

BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email info.bmjopen@bmj.com

BMJ Open

The handling of missing data with multiple imputation in observational studies that address causal questions: Protocol for a scoping review

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2022-065576
Article Type:	Protocol
Date Submitted by the Author:	10-Jun-2022
Complete List of Authors:	Mainzer, Rheanna; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Moreno-Betancur, Margarita; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Nguyen, Cattram; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Simpson, Julie; University of Melbourne School of Population and Global Health Carlin, John; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Lee, Katherine; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics
Keywords:	EPIDEMIOLOGY, STATISTICS & RESEARCH METHODS, Public health < INFECTIOUS DISEASES

SCHOLARONE™
Manuscripts

The handling of missing data with multiple imputation in observational studies that address causal questions: Protocol for a scoping review

Rheanna M. Mainzer^{*1,2}, Margarita Moreno-Betancur^{1,2}, Cattram D. Nguyen^{1,2}, Julie A. Simpson³, John B. Carlin^{1,2,3}, Katherine J. Lee^{1,2}

1. Clinical Epidemiology and Biostatistics Unit, Murdoch Children's Research Institute, Parkville, Victoria 3052, Australia
2. Department of Paediatrics, The University of Melbourne, Parkville, Victoria 3052, Australia
3. Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, The University of Melbourne, Parkville, Victoria 3052, Australia

*Corresponding author: Rheanna Mainzer; rheanna.mainzer@mcri.edu.au

ABSTRACT

Introduction Observational studies in health-related research often aim to answer causal questions. Missing data are common in such studies and can occur in the exposure, outcome and/or variables used to control for confounding. The standard classification of all missing data as missing completely at random (MCAR), missing at random (MAR) or missing not at random (MNAR), does not allow for a clear assessment of missingness assumptions when missingness arises in more than one variable. This presents challenges for selecting an analytic approach and determining when a sensitivity analysis under plausible alternative missing data assumptions is required. This is particularly pertinent with multiple imputation (MI), which is often justified by assuming data are MAR. The objective of this scoping review is to examine the use of MI in observational studies that address causal questions, with a focus on (i) how missingness assumptions are expressed and assessed, (ii) the connection between missingness assumptions and the use of MI or other approaches for handling missing data, and (iii) the conduct of sensitivity analyses under alternative plausible missingness mechanisms.

Methods and analysis We will systematically review observational studies that aim to answer causal questions using MI, published between January 2019 and December 2021 in five top general epidemiology journals. Studies will be identified using a full text search for the term "multiple imputation". Information extracted from eligible studies will include details about the study characteristics, missing data, missingness assumptions, analysis methods and MI implementation. Systematic review methods will be used to screen, review and extract data. Data will be summarised using descriptive statistics.

Ethics and dissemination Ethics approval is not required for this review because data will be collected only from published studies. The results will be disseminated through a peer reviewed publication and conference presentations.

Registration This protocol is registered on figshare (<https://doi.org/10.6084/m9.figshare.20010497.v1>).

Strengths and limitations of this study

- A targeted review of observational studies published in the five top-ranked epidemiology journals will benchmark the current state of practice for handling multivariable missingness with multiple imputation. Although our targeted review will not include all relevant studies, we expect that included studies will be sufficient to provide insight and general trends on the application and reporting of multiple imputation in observational studies.
- Screening, reviewing and data extraction will be performed systematically, with double data extraction for a subset of articles and any discrepancies resolved by a panel.
- All data and code will be made publicly available, enabling our analysis to be entirely reproducible.
- It is likely that some of the information sought will be unclear or not reported. To accommodate this, we have specified how anticipated challenges with data extraction will be handled if they arise.
- Results from the review will be reported according to best practice, using the Preferred Reporting Items for Systematic reviews and Meta-Analyses extension for Scoping Reviews (PRISMA-ScR).

INTRODUCTION

Observational studies in clinical and health-related research often aim to answer causal questions, i.e. to estimate the effect of an exposure on an outcome.⁽¹⁾ In such studies missing data are common and can occur in the exposure, the outcome and/or the variables used to control for confounding. Restricting statistical analysis to individuals with available data (complete case analysis, CCA) can lead to bias and/or loss of precision in estimates of the average causal effect. ⁽²⁾ Multiple imputation (MI) is a popular and flexible approach for estimating target quantities in the presence of incomplete data.^(3, 4) In the first stage of MI, missing data are imputed multiple times with random draws from the predictive distribution of the missing values given the observed data and a specified imputation model. In the second stage, the statistical analysis of interest is applied to each imputed dataset and the results are combined using Rubin's rules to obtain a single estimate with associated standard error.⁽³⁾

Standard implementations of MI are known to provide consistent estimation of target parameters under certain unverifiable assumptions about the mechanism leading to missing data. These assumptions are usually expressed using Rubin's classification of missing data mechanisms into missing completely at random (MCAR, where the probability of data being missing does not depend on the observed or unobserved data), missing at random (MAR, where the probability of data being missing does not depend on the unobserved data, conditional on the observed data) and missing not at random (MNAR, where the probability of data being missing depends on the unobserved data, even after conditioning on the observed data).⁽⁵⁾ While this framework is useful if missing data occur in a single variable, it is poorly understood and does not allow for a transparent assessment of missingness assumptions when missingness arises in more than one variable.⁽⁶⁾ For example, MAR is a sufficient but not necessary condition for the validity of standard MI estimates.⁽⁷⁾ Further, because one cannot be sure about the true missing data mechanism, sensitivity analyses to examine the robustness of results to alternative plausible missingness mechanisms (hereafter, "sensitivity analyses") are strongly recommended.⁽⁸⁾ As stated by the US National Research Council, "the usefulness of a sensitivity analysis ultimately depends on the transparency and plausibility of the unverifiable assumptions."⁽⁸⁾ The inherent difficulty in assessing missingness assumptions when framed in the traditional MCAR/MAR/MNAR manner and their lack of one-to-one correspondence with analytic approaches in the presence of multivariable missingness leads to further complications when planning and conducting sensitivity analyses.

Most reviews of the handling and reporting of missing data, and the implementation and documentation of MI, have been carried out in the context of randomised controlled trials (RCTs) with missing outcome data.⁽⁹⁻¹⁵⁾ For trials, typically only the outcome variable is incomplete, while the intervention and other key variables are observed for all participants. In this setting where there is missing data in a single variable, the MCAR/MAR/MNAR framework is more transparent and guidance on sensitivity analyses has been well-developed (see, for example, (12, 16)). In contrast, there have been few reviews concerned with how missing data are handled in observational studies where there is the additional complication of multivariable missingness. A review by Mackinnon published in 2010 found that only two (4%) out of 50 non-RCT studies reviewed carried out an additional analysis that was described as a sensitivity analysis.⁽¹⁷⁾ Similarly, Rezvan et al. (2015) found that none of the 30 observational studies reviewed conducted a sensitivity analysis to departures from the missingness assumptions following MI.⁽¹⁸⁾ Even when they are carried out, what is meant by a "sensitivity analysis" is often unclear. Confusion between sensitivity analyses and secondary analyses has been observed, (17, 19) and the logic behind applying MI as a sensitivity analysis to a CCA (or vice versa) is unsound.⁽¹⁷⁾ While the reviews by Mackinnon and Rezvan et al. provide useful insight into the problem, neither focused specifically on observational studies and the issues described above. In addition, subsequent to publication of these reviews there have been important developments in the theory and application of missingness directed acyclic graphs (m-DAGs), also known as m-graphs, a tool for the formulation of causal assumptions in the presence of multivariable missingness.⁽⁷⁾ M-DAGs can aid the depiction and assessment of missingness assumption, which is important since transparency in the assumed causal mechanisms underlying the missing data facilitates the choice of analytical approach.⁽²⁰⁾ Although, it is currently unclear how much m-DAGs are being used in the literature.

1
2
3 The aim of this scoping review is to systematically review the epidemiological literature to examine the use of
4 MI in observational studies that address causal questions, which is typically the focus of such studies even
5 when this may not be very clearly articulated.(21) These studies often face missingness in multiple variables
6 required for analysis. We will examine (i) how missingness assumptions are expressed, (ii) their connection to
7 the justification for the use of MI or other approaches for handling missing data, and (iii) the conduct of
8 sensitivity analyses to alternative plausible missingness mechanisms. We will also examine how MI is
9 implemented. This review will be used to document the current state of practice, to identify areas for
10 improvement of reporting on the handling of missing data with MI in observational studies, and to
11 subsequently develop guidance for researchers.
12

13 **METHODS AND ANALYSIS**

14
15 In this section we provide a full description of the study design, including how articles will be selected, what
16 outcomes will be measured, and how data will be extracted and analysed. The anticipated start date of this
17 review is 13th June 2022 and the anticipated completion date is 30th November 2022.
18

19 **Search strategy**

20
21 We will systematically search five general epidemiology journals for observational studies published between
22 January 2019 and December 2021 that aim to answer at least one causal research question using MI. The
23 general epidemiology journals that will be included in this search are: *International Journal of Epidemiology*,
24 *American Journal of Epidemiology*, *European Journal of Epidemiology*, *Journal of Clinical Epidemiology* and
25 *Epidemiology*. These journals were chosen because they are high ranking, general journals in epidemiology
26 that publish original research from observational studies. As such, articles from these journals should capture
27 the current best practice in the use of MI to handle missing data when answering causal questions using
28 observational data. They have also been used previously in a systematic review of epidemiologic practice.(22)
29 Original research articles will be identified using the full-text search term “multiple imputation” on each
30 journal’s website. This search strategy is similar to that used in previous scoping reviews in this area.(17, 18)
31

