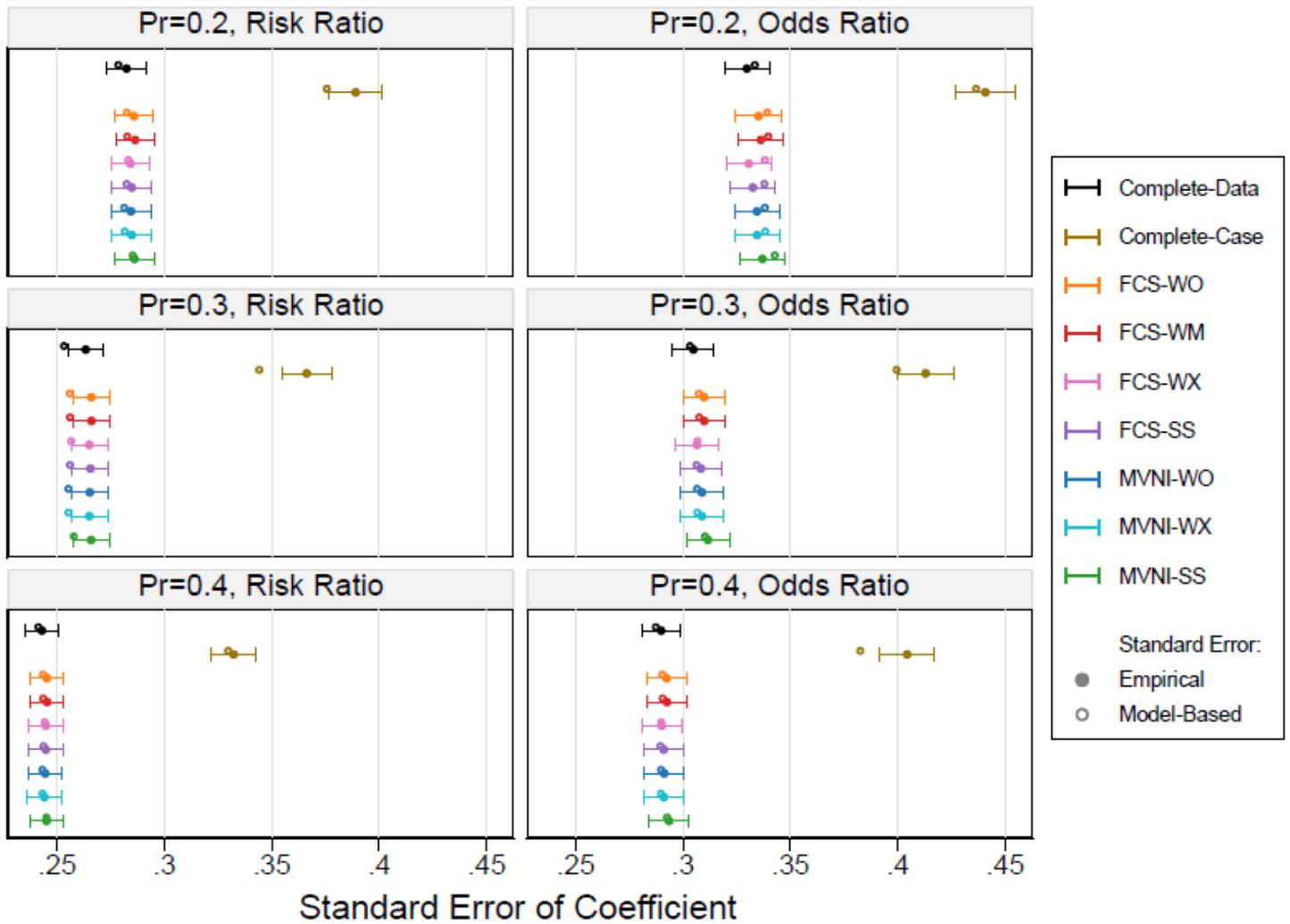


Figure 3

Empirical standard error and model based standard error under the extreme scenarios with 30% missing covariate information. Error bars represent 1.96xMonte Carlo standard errors.