32 **Inclusion criteria**

33 We will include original research articles that were published between January 2019 and December 2021, and
34 aim to answer at least one causal question using MI to handle the missing data. We will determine that a study
35 has aimed to answer a causal question if at least one of the following criteria is satisfied:
36

- 37 1. the authors explicitly stated they were estimating a causal effect;
- 38 2. the study estimated an effect that was given (at least implicitly) a causal interpretation, i.e., an
39 interpretation which suggested that intervening on the exposure could change the outcome (e.g.,
40 increasing coffee consumption may be protective against stroke). This will be determined by wording
41 in conclusions and typically signalled by the identification of confounders, the inclusion of a DAG to
42 illustrate causal assumption made in the analysis, and/or analytical approaches incorporating
43 adjustment for confounders (for example, estimating an effect using a regression model that was
44 adjusted for a set of covariates).
45

46
47 All disease areas/medical conditions will be considered and there will be no restrictions on the study
48 participants.
49

50 **Exclusion criteria**

51 Studies will be excluded from the review if they meet any of the following criteria:
52

- 53 • *No causal question.* The article did not aim to answer a causal question, for example, the aim of the
54 study was to validate a predictive model or to estimate a disease burden.
- 55 • *Unclear type of question.* A clear research goal could not be identified. In other words, it was unclear
56 whether the study aimed to answer a descriptive, predictive or causal question.
- 57 • *The analysis did not use MI.*
- 58 • *Methodological research.* The primary purpose of the article was methodological development, for
59 example, using a simulation study to compare the performance of methods or mathematical
60

derivations to develop a new method or model. While these articles often include comprehensive case studies, they may not be representative of published studies that aim to answer causal research questions.

- *Aggregate-level data.* The analysis was based on aggregated data where MI could not be applied at the participant level, as is common in meta-analysis or interrupted time series analysis.
- *Qualitative research.* The article provided a commentary, review, opinion, study protocol, study profile or description only.
- *Trial.* The study intervention was assigned to participants by the trial investigators.

Sample size

We will require at least 100 studies to estimate the percentage of studies with a particular element (e.g., studies that justify their missingness assumptions) to within a maximum margin of error (two standard errors) of 10%. Assuming a prevalence of 50%, this would give a 95% confidence interval from 40% to 60%. For a prevalence greater than or less than 50%, the 95% confidence interval will be narrower. This sample size is similar to the sample size used in the first review of MI in medical research ($n = 99$, (17)), and many of the subsequent reviews in this area (e.g., $n = 103$ in (18), 77 in (12) and 118 in (9)). We expect to identify at least 100 eligible studies given the three-year publication time frame. All eligible studies will be included in the review.

Study selection

The search of the journal databases will be performed by a single researcher (RM). The title, abstract and date of each article will be screened for eligibility. When a decision about the eligibility of an article cannot be reached based on the title, abstract and publication date alone, the full text will be screened for eligibility. A second researcher (CN) will independently screen articles when there is uncertainty about the inclusion criteria. Disagreements about inclusion criteria will be resolved by discussion in meetings with at least three researchers (RM, CN and at least one of JC, JS, KL or MMB).

Data extraction and management

Covidence, a web-based tool for systematic review management, will be used to perform the review.⁽²³⁾ The data extraction questionnaire was developed and tested for use by RM and KL using a sample of 10 articles. All eligible studies will be extracted and reviewed by RM. The supplementary material of all eligible studies will also be extracted and reviewed. We will use double data extraction (performed by CN) for a random selection of 10% of articles and additionally when there is uncertainty about the information being extracted. Discrepancies and uncertainties will be resolved by discussion in meetings with at least three researchers (RM, CN and at least one of JC, JS, KL or MMB).

Outcomes measured

We will extract data pertaining to the study characteristics, the amount of missing data and in which variables, missingness assumptions, methods for handling missing data and implementation of multiple imputation. Data extraction items are summarised in Table 1. Because we anticipate difficulties in extracting some items (such as the percentage of complete cases), in Supplementary Table 1 we list potential challenges in extracting data and any assumptions or simplifications that will be made if these challenges arise. Any post-hoc assumptions or simplifications for unanticipated challenges will be recorded and reported as part of the analysis.

Table 1. Summary of items to be extracted from each article.

Category	Summary of data extraction items
Study characteristics	<ul style="list-style-type: none"> • Title • Authors • Publication date • Journal • Type of study design
Missing data	<ul style="list-style-type: none"> • Percentage of complete cases

	<ul style="list-style-type: none"> Percentage of missing values in the exposure and outcome Number of incomplete covariates
Missingness assumptions	<ul style="list-style-type: none"> Statement of missingness data assumptions (including whether the study used m-DAGs or the MCAR/MAR/MNAR framework) Justification of missingness assumptions
Analysis methods	<ul style="list-style-type: none"> The primary analysis method used to answer the key causal question, e.g. MI or CCA Whether the primary analysis was justified on the basis of missingness assumptions If applicable, any other analyses conducted to answer the key causal question that handle the missing data differently (e.g. a CCA or a delta-adjusted MI analysis, where imputations are shifted by a parameter “delta” representing the difference between the observed and unobserved data(24)) Whether the alternative analysis was justified If a delta-adjusted MI analysis was used, whether external information elicited from subject-matter experts was used to choose the value(s) of the delta parameter
MI implementation	<ul style="list-style-type: none"> The method used for MI, for example, multivariate normal imputation or multiple imputation by chained equations The statistical software used for MI The number of imputations performed, Whether all analysis variables were included in the imputation model Whether auxiliary variables (i.e. variables defined as potential predictors of missingness and/or the variable(s) with missing data, but are not included in the target analysis) were included in the imputation model Whether interactions were included in the imputation model

Analysis

The questionnaire data will be cleaned and analysed in R. Descriptive statistics will be used to summarise the data. Frequencies and percentages will be presented for categorical data, for example, the method used to obtain the primary results. Median and interquartile range will be presented for continuous data, for example, the percentage of complete cases in each observational study. All data and code will be made publicly available on GitHub.

Reporting

Findings from this review will be reported using the Preferred Reporting Items for Systematic reviews and Meta-Analyses extension for Scoping Reviews (PRISMA-ScR) checklist.(25)

Patient and public involvement

There will be no patient or public involvement in this project because data will be collected only from published studies.

DISCUSSION

Previous reviews of the handling of missing data have primarily focused on RCTs with incomplete outcomes. Observational studies are subject to greater challenges than RCTs in terms of missing data as they often face missing data in multiple variables (exposure, outcome and/or confounders). This paper describes a protocol for a scoping review of how MI is used to handle missing data in observational studies that answer causal questions.

Strengths and limitations

There are several strengths to our study. A targeted review of observational studies in top epidemiology journals publishing general research will benchmark the current state of practice for handling multivariable

1
2
3 missingness with MI. Screening, reviewing and data extraction will be performed systematically. All data and
4 code will be made publicly available, enabling our analysis to be entirely reproducible. Results from the review
5 will be reported according to best practice, using PRISMA-ScR.
6

7 There are also limitations. Identifying whether the aim of the research was to answer a descriptive, causal or
8 predictive question is somewhat subjective because many researchers have not adopted this classification of
9 research questions.(1) Although our targeted review will not include studies from all epidemiology journals,
10 we expect that included studies (expected to be > 100 studies from five major epidemiology journals) will be
11 sufficient to provide insight and general trends on the methods of interest. It is likely that some of the
12 information sought will be unclear or not reported. To accommodate this, we have specified how anticipated
13 challenges with data extraction will be handled if they arise.
14

15 **Implications of this research**

16
17 In addition to critically appraising the current state of the literature regarding the use and reporting of
18 analyses using MI to handle missing data, this review will identify areas for improvement in the handling and
19 reporting of missing data in observational studies. The results of this review will be used to develop practical
20 guidance for researchers and promote the formulation of missingness assumptions in a clear and transparent
21 manner.
22

23 **Funding sources / sponsors**

24
25 This work was supported by an Australian National Health and Medical Research Council (NHMRC) Career
26 Development Fellowship (CDF) Level 2 Grant (grant 1127984 awarded to KJL), a NHMRC Investigator Grant
27 Leadership Level 1 (grant 1196068 awarded to JAS), a NHMRC Investigator Grant Emerging Leadership Level 2
28 (grant 2009572 awarded to MMB) and a NHMRC Project Grant (grant 1166023). Research at the Murdoch
29 Children's Research Institute is supported by the Victorian Government's Operational Infrastructure Support
30 Program.
31

32 **Authors' contributions**

33
34 RM conceived the study and wrote the first draft of the manuscript. All authors contributed to the design of
35 the study, revision of the manuscript and take public responsibility for its content.
36

37 **Competing interests statement**

38
39 None declared.
40

41 **REFERENCES**

- 42 1. Hernán MA, Hsu J, Healy B. A second chance to get causal inference right: a classification of data
43 science tasks. *Chance*. 2019;32(1):42-9.
- 44 2. Little RJ, Rubin DB. *Statistical analysis with missing data*: John Wiley & Sons; 2019.
- 45 3. Rubin DB. *Multiple imputation for nonresponse in surveys*: John Wiley & Sons; 2004.
- 46 4. Van Buuren S. *Flexible imputation of missing data*: CRC press; 2018.
- 47 5. Sterne JAC, White IR, Carlin JB, Spratt M, Royston P, Kenward MG, et al. Multiple imputation for
48 missing data in epidemiological and clinical research: potential and pitfalls. *BMJ*. 2009;338:b2393.
- 49 6. Seaman S, Galati J, Jackson D, Carlin J. What is meant by "missing at random"? *Statistical Science*.
50 2013;28(2):257-68.
- 51 7. Mohan K, Pearl J. Graphical models for processing missing data. *Journal of the American Statistical*
52 *Association*. 2021;116(534):1023-37.
- 53 8. National Research Council. *The prevention and treatment of missing data in clinical trials*. 2010.
- 54 9. Tan P-T, Cro S, Van Vogt E, Szigeti M, Cornelius VR. A review of the use of controlled multiple
55 imputation in randomised controlled trials with missing outcome data. *BMC medical research methodology*.
56 2021;21(1):1-17.
- 57 10. Rabe BA, Day S, Fiero MH, Bell ML. Missing data handling in non-inferiority and equivalence trials: A
58 systematic review. *Pharmaceutical statistics*. 2018;17(5):477-88.
- 59 11. Fiero MH, Huang S, Oren E, Bell ML. Statistical analysis and handling of missing data in cluster
60 randomized trials: a systematic review. *Trials*. 2016;17(1):1-10.

12. Bell ML, Fiero M, Horton NJ, Hsu C-H. Handling missing data in RCTs; a review of the top medical journals. *BMC medical research methodology*. 2014;14(1):1-8.
13. Powney M, Williamson P, Kirkham J, Kolamunnage-Dona R. A review of the handling of missing longitudinal outcome data in clinical trials. *Trials*. 2014;15(1):1-11.
14. Ibrahim F, Tom BD, Scott DL, Prevost AT. A systematic review of randomised controlled trials in rheumatoid arthritis: the reporting and handling of missing data in composite outcomes. *Trials*. 2016;17(1):1-8.
15. Rombach I, Rivero-Arias O, Gray AM, Jenkinson C, Burke O. The current practice of handling and reporting missing outcome data in eight widely used PROMs in RCT publications: a review of the current literature. *Quality of Life Research*. 2016;25(7):1613-23.
16. White IR, Horton NJ, Carpenter J, statistics rim, social, Pocock SJ. Strategy for intention to treat analysis in randomised trials with missing outcome data. *BMJ*. 2011;342:d40.
17. Mackinnon A. The use and reporting of multiple imputation in medical research—a review. *Journal of internal medicine*. 2010;268(6):586-93.
18. Hayati Rezvan P, Lee KJ, Simpson JA. The rise of multiple imputation: a review of the reporting and implementation of the method in medical research. *BMC medical research methodology*. 2015;15(1):1-14.
19. Rehal S, Morris TP, Fielding K, Carpenter JR, Phillips PP. Non-inferiority trials: are they inferior? A systematic review of reporting in major medical journals. *BMJ open*. 2016;6(10):e012594.
20. Moreno-Betancur M, Lee KJ, Leacy FP, White IR, Simpson JA, Carlin JB. Canonical Causal Diagrams to Guide the Treatment of Missing Data in Epidemiologic Studies. *American Journal of Epidemiology*. 2018;187(12):2705-15.
21. Hernán MA. The C-word: scientific euphemisms do not improve causal inference from observational data. *American journal of public health*. 2018;108(5):616-9.
22. de Vries BBP, van Smeden M, Rosendaal FR, Groenwold RH. , abstract, and keyword searching resulted in poor recovery of articles in systematic reviews of epidemiologic practice. *Journal of Clinical Epidemiology*. 2020;121:55-61.
23. Veritas Health Innovation. Covidence systematic review software. Melbourne, Australia.
24. Cro S, Morris TP, Kenward MG, Carpenter JR. Sensitivity analysis for clinical trials with missing continuous outcome data using controlled multiple imputation: a practical guide. *Statistics in medicine*. 2020;39(21):2815-42.
25. Tricco AC, Lillie E, Zarin W, O'Brien KK, Colquhoun H, Levac D, et al. PRISMA extension for scoping reviews (PRISMA-ScR): checklist and explanation. *Annals of internal medicine*. 2018;169(7):467-73.

The handling of missing data with multiple imputation in observational studies that address causal questions: Protocol for a scoping review

Rheanna M. Mainzer^{*1,2}, Margarita Moreno-Betancur^{1,2}, Cattram D. Nguyen^{1,2}, Julie A. Simpson³, John B. Carlin^{1,2,3}, Katherine J. Lee^{1,2}

1. Clinical Epidemiology and Biostatistics Unit, Murdoch Children's Research Institute, Parkville, Victoria 3052, Australia
2. Department of Paediatrics, The University of Melbourne, Parkville, Victoria 3052, Australia
3. Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, The University of Melbourne, Parkville, Victoria 3052, Australia

*Corresponding author: Rheanna Mainzer; rheanna.mainzer@mcri.edu.au

Supplementary Table 1. Anticipated challenges with data extraction and how they will be handled.

Challenge for data extraction	Category of items affected	How challenge will be handled
Articles may have more than one publication date, for example, the date the article first appeared online and when it was published in-print.	Inclusion criteria	Only one publication date is required to be between January 2019 and December 2021. If two or more publication dates are between January 2019 and December 2021, the earlier date will be recorded.
There are multiple causal questions, exposures or outcomes.	Missing data	We will identify the primary causal question based on the research aims and conclusion. The proportion of missing data in the exposure, outcome and confounders used to answer this primary question will be recorded. This is expected to be acceptable in most cases. If the primary causal question cannot be identified due to multiple outcomes, we will report the missing data details for the first outcome listed in the methods section. (This is comparable to the strategy taken by Fiero et al. (1)) Similarly, if the primary causal question cannot be identified due to multiple exposures, we will report the missing data details for the first exposure listed in the methods section.
Multiple sets of covariates are used for adjustment.	Missing data	The largest adjustment set will be considered. The number of incomplete covariates will be recorded categorically (no incomplete covariates, 1 incomplete covariate, 2 or more incomplete covariates, not stated or unable to establish). This categorisation has been chosen to enable determination of multivariable missingness.
Not clear whether all variables in the target analysis were	MI implementation	If some (but not all) analysis variables were reported as being included in the imputation model then we will assume

1 2 3 4 5 6 7 8	included in the imputation model.		that the analysis variables not explicitly mentioned were excluded from the imputation model. If there was no description of the imputation model, then we will categorise this as “unclear”.
9 10 11 12 13 14	Not clear whether auxiliary variables or interactions were included in the imputation model.	MI implementation	If it is not explicitly stated that these were included in the imputation model, we will assume they were excluded. If there was no mention of the imputation model then we will categorise this as “unclear”.
15 16 17 18 19 20 21 22 23 24	Imputation method used not explicitly stated.	MI implementation	If the imputation method used (e.g. multivariate normal imputation or multiple imputation by chained equations) is not provided, we will infer the method used, where possible, from the statistical software procedures listed in the main paper or supplementary material. If the method is unable to be inferred, we will categorise this as “unclear”.

REFERENCE

1. Fiero MH, Huang S, Oren E, Bell ML. Statistical analysis and handling of missing data in cluster randomized trials: a systematic review. *Trials*. 2016;17(1):1-10.

BMJ Open

The handling of missing data with multiple imputation in observational studies that address causal questions: Protocol for a scoping review

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2022-065576.R1
Article Type:	Protocol
Date Submitted by the Author:	15-Nov-2022
Complete List of Authors:	Mainzer, Rheanna; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Moreno-Betancur , Margarita ; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Nguyen, Cattram; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Simpson, Julie; University of Melbourne School of Population and Global Health Carlin, John; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Lee, Katherine; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics
Primary Subject Heading:	Epidemiology
Secondary Subject Heading:	Public health, Research methods
Keywords:	EPIDEMIOLOGY, STATISTICS & RESEARCH METHODS, Public health < INFECTIOUS DISEASES

SCHOLARONE™
Manuscripts

The handling of missing data with multiple imputation in observational studies that address causal questions: Protocol for a scoping review

Rheanna M. Mainzer*^{1,2}, Margarita Moreno-Betancur^{1,2}, Cattram D. Nguyen^{1,2}, Julie A. Simpson³, John B. Carlin^{1,2,3}, Katherine J. Lee^{1,2}

1. Clinical Epidemiology and Biostatistics Unit, Murdoch Children's Research Institute, Parkville, Victoria 3052, Australia
2. Department of Paediatrics, The University of Melbourne, Parkville, Victoria 3052, Australia
3. Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, The University of Melbourne, Parkville, Victoria 3052, Australia

*Corresponding author: Rheanna Mainzer; rheanna.mainzer@mcri.edu.au

ABSTRACT

Introduction Observational studies in health-related research often aim to answer causal questions. Missing data are common in these studies and often occur in multiple variables, such as the exposure, outcome and/or variables used to control for confounding. The standard classification of missing data as missing completely at random (MCAR), missing at random (MAR) or missing not at random (MNAR), does not allow for a clear assessment of missingness assumptions when missingness arises in more than one variable. This presents challenges for selecting an analytic approach and determining when a sensitivity analysis under plausible alternative missing data assumptions is required. This is particularly pertinent with multiple imputation (MI), which is often justified by assuming data are MAR. The objective of this scoping review is to examine the use of MI in observational studies that address causal questions, with a focus on if and how (i) missingness assumptions are expressed and assessed, (ii) missingness assumptions are used to justify the choice of a complete case analysis and/or MI for handling missing data, and (iii) sensitivity analyses under alternative plausible assumptions about the missingness mechanism are conducted.

Methods and analysis We will systematically review observational studies that aim to answer causal questions and use MI, published between January 2019 and December 2021 in five top general epidemiology journals. Studies will be identified using a full text search for the term "multiple imputation" and then assessed for eligibility. Information extracted will include details about the study characteristics, missing data, missingness assumptions and MI implementation. Data will be summarised using descriptive statistics.