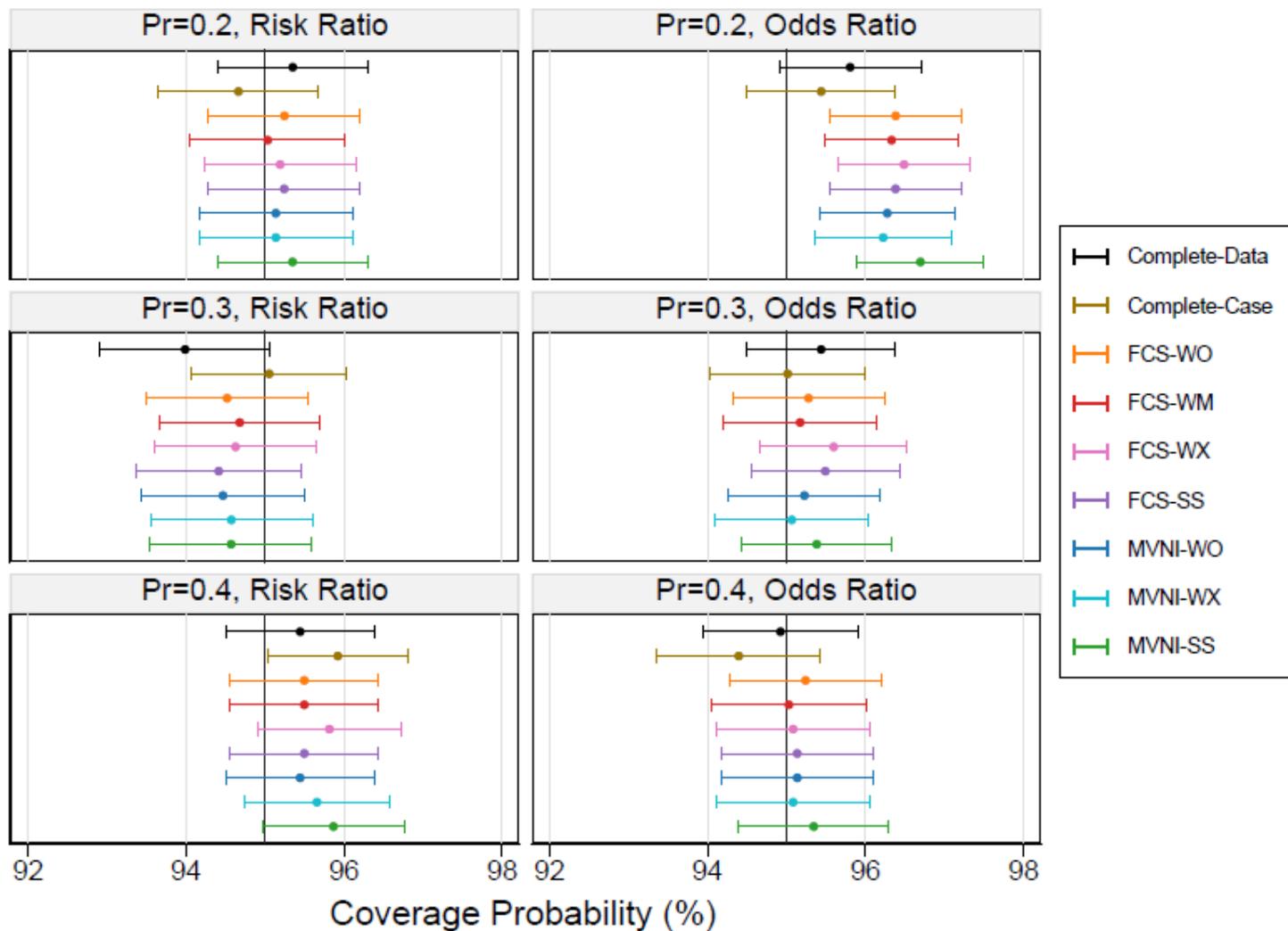


Figure 4

Coverage probability across 2000 simulations under the extreme scenarios with 30% missing covariate information. Error bars represent $1.96 \times$ Monte Carlo standard errors.