Ethics and dissemination Ethics approval is not required for this review because data will be collected only from published studies. The results will be disseminated through a peer reviewed publication and conference presentations.

Registration This protocol is registered on figshare (<https://doi.org/10.6084/m9.figshare.20010497.v1>).

Strengths and limitations of this study

- A targeted review of observational studies published in the five top-ranked epidemiology journals will benchmark the current state of practice for handling multivariable missingness with multiple imputation in causal analyses.
- Screening, reviewing and data extraction will be performed systematically, with double data extraction for a subset of articles and any discrepancies resolved by a panel.
- It is likely that some of the information sought will be ambiguously reported or not reported.
- Potential challenges with data extraction have been considered and a strategy for handling these challenges has been put in place.
- All extracted data and code will be made publicly available, enabling our descriptive analysis to be entirely reproducible.

1 INTRODUCTION

2 Observational studies in clinical and health-related research often aim to answer causal questions, even if this
3 intent is only implicit.(1, 2) This aim is usually addressed by estimation of a target parameter to quantify the
4 impact of intervening on an exposure on an outcome of interest, in a given population. In observational
5 studies missing data are common and can occur in multiple variables, such as the exposure, the outcome
6 and/or the variables used to control for confounding. Restricting statistical analysis to individuals with
7 complete data on all analysis variables, i.e., conducting a “complete case analysis” (CCA), can lead to bias
8 and/or loss of precision in estimates of the target parameter.(3) Multiple imputation (MI) is a popular and
9 flexible approach for estimating a target parameter in the presence of incomplete data.(4, 5) In the first stage
10 of MI, missing data are imputed multiple times with random draws from the predictive distribution of the
11 missing values given the observed data and a specified imputation model. In the second stage, the statistical
12 analysis of interest is applied to each imputed dataset and the results are combined using Rubin’s rules to
13 obtain a single estimate of the target parameter with associated standard error.(4)

14 Standard implementations of MI are known to provide consistent estimation of target parameters under
15 certain (unverifiable) assumptions about the mechanism leading to missing data. Assumptions about missing
16 data are usually expressed using Rubin’s classification of missing data mechanisms into missing completely at
17 random (MCAR, where the probability of data being missing does not depend on the observed or unobserved
18 data), missing at random (MAR, where the probability of data being missing does not depend on the
19 unobserved data, conditional on the observed data) and missing not at random (MNAR, where the probability
20 of data being missing depends on the unobserved data, even after conditioning on the observed data).(6)
21 While this framework is useful if missing data occur in a single variable, it raises issues when missingness arises
22 in more than one variable. First, what these mechanisms mean with multivariable missingness is poorly
23 understood and does not allow for a transparent assessment of missingness assumptions.(7) Second, based on
24 our experience researching, teaching and applying MI, these mechanisms have become widely
25 (mis)understood as synonymous with methods. For example, researchers often use MI under the assumption
26 that data are MAR, but this is only a sufficient and not necessary condition for standard MI to be consistent.(8)
27 Both a CCA and a MI analysis could be unbiased under a range of multivariable missingness mechanisms (even
28 those considered to be MNAR).(9) Likewise, there are missingness mechanisms in which neither MI nor a CCA
29 can be used to estimate an exposure-outcome association without bias, and a different approach would be
30 needed for unbiased estimation.

31 Because one cannot verify from the observed data what the true missing data mechanism is, sensitivity
32 analyses to examine the robustness of results to alternative plausible assumptions about the missingness
33 mechanism (hereafter, “sensitivity analyses”) are strongly recommended.(10) However, as stated by the US
34 National Research Council, “the usefulness of a sensitivity analysis ultimately depends on the transparency and
35 plausibility of the unverifiable assumptions.”(10) The inherent difficulty in assessing missingness assumptions
36 when framed in the traditional MCAR/MAR/MNAR manner is an obvious obstacle to this. Furthermore, the
37 mistakenly assumed one-to-one correspondence with analytic approaches in the presence of multivariable
38 missingness leads to misguided practices. For example, from our observation, MI is routinely applied as a
39 sensitivity analysis to a CCA. However, the logic behind applying MI as a sensitivity analysis to a CCA (or vice
40 versa) without first considering one’s assumptions about the missingness mechanism is unsound.(11)
41 Obtaining similar or different estimates from these analyses does not provide insight into the impact of
42 alternative plausible assumptions about the missingness mechanism on the study results unless one has first
43 made their missingness assumptions explicit and identified these two approaches as appropriate for
44 estimating the target parameter under those explicit assumptions.

45 Most reviews of the handling and reporting of missing data, and the implementation and documentation of
46 MI, have been carried out in the context of randomised controlled trials (RCTs).(12-18) For trials, typically only
47 the outcome variable is incomplete, while the intervention and other key variables (typically baseline
48 variables) are observed for all participants. In this setting where there are missing data in a single variable, the
49 MCAR/MAR/MNAR framework is more transparent and guidance on sensitivity analyses has been well-
50 developed (see, for example, (15, 19)). In contrast, there have been few reviews concerned with how missing

1 data are handled in observational studies where there is the additional complication of multivariable
2 missingness. A review by Mackinnon published in 2010 found that only two (4%) out of 50 non-RCT studies
3 reviewed carried out an additional analysis that was described as a sensitivity analysis.(11) Similarly, Rezvan et
4 al. (2015) found that none of the 30 observational studies reviewed conducted a sensitivity analysis to
5 departures from the missingness assumptions following MI.(20)

6 While the reviews by Mackinnon and Rezvan et al. provide useful insight into the problem, neither focused
7 specifically on observational studies and the issues described above. In addition, subsequent to publication of
8 these reviews there have been important developments in the theory and application of missingness directed
9 acyclic graphs (m-DAGs), also known as m-graphs, a tool for the formulation of causal assumptions in the
10 presence of multivariable missingness.(8) M-DAGs aid the depiction and assessment of missingness
11 assumptions. Clarity regarding each plausible causal mechanism underlying the missing data then facilitates
12 the choice of analytical approach. For example, the application of DAG theory allows one to determine
13 whether a target parameter can be estimated without bias from the available data using an approach like CCA
14 or MI, or whether additional assumptions and a more sophisticated analysis is required (such as a delta-
15 adjusted MI approach, where imputations are shifted by a parameter “delta” representing the difference
16 between the observed and unobserved data).(9, 21-23)

17 The aim of this scoping review is to systematically review the epidemiological literature to examine the use of
18 MI in observational studies that address causal questions, which is typically the focus of such studies even
19 when this may not be very clearly articulated.(2) These studies often face missingness in multiple variables
20 required for analysis. We will examine (i) how missingness assumptions are expressed, (ii) if and how
21 missingness assumptions are used to justify the choice of a CCA and/or MI for handling missing data, and (iii)
22 the conduct of sensitivity analyses to alternative plausible assumptions about the missingness mechanism. We
23 will also examine how MI is implemented. This review will be used to document the current state of practice,
24 to identify areas for improvement in the handling and reporting of missing data with MI in observational
25 studies, and to subsequently develop guidance on these key components for researchers.

26 **METHODS AND ANALYSIS**

27 In this section we provide a full description of the study design, including how articles will be selected, which
28 variables will be extracted, and how data will be analysed. The review described in this protocol began in June
29 2022 and we anticipate it will be completed by June 2023.

30 **Search strategy**

31 We will systematically search five general epidemiology journals for observational studies published between
32 January 2019 and December 2021 that aim to answer at least one causal research question using MI. The
33 general epidemiology journals that will be included in this search are: *International Journal of Epidemiology*,
34 *American Journal of Epidemiology*, *European Journal of Epidemiology*, *Journal of Clinical Epidemiology* and
35 *Epidemiology*. These journals were chosen because they are high ranking, general journals in epidemiology
36 that publish original research from observational studies. As such, articles from these journals should capture
37 the current best practice in the use of MI to handle missing data when answering causal questions using
38 observational data. They have also been used previously in a systematic review of epidemiologic practice.(24)
39 Original research articles will be identified using the full-text search term “multiple imputation” on each
40 journal’s website. This search strategy is similar to that used in previous scoping reviews in this area.(11, 20)

41 **Inclusion criteria**

42 We will include original research articles published between January 2019 and December 2021 that aim to
43 answer at least one causal question using MI to handle missing data. We will determine that a study has aimed
44 to answer a causal question if at least one of the following criteria is satisfied:

- 45 1. the authors explicitly stated they were estimating a causal effect;
- 46 2. the study estimated an effect that was given (at least implicitly) a causal interpretation, i.e., an
47 interpretation which suggested that intervening on the exposure could change the outcome (e.g.,
48 increasing coffee consumption may be protective against stroke). This will be determined by wording

in conclusions. If it is not clear from this wording alone, investigation of the following three typical signals of causal analyses will be used to aid in the determining: identification of confounders, the inclusion of a DAG to illustrate causal assumption made in the analysis, and analytical approaches incorporating adjustment for confounders (for example, estimating an effect using a regression model that was adjusted for a set of covariates).

Studies on all disease areas/medical conditions and any target population will be considered.