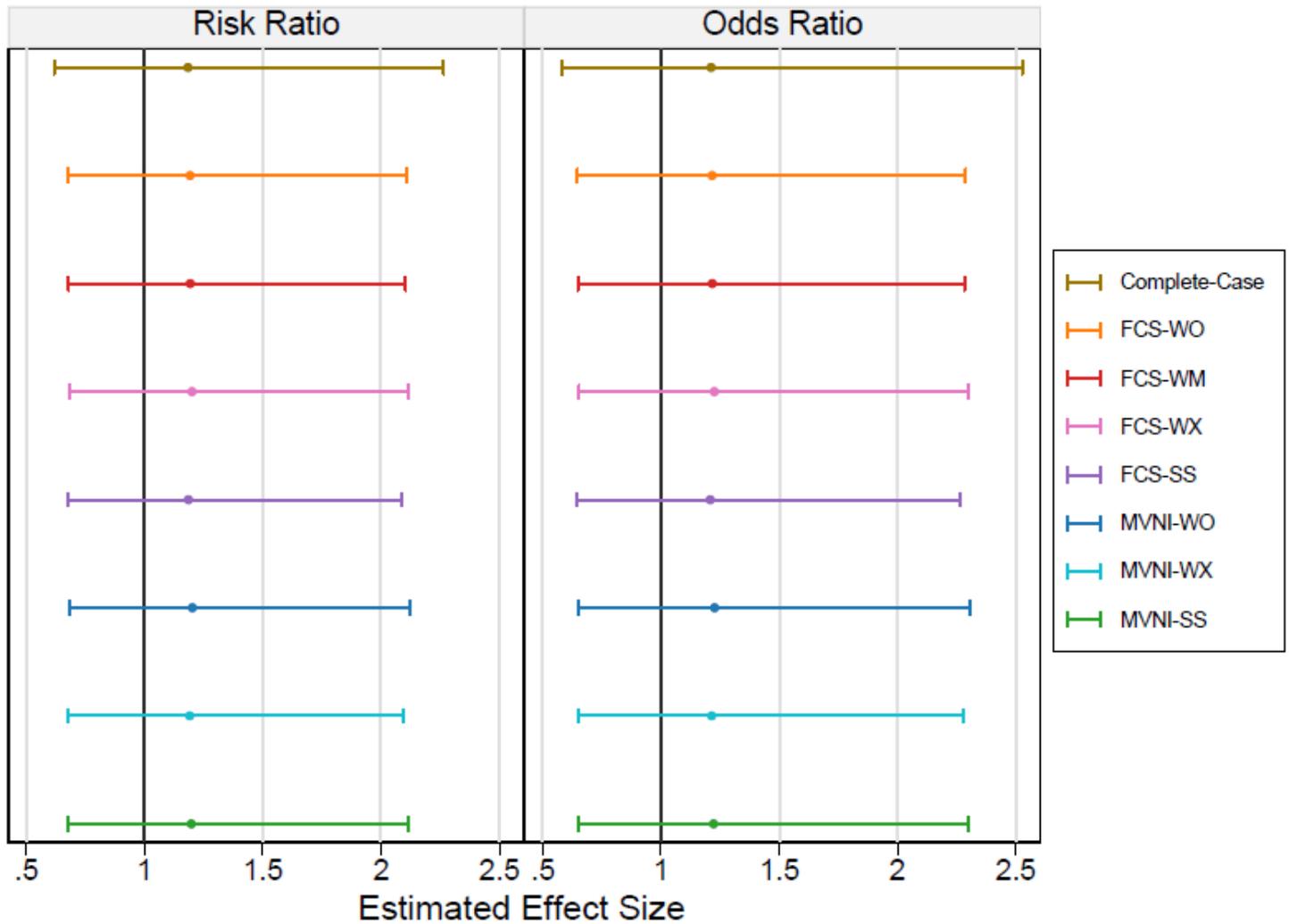


Figure 5

Estimated parameter value with 95% confidence interval in case study dataset

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [AdditionalFile1p1.docx](#)
- [AdditionalFile2p1.docx](#)
- [AdditionalFile3p1.docx](#)