Exclusion criteria

Studies will be excluded from the review if they meet any of the following criteria:

- *No causal question.* The article did not aim to answer a causal question, for example, the aim of the study was to develop a predictive model or to estimate a disease burden.
- *Unclear type of question.* A clear research goal could not be identified. In other words, it was unclear whether the study aimed to answer a descriptive, predictive or causal question.
- *The analysis did not use MI.*
- *Methodological research.* The primary purpose of the article was methodological development, for example, using a simulation study to compare the performance of methods or mathematical derivations to develop a new method or model. While these articles often include comprehensive case studies, they may not be representative of empirical studies aiming primarily to answer causal research questions.
- *Aggregate-level data.* The analysis was based on aggregated data where MI could not be applied at the participant level, as is common in meta-analysis or interrupted time series analysis.
- *Qualitative research.* The article provided a commentary, review, opinion, study protocol, study profile or description only.
- *Trial.* The study intervention was assigned to participants by the study investigators.

Sample size

We will require at least 100 studies to estimate the percentage of studies with a particular element (e.g., studies that justify their missingness assumptions) to within a maximum margin of error (two standard errors) of 10%. Assuming a prevalence of 50%, this would give a 95% confidence interval from 40% to 60%. For a prevalence greater than or less than 50%, the 95% confidence interval will be narrower. This sample size is similar to the sample size used in the first review of MI in medical research ($n = 99$, (11)), and many of the subsequent reviews in this area (e.g., $n = 103$ in (20), 77 in (15) and 118 in (12)). We expect to identify at least 100 eligible studies given the three-year publication time frame. All eligible studies will be included in the review.

Study selection

The search of the journal databases and selection of studies for inclusion in the review will be performed primarily by a single researcher (RM) in two steps. First, the title, abstract and date of each article will be screened to rule out studies that are clearly not eligible for the review. Second, the full text of the remaining studies will be reviewed to confirm if studies are eligible for the review. If a decision about the eligibility of an article cannot be reached by RM (for example, due to uncertainty about the inclusion criteria), a second researcher (CN) will independently review the full text. Disagreements about inclusion criteria will be resolved by discussion in meetings with at least three researchers (RM, CN and at least one of JC, JS, KL or MMB).

Data extraction and management

Covidence, a web-based tool for systematic review management, will be used to perform the review. (25) The data extraction questionnaire was developed and tested for use by RM and KL using a sample of 10 articles. Data from all eligible studies will be extracted by RM. The supplementary material of all eligible studies will also be reviewed. We will use double data extraction (performed by KL) for a random selection of 10% of articles and additionally when there is uncertainty about the information being extracted. Discrepancies and

1
2
3 1 uncertainties will be resolved by discussion in meetings with at least three researchers (RM, KL and at least
4 2 one of JC, JS, CN or MMB).

3 **Outcomes measured**

4 We will extract data pertaining to the study characteristics, the amount of missing data and in which variables
5 it occurs, missingness assumptions, methods for handling missing data and implementation of multiple
6 imputation. Data extraction items are summarised in Table 1. Because we anticipate difficulties in extracting
7 some items (such as the percentage of complete cases), in Supplementary Table 1 we list potential challenges
8 in extracting data and any assumptions or simplifications that will be made if these challenges arise. Any post-
9 hoc assumptions or simplifications for unanticipated challenges will be recorded and reported as part of the
10 analysis.

11 **Table 1.** Summary of items to be extracted from each article.

Category	Summary of data extraction items
Study characteristics	<ul style="list-style-type: none"> • First author's last name • Publication date • Journal • Type of study design
Missing data	<ul style="list-style-type: none"> • Percentage of complete cases • Percentage of missing values in the exposure and outcome • Number of incomplete covariates
Missingness assumptions	<ul style="list-style-type: none"> • Statement of missingness data assumptions (including whether the study used m-DAGs or the MCAR/MAR/MNAR framework) • Justification of missingness assumptions
Analysis methods	<ul style="list-style-type: none"> • The primary analysis method used to answer the key causal question, e.g. MI or CCA • Whether the primary analysis was justified on the basis of missingness assumptions • If applicable, any other analyses conducted to answer the key causal question that handle the missing data differently (e.g. a CCA or a delta-adjusted MI analysis) • Whether the alternative analysis was justified on the basis of missingness assumptions • If a delta-adjusted MI analysis was used, whether external information elicited from subject-matter experts was used to choose the value(s) of the delta parameter
MI implementation	<ul style="list-style-type: none"> • The method used for MI, for example, multivariate normal imputation or multiple imputation by chained equations • The statistical software used for MI • The number of imputations performed • Whether all analysis variables were included in the imputation model • Whether auxiliary variables (i.e. variables defined as potential predictors of the variable(s) with missing data that are not included in the target analysis) were included in the imputation model • Whether interactions were included in the imputation model

12

13 **Analysis**

14 The questionnaire data will be cleaned and analysed in R. Descriptive statistics will be used to summarise the
15 data. Frequencies and percentages will be presented for categorical data, for example, the method used to
16 obtain the primary results. Median and interquartile range will be presented for continuous data, for example,
17 the percentage of complete cases in each observational study. All data and code will be made publicly
18 available on GitHub.

19

20

21

22

23

24

25

26

1 Reporting

2 Findings from this review will be reported using the Preferred Reporting Items for Systematic reviews and
3 Meta-Analyses extension for Scoping Reviews (PRISMA-ScR) checklist.(26)

4 Patient and public involvement

5 There will be no patient or public involvement in this project because data will be collected only from
6 published studies.

7 ETHICS AND DISSEMINATION

8 Ethics approval is not required for this review because data will be collected only from published studies. The
9 results will be disseminated through a peer-review publication and conference presentations.

10 DISCUSSION

11 Previous reviews of the handling of missing data have primarily focused on RCTs with incomplete outcome
12 data. Observational studies that answer causal questions are common and subject to greater challenges than
13 RCTs in terms of missing data as they often face missing data in multiple variables (exposure, outcome and/or
14 confounders). This paper describes a protocol for a scoping review of how MI is used to handle missing data in
15 these studies.

16 Strengths and limitations

17 There are several strengths to our study. A targeted review of observational studies in top epidemiology
18 journals publishing general research will benchmark the current state of practice for handling multivariable
19 missingness with MI in causal analyses. Screening, reviewing and data extraction will be performed
20 systematically. All data and code will be made publicly available, enabling our analysis to be entirely
21 reproducible. Results from the review will be reported according to best practice, using PRISMA-ScR.

22 There are also limitations. Identifying whether the aim of the research was to answer a descriptive, causal or
23 predictive question is somewhat subjective because many researchers have not adopted this classification of
24 research questions.(1) Although our targeted review will not include studies from all epidemiology journals,
25 we expect that included studies (expected to be > 100 studies from five major epidemiology journals) will be
26 sufficient to provide insight and general trends on the methods of interest. It is likely that some of the
27 information sought will be unclear or not reported. To accommodate this, we have specified how anticipated
28 challenges with data extraction will be handled if they arise.

29 Implications of this research

30 In addition to critically appraising the current state of the literature regarding the use and reporting of causal
31 analyses using MI to handle missing data in observational studies, this review will identify areas for
32 improvement in the handling and reporting of missing data in these studies. The results of this review will be
33 used to develop practical guidance for researchers and inform future research in these areas.

34 Funding sources / sponsors

35 This work was supported by an Australian National Health and Medical Research Council (NHMRC) Career
36 Development Fellowship (CDF) Level 2 Grant (grant 1127984 awarded to KJL), a NHMRC Investigator Grant
37 Leadership Level 1 (grant 1196068 awarded to JAS), a NHMRC Investigator Grant Emerging Leadership Level 2
38 (grant 2009572 awarded to MMB) and a NHMRC Project Grant (grant 1166023). Research at the Murdoch
39 Children's Research Institute is supported by the Victorian Government's Operational Infrastructure Support
40 Program.

41 Authors' contributions

42 RM conceived the study idea, developed the methodology, designed the data extraction tool, drafted and
43 revised the paper. KL developed the study idea, methodology, data extraction tool and revised the paper.
44 MMB and JS developed the study idea, methodology and revised the paper. CN developed the study idea,

1 methodology and data extraction tool. JC developed the study idea, methodology, data extraction tool and
2 revised the paper.

3 **Competing interests statement**

4 None declared.

5 **REFERENCES**

- 6 1. Hernán MA, Hsu J, Healy B. A second chance to get causal inference right: a classification of data
7 science tasks. *Chance*. 2019;32(1):42-9.
- 8 2. Hernán MA. The C-word: scientific euphemisms do not improve causal inference from observational
9 data. *American journal of public health*. 2018;108(5):616-9.
- 10 3. Little RJ, Rubin DB. *Statistical analysis with missing data*: John Wiley & Sons; 2019.
- 11 4. Rubin DB. *Multiple imputation for nonresponse in surveys*: John Wiley & Sons; 2004.
- 12 5. Van Buuren S. *Flexible imputation of missing data*: CRC press; 2018.
- 13 6. Sterne JAC, White IR, Carlin JB, Spratt M, Royston P, Kenward MG, et al. Multiple imputation for
14 missing data in epidemiological and clinical research: potential and pitfalls. *BMJ*. 2009;338:b2393.
- 15 7. Seaman S, Galati J, Jackson D, Carlin J. What is meant by “missing at random”? *Statistical Science*.
16 2013;28(2):257-68.
- 17 8. Mohan K, Pearl J. Graphical models for processing missing data. *Journal of the American Statistical*
18 *Association*. 2021;116(534):1023-37.
- 19 9. Moreno-Betancur M, Lee KJ, Leacy FP, White IR, Simpson JA, Carlin JB. Canonical Causal Diagrams to
20 Guide the Treatment of Missing Data in Epidemiologic Studies. *American Journal of Epidemiology*.
21 2018;187(12):2705-15.
- 22 10. National Research Council. *The prevention and treatment of missing data in clinical trials*. 2010.
- 23 11. Mackinnon A. The use and reporting of multiple imputation in medical research—a review. *Journal of*
24 *internal medicine*. 2010;268(6):586-93.
- 25 12. Tan P-T, Cro S, Van Vogt E, Szigeti M, Cornelius VR. A review of the use of controlled multiple
26 imputation in randomised controlled trials with missing outcome data. *BMC medical research methodology*.
27 2021;21(1):1-17.
- 28 13. Rabe BA, Day S, Fiero MH, Bell ML. Missing data handling in non-inferiority and equivalence trials: A
29 systematic review. *Pharmaceutical statistics*. 2018;17(5):477-88.
- 30 14. Fiero MH, Huang S, Oren E, Bell ML. Statistical analysis and handling of missing data in cluster
31 randomized trials: a systematic review. *Trials*. 2016;17(1):1-10.
- 32 15. Bell ML, Fiero M, Horton NJ, Hsu C-H. Handling missing data in RCTs; a review of the top medical
33 journals. *BMC medical research methodology*. 2014;14(1):1-8.
- 34 16. Powney M, Williamson P, Kirkham J, Kolamunnage-Dona R. A review of the handling of missing
35 longitudinal outcome data in clinical trials. *Trials*. 2014;15(1):1-11.
- 36 17. Ibrahim F, Tom BD, Scott DL, Prevost AT. A systematic review of randomised controlled trials in
37 rheumatoid arthritis: the reporting and handling of missing data in composite outcomes. *Trials*. 2016;17(1):1-8.
- 38 18. Rombach I, Rivero-Arias O, Gray AM, Jenkinson C, Burke O. The current practice of handling and
39 reporting missing outcome data in eight widely used PROMs in RCT publications: a review of the current
40 literature. *Quality of Life Research*. 2016;25(7):1613-23.
- 41 19. White IR, Horton NJ, Carpenter J, statistics rim, social, Pocock SJ. Strategy for intention to treat
42 analysis in randomised trials with missing outcome data. *BMJ*. 2011;342:d40.
- 43 20. Hayati Rezvan P, Lee KJ, Simpson JA. The rise of multiple imputation: a review of the reporting and
44 implementation of the method in medical research. *BMC medical research methodology*. 2015;15(1):1-14.
- 45 21. Tompsett DM, Leacy F, Moreno-Betancur M, Heron J, White IR. On the use of the not-at-random fully
46 conditional specification (NARFCS) procedure in practice. *Statistics in medicine*. 2018;37(15):2338-53.
- 47 22. Cro S, Morris TP, Kenward MG, Carpenter JR. Sensitivity analysis for clinical trials with missing
48 continuous outcome data using controlled multiple imputation: a practical guide. *Statistics in medicine*.
49 2020;39(21):2815-42.
- 50 23. Rezvan PH, Lee KJ, Simpson JA. Sensitivity analysis within multiple imputation framework using delta-
51 adjustment: Application to Longitudinal Study of Australian Children. *Longitudinal and Life Course Studies*.
52 2018;9(3):259-78.

- 1
2
3 1 24. de Vries BBP, van Smeden M, Rosendaal FR, Groenwold RH. Title, abstract, and keyword searching
4 2 resulted in poor recovery of articles in systematic reviews of epidemiologic practice. Journal of Clinical
5 3 Epidemiology. 2020;121:55-61.
6 4 25. Veritas Health Innovation. Covidence systematic review software. Melbourne, Australia.
7 5 26. Tricco AC, Lillie E, Zarin W, O'Brien KK, Colquhoun H, Levac D, et al. PRISMA extension for scoping
8 6 reviews (PRISMA-ScR): checklist and explanation. Annals of internal medicine. 2018;169(7):467-73.
9
10 7
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

The handling of missing data with multiple imputation in observational studies that address causal questions: Protocol for a scoping review

Rheanna M. Mainzer^{*1,2}, Margarita Moreno-Betancur^{1,2}, Cattram D. Nguyen^{1,2}, Julie A. Simpson³, John B. Carlin^{1,2,3}, Katherine J. Lee^{1,2}

1. Clinical Epidemiology and Biostatistics Unit, Murdoch Children's Research Institute, Parkville, Victoria 3052, Australia
2. Department of Paediatrics, The University of Melbourne, Parkville, Victoria 3052, Australia
3. Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, The University of Melbourne, Parkville, Victoria 3052, Australia

*Corresponding author: Rheanna Mainzer; rheanna.mainzer@mcri.edu.au

Supplementary Table 1. Anticipated challenges with data extraction and how they will be handled.

Challenge for data extraction	Category of items affected	How challenge will be handled
Articles may have more than one publication date, for example, the date the article first appeared online and when it was published in-print.	Inclusion criteria	Only one publication date is required to be between January 2019 and December 2021. If two or more publication dates are between January 2019 and December 2021, the earlier date will be recorded.
There are multiple causal questions, exposures or outcomes.	Missing data	We will identify the primary causal question based on the research aims and conclusion. The proportion of missing data in the exposure, outcome and confounders used to answer this primary question will be recorded. This is expected to be acceptable in most cases. If the primary causal question cannot be identified due to multiple outcomes, we will report the missing data details for the first outcome listed in the methods section. (This is comparable to the strategy taken by Fiero et al. (1)) Similarly, if the primary causal question cannot be identified due to multiple exposures, we will report the missing data details for the first exposure listed in the methods section.
Multiple sets of covariates are used for adjustment.	Missing data	The largest adjustment set will be considered. The number of incomplete covariates will be recorded categorically (no incomplete covariates, 1 incomplete covariate, 2 or more incomplete covariates, not stated or unable to establish). This categorisation has been chosen to enable determination of multivariable missingness.
Not clear whether all variables in the target analysis were	MI implementation	If some (but not all) analysis variables were reported as being included in the imputation model then we will assume

1 2 3 4 5 6 7 8	included in the imputation model.		that the analysis variables not explicitly mentioned were excluded from the imputation model. If there was no description of the imputation model, then we will categorise this as “unclear”.
9 10 11 12 13 14	Not clear whether auxiliary variables or interactions were included in the imputation model.	MI implementation	If it is not explicitly stated that these were included in the imputation model, we will assume they were excluded. If there was no mention of the imputation model then we will categorise this as “unclear”.
15 16 17 18 19 20 21 22 23 24	Imputation method used not explicitly stated.	MI implementation	If the imputation method used (e.g. multivariate normal imputation or multiple imputation by chained equations) is not provided, we will infer the method used, where possible, from the statistical software procedures listed in the main paper or supplementary material. If the method is unable to be inferred, we will categorise this as “unclear”.

REFERENCE

1. Fiero MH, Huang S, Oren E, Bell ML. Statistical analysis and handling of missing data in cluster randomized trials: a systematic review. *Trials*. 2016;17(1):1-10.

BMJ Open

The handling of missing data with multiple imputation in observational studies that address causal questions: protocol for a scoping review

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2022-065576.R2
Article Type:	Protocol
Date Submitted by the Author:	03-Jan-2023
Complete List of Authors:	Mainzer, Rheanna; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Moreno-Betancur , Margarita ; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Nguyen, Cattram; Murdoch Childrens Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Simpson, Julie; University of Melbourne School of Population and Global Health Carlin, John; Murdoch Childrens Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics Lee, Katherine; Murdoch Children's Research Institute, Clinical Epidemiology and Biostatistics Unit; The University of Melbourne, Department of Paediatrics
Primary Subject Heading:	Epidemiology
Secondary Subject Heading:	Public health, Research methods
Keywords:	EPIDEMIOLOGY, STATISTICS & RESEARCH METHODS, Public health < INFECTIOUS DISEASES

SCHOLARONE™
Manuscripts

The handling of missing data with multiple imputation in observational studies that address causal questions: protocol for a scoping review

Rheanna M. Mainzer^{*1,2}, Margarita Moreno-Betancur^{1,2}, Cattram D. Nguyen^{1,2}, Julie A. Simpson³, John B. Carlin^{1,2,3}, Katherine J. Lee^{1,2}

1. Clinical Epidemiology and Biostatistics Unit, Murdoch Children's Research Institute, Parkville, Victoria 3052, Australia
2. Department of Paediatrics, The University of Melbourne, Parkville, Victoria 3052, Australia
3. Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, The University of Melbourne, Parkville, Victoria 3052, Australia

*Correspondence to:

Rheanna Mainzer

rheanna.mainzer@mcri.edu.au

ABSTRACT

Introduction: Observational studies in health-related research often aim to answer causal questions. Missing data are common in these studies and often occur in multiple variables, such as the exposure, outcome and/or variables used to control for confounding. The standard classification of missing data as missing completely at random (MCAR), missing at random (MAR) or missing not at random (MNAR), does not allow for a clear assessment of missingness assumptions when missingness arises in more than one variable. This presents challenges for selecting an analytic approach and determining when a sensitivity analysis under plausible alternative missing data assumptions is required. This is particularly pertinent with multiple imputation (MI), which is often justified by assuming data are MAR. The objective of this scoping review is to examine the use of MI in observational studies that address causal questions, with a focus on if and how (i) missingness assumptions are expressed and assessed, (ii) missingness assumptions are used to justify the choice of a complete case analysis and/or MI for handling missing data, and (iii) sensitivity analyses under alternative plausible assumptions about the missingness mechanism are conducted.

Methods and analysis: We will review observational studies that aim to answer causal questions and use MI, published between January 2019 and December 2021 in five top general epidemiology journals. Studies will be identified using a full text search for the term "multiple imputation" and then assessed for eligibility. Information extracted will include details about the study characteristics, missing data, missingness assumptions and MI implementation. Data will be summarised using descriptive statistics.

Ethics and dissemination: Ethics approval is not required for this review because data will be collected only from published studies. The results will be disseminated through a peer reviewed publication and conference presentations.

Study registration: This protocol is registered on figshare (<https://doi.org/10.6084/m9.figshare.20010497.v1>).

Strengths and limitations of this study

- A targeted review of observational studies published in the five top-ranked epidemiology journals will benchmark the current state of practice for handling multivariable missingness with multiple imputation in causal analyses.
- Screening, reviewing and data extraction will be performed systematically, with double data extraction for a subset of articles and any discrepancies resolved by a panel.
- It is likely that some of the information sought will be ambiguously reported or not reported.
- Potential challenges with data extraction have been considered and a strategy for handling these challenges has been put in place.

- All extracted data and code will be made publicly available, enabling our descriptive analysis to be entirely reproducible.

INTRODUCTION

Observational studies in clinical and health-related research often aim to answer causal questions, even if this intent is only implicit.^(1, 2) This aim is usually addressed by estimation of a target parameter to quantify the impact of intervening on an exposure on an outcome of interest, in a given population. In observational studies missing data are common and can occur in multiple variables, such as the exposure, the outcome and/or the variables used to control for confounding. Restricting the statistical analysis to individuals with complete data on all analysis variables, i.e., conducting a “complete case analysis” (CCA), can lead to bias and/or loss of precision in estimates of the target parameter.⁽³⁾ Multiple imputation (MI) is a popular and flexible approach for estimating a target parameter in the presence of incomplete data.^(4, 5) In the first stage of MI, missing data are imputed multiple times with random draws from the predictive distribution of the missing values given the observed data and a specified imputation model. In the second stage, the statistical analysis of interest is applied to each imputed dataset and the results are combined using Rubin’s rules to obtain a single estimate of the target parameter with associated standard error.⁽⁴⁾

Standard implementations of MI are known to provide consistent estimation of target parameters under certain (unverifiable) assumptions about the mechanism leading to missing data. Assumptions about missing data are usually expressed using Rubin’s classification of missing data mechanisms into missing completely at random (MCAR, where the probability of data being missing does not depend on the observed or unobserved data), missing at random (MAR, where the probability of data being missing does not depend on the unobserved data, conditional on the observed data) and missing not at random (MNAR, where the probability of data being missing depends on the unobserved data, even after conditioning on the observed data).⁽⁶⁾ While this framework is useful if missing data occur in a single variable, it raises issues when missingness arises in more than one variable. First, what these mechanisms mean with multivariable missingness is poorly understood and does not allow for a transparent assessment of missingness assumptions.⁽⁷⁾ Second, based on our experience researching, teaching and applying MI, these mechanisms have become widely (mis)understood as synonymous with methods. For example, researchers often use MI under the assumption that data are MAR, but this is only a sufficient and not necessary condition for standard MI to be consistent.⁽⁸⁾ Both a CCA and a MI analysis could be unbiased under a range of multivariable missingness mechanisms (even those considered to be MNAR).⁽⁹⁾ Likewise, there are missingness mechanisms in which neither MI nor a CCA can be used to estimate an exposure-outcome association without bias, and a different approach would be needed for unbiased estimation.

The primary analysis in a study would ideally be conducted under the missing data assumptions that the researcher believes to be most likely. However, because one cannot verify from the observed data what the true missing data mechanism is, sensitivity analyses to examine how results differ under other plausible assumptions about the missingness mechanism (hereafter, “sensitivity analyses”) are strongly recommended.⁽¹⁰⁾ Such an analysis could be carried out by estimating the target parameter under the other mechanism(s) that the researcher has identified as likely. As stated by the US National Research Council, “the usefulness of a sensitivity analysis ultimately depends on the transparency and plausibility of the unverifiable assumptions.”⁽¹⁰⁾ The inherent difficulty in assessing missingness assumptions when framed in the traditional MCAR/MAR/MNAR manner is an obvious obstacle to conducting sensitivity analyses. Furthermore, from our observation, MI is routinely applied as a sensitivity analysis to a CCA. However, this practice is flawed without considering one’s plausible assumptions regarding the missingness mechanism, ⁽¹¹⁾ as neither of these approaches may be valid under particular assumptions regarding the missingness mechanism. If this is the case, obtaining similar results from a CCA and MI is not informative.

1
2
3 1 Most reviews of the handling and reporting of missing data, and the implementation and documentation of
4 2 MI, have been carried out in the context of randomised controlled trials (RCTs).(12-18) For trials, typically only
5 3 the outcome variable is incomplete, while the intervention and other key variables (typically baseline
6 4 variables) are observed for all participants. In this setting where there are missing data in a single variable, the
7 5 MCAR/MAR/MNAR framework is more transparent and guidance on sensitivity analyses has been well-
8 6 developed (see, for example, (15, 19)). In contrast, there have been few reviews concerned with how missing
9 7 data are handled in observational studies where there is the additional complication of multivariable
10 8 missingness. A review by Mackinnon published in 2010 found that only two (4%) out of 50 non-RCT studies
11 9 reviewed carried out an additional analysis that was described as a sensitivity analysis.(11) Similarly, Rezvan et
12 10 al. (2015) found that none of the 30 observational studies reviewed conducted a sensitivity analysis to
13 11 departures from the missingness assumptions following MI.(20)

14 12 While the reviews by Mackinnon and Rezvan et al. provide useful insight into the problem, neither focused
15 13 specifically on observational studies and the issues described above. In addition, subsequent to publication of
16 14 these reviews there have been important developments in the theory and application of missingness directed
17 15 acyclic graphs (m-DAGs), also known as m-graphs, a tool for the formulation of causal assumptions in the
18 16 presence of multivariable missingness.(8) M-DAGs aid the depiction and assessment of missingness
19 17 assumptions. Clarity regarding each plausible causal mechanism underlying the missing data then facilitates
20 18 the choice of analytical approach. For example, the application of DAG theory allows one to determine
21 19 whether a target parameter can be estimated without bias from the available data using an approach like CCA
22 20 or MI, or whether additional assumptions and a more sophisticated analysis is required (such as a delta-
23 21 adjusted MI approach, where imputations are shifted by a parameter “delta” representing the difference
24 22 between the observed and unobserved data).(9, 21-23)

25 23 The aim of this scoping review is to examine the use of MI in observational studies that address causal
26 24 questions relating to health. Addressing causal questions is typically the focus of epidemiological studies even
27 25 when this may not be very clearly articulated.(2) These studies often face missingness in multiple variables
28 26 required for analysis. We will examine (i) how missingness assumptions are expressed, (ii) if and how
29 27 missingness assumptions are used to justify the choice of a CCA and/or MI for handling missing data, and (iii)
30 28 the conduct of sensitivity analyses under alternative plausible assumptions about the missingness mechanism.
31 29 We will also examine how MI is implemented. This review will be used to document the current state of
32 30 practice, to identify areas for improvement in the handling and reporting of missing data with MI in
33 31 observational studies, and to subsequently develop guidance on these key components for researchers.

32 **METHODS AND ANALYSIS**

33 33 In this section we provide a full description of the study design, including how articles will be selected, what
34 34 information will be extracted, and how extracted data will be analysed. The review described in this protocol
35 35 began in June 2022 and we anticipate it will be completed by June 2023.

36 **Search strategy**

37 37 We will search five general epidemiology journals for observational studies published between January 2019
38 38 and December 2021 that aim to answer at least one causal research question using MI. The general
39 39 epidemiology journals that will be included in this search are: *International Journal of Epidemiology*, *American*
40 40 *Journal of Epidemiology*, *European Journal of Epidemiology*, *Journal of Clinical Epidemiology* and *Epidemiology*.
41 41 These journals were chosen because they are high ranking, general journals in epidemiology that publish
42 42 original research from observational studies. As such, articles from these journals should capture the current
43 43 best practice in the use of MI to handle missing data when answering causal questions using observational
44 44 data. They have also been used previously in a review of epidemiologic practice.(24) Original research articles
45 45 will be identified using the full-text search term “multiple imputation” on each journal’s website. This search
46 46 strategy is similar to that used in previous scoping reviews in this area.(11, 20)

47 **Inclusion criteria**

48 48 We will include original research articles published between January 2019 and December 2021 that aim to

1
2
3 1 answer at least one causal question using MI to handle missing data. We will determine that a study has aimed
4 2 to answer a causal question if at least one of the following criteria is satisfied:

- 5 3 1. the authors explicitly stated they were estimating a causal effect;
- 6 4 2. the study estimated an effect that was given (at least implicitly) a causal interpretation, i.e., an
7 5 interpretation which suggested that intervening on the exposure could change the outcome (e.g.,
8 6 increasing coffee consumption may be protective against stroke). This will be determined by wording
9 7 in conclusions. If it is not clear from this wording alone, investigation of the following three typical
10 8 signals of causal analyses will be used to aid in the determining: identification of confounders, the
11 9 inclusion of a DAG to illustrate causal assumption made in the analysis, and analytical approaches
12 10 incorporating adjustment for confounders (for example, estimating an effect using a regression model
13 11 that was adjusted for a set of covariates).

14 12 Studies on all disease areas/medical conditions and any target population will be considered.

15 13 **Exclusion criteria**

16 14 Studies will be excluded from the review if they meet any of the following criteria:

- 17 15 • *No causal question.* The article did not aim to answer a causal question, for example, the aim of the
18 16 study was to develop a predictive model or to estimate a disease burden.
- 19 17 • *Unclear type of question.* A clear research goal could not be identified. In other words, it was unclear
20 18 whether the study aimed to answer a descriptive, predictive or causal question.
- 21 19 • *The analysis did not use MI.*
- 22 20 • *Methodological research.* The primary purpose of the article was methodological development, for
23 21 example, using a simulation study to compare the performance of methods or mathematical
24 22 derivations to develop a new method or model. While these articles often include comprehensive
25 23 case studies, they may not be representative of empirical studies aiming primarily to answer causal
26 24 research questions.
- 27 25 • *Aggregate-level data.* The analysis was based on aggregated data where MI could not be applied at
28 26 the participant level, as is common in meta-analysis or interrupted time series analysis.
- 29 27 • *Qualitative research.* The article provided a commentary, review, opinion, study protocol, study
30 28 profile or description only.
- 31 29 • *Trial.* The study intervention was assigned to participants by the study investigators.

32 30 **Sample size**

33 31 We will require at least 100 studies to estimate the percentage of studies with a particular element (e.g.,
34 32 studies that justify their missingness assumptions) to within a maximum margin of error (two standard errors)
35 33 of 10%. Assuming a prevalence of 50%, this would give a 95% confidence interval from 40% to 60%. For a
36 34 prevalence greater than or less than 50%, the 95% confidence interval will be narrower. This sample size is
37 35 similar to the sample size used in the first review of MI in medical research ($n = 99$, (11)), and many of the
38 36 subsequent reviews in this area (e.g., $n = 103$ in (20), 77 in (15) and 118 in (12)). We expect to identify at least
39 37 100 eligible studies given the three-year publication time frame. All eligible studies will be included in the
40 38 review.

41 39 **Study selection**

42 40 The search of the journal databases and selection of studies for inclusion in the review will be performed
43 41 primarily by a single researcher (RM) in two steps. First, the title, abstract and date of each article will be
44 42 screened to rule out studies that are clearly not eligible for the review. Second, the full text of the remaining
45 43 studies will be reviewed to confirm if studies are eligible for the review. If a decision about the eligibility of an
46 44 article cannot be reached by RM (for example, due to uncertainty about the inclusion criteria), a second
47 45 researcher (CN) will independently review the full text. Disagreements about inclusion criteria will be resolved
48 46 by discussion in meetings with at least three researchers (RM, CN and at least one of JC, JS, KL or MMB).

49 47 **Data extraction and management**

Covidence, a web-based tool for systematic review management, will be used to perform the review.⁽²⁵⁾ The data extraction questionnaire was developed and tested for use by RM and KL using a sample of 10 articles. Data from all eligible studies will be extracted by RM. The supplementary material of all eligible studies will also be reviewed. We will use double data extraction (performed by KL) for a random selection of 10% of articles and additionally when there is uncertainty about the information being extracted. Discrepancies and uncertainties will be resolved by discussion in meetings with at least three researchers (RM, KL and at least one of JC, JS, CN or MMB).

8 Outcomes measured

We will extract data pertaining to the study characteristics, the amount of missing data and in which variables it occurs, missingness assumptions, methods for handling missing data and implementation of multiple imputation. Data extraction items are summarised in Table 1 and a copy of the data extraction questionnaire is provided in the Supplementary Material. Because we anticipate difficulties in extracting some items (such as the percentage of complete cases), in Supplementary Table 1 we list potential challenges in extracting data and any assumptions or simplifications that will be made if these challenges arise. Any post-hoc assumptions or simplifications for unanticipated challenges will be recorded and reported as part of the analysis.

16 **Table 1.** Summary of items to be extracted from each article

Category	Summary of data extraction items
Study characteristics	<ul style="list-style-type: none"> • First author's last name • Publication date • Journal • Type of study design
Missing data	<ul style="list-style-type: none"> • Percentage of complete cases • Percentage of missing values in the exposure and outcome • Number of incomplete covariates
Missingness assumptions	<ul style="list-style-type: none"> • Statement of missingness data assumptions (including whether the study used m-DAGs or the MCAR/MAR/MNAR framework) • Justification of missingness assumptions
Analysis methods	<ul style="list-style-type: none"> • The primary analysis method used to answer the key causal question, e.g. MI or CCA • Whether the primary analysis was justified on the basis of missingness assumptions • If applicable, any other analyses conducted to answer the key causal question that handle the missing data differently (e.g. a CCA or a delta-adjusted MI analysis) • Whether the alternative analysis was justified on the basis of missingness assumptions • If a delta-adjusted MI analysis was used, whether external information elicited from subject-matter experts was used to choose the value(s) of the delta parameter
MI implementation	<ul style="list-style-type: none"> • The method used for MI, for example, multivariate normal imputation or multiple imputation by chained equations • The statistical software used for MI • The number of imputations performed • Whether all analysis variables were included in the imputation model • Whether auxiliary variables (i.e. variables defined as potential predictors of the variable(s) with missing data and possibly also the missingness in these variables that are not included in the target analysis) were included in the imputation model • Whether interactions were included in the imputation model

18 Analysis

1
2
3 1 The questionnaire data will be cleaned and analysed in R. Descriptive statistics will be used to summarise the
4 2 extracted data. Frequencies and percentages will be presented for categorical data, for example, the method
5 3 used to obtain the primary results. Median and interquartile range will be presented for continuous data, for
6 4 example, the percentage of complete cases in each observational study. We are also collecting free-text data
7 5 on certain aspects of missing data handling to capture information that may be difficult to capture otherwise,
8 6 such as the details of the justification provided for the missingness assumptions. We will examine the free-text
9 7 data for themes and patterns. If possible, we will group responses into common themes and summarise these
10 8 themes using frequencies and percentages. If this is not possible, we will summarise the results in text. All data
11 9 and code will be made publicly available on GitHub.

10 **Reporting**

11 Findings from this review will be reported using the Preferred Reporting Items for Systematic reviews and
12 Meta-Analyses extension for Scoping Reviews (PRISMA-ScR) checklist.(26)

13 **Patient and public involvement**

14 None.

15 **ETHICS AND DISSEMINATION**

16 Ethics approval is not required for this review because data will be collected only from published studies. The
17 results will be disseminated through a peer-review publication and conference presentations.

18 **DISCUSSION**

19 Previous reviews of the handling of missing data have primarily focused on RCTs with incomplete outcome
20 data. Observational studies that answer causal questions are common and subject to greater challenges than
21 RCTs in terms of missing data as they often face missing data in multiple variables (exposure, outcome and/or
22 confounders). This paper describes a protocol for a scoping review of how MI is used to handle missing data in
23 these studies.

24 **Strengths and limitations**

25 There are several strengths to our study. A targeted review of observational studies in top epidemiology
26 journals publishing general research will benchmark the current state of practice for handling multivariable
27 missingness with MI in causal analyses. Screening, reviewing and data extraction will be performed
28 systematically. All data and code will be made publicly available, enabling our analysis to be entirely
29 reproducible. Results from the review will be reported according to best practice, using PRISMA-ScR.

30 There are also limitations. Identifying whether the aim of the research was to answer a descriptive, causal or
31 predictive question is somewhat subjective because many researchers have not adopted this classification of
32 research questions.(1) Although our targeted review will not include studies from all epidemiology journals,
33 we expect that included studies (expected to be > 100 studies from five major epidemiology journals) will be
34 sufficient to provide insight and general trends on the methods of interest. It is likely that some of the
35 information sought will be unclear or not reported. To accommodate this, we have specified how anticipated
36 challenges with data extraction will be handled if they arise.

37 **Implications of this research**

38 In addition to critically appraising the current state of the literature regarding the use and reporting of causal
39 analyses using MI to handle missing data in observational studies, this review will identify areas for
40 improvement in the handling and reporting of missing data in these studies. The results of this review will be
41 used to develop practical guidance for researchers and inform future research in these areas.

1 Funding

2 This work was supported by an Australian National Health and Medical Research Council (NHMRC) Career
3 Development Fellowship (CDF) Level 2 Grant (grant 1127984 awarded to KJL), a NHMRC Investigator Grant
4 Leadership Level 1 (grant 1196068 awarded to JAS), a NHMRC Investigator Grant Emerging Leadership Level 2
5 (grant 2009572 awarded to MMB) and a NHMRC Project Grant (grant 1166023). Research at the Murdoch
6 Children's Research Institute is supported by the Victorian Government's Operational Infrastructure Support
7 Program.

8 Contributors

9 RM conceived the study idea, developed the methodology, designed the data extraction tool, drafted and
10 revised the paper. KL developed the study idea, methodology, data extraction tool and revised the paper.
11 MMB and JS developed the study idea, methodology and revised the paper. CN developed the study idea,
12 methodology and data extraction tool. JC developed the study idea, methodology, data extraction tool and
13 revised the paper.

14 Competing interests

15 None declared.

16 REFERENCES

- 17 1. Hernán MA, Hsu J, Healy B. A second chance to get causal inference right: a classification of data
18 science tasks. *Chance*. 2019;32(1):42-9.
- 19 2. Hernán MA. The C-word: scientific euphemisms do not improve causal inference from observational
20 data. *American journal of public health*. 2018;108(5):616-9.
- 21 3. Little RJ, Rubin DB. *Statistical analysis with missing data*: John Wiley & Sons; 2019.
- 22 4. Rubin DB. *Multiple imputation for nonresponse in surveys*: John Wiley & Sons; 2004.
- 23 5. Van Buuren S. *Flexible imputation of missing data*: CRC press; 2018.
- 24 6. Sterne JAC, White IR, Carlin JB, Spratt M, Royston P, Kenward MG, et al. Multiple imputation for
25 missing data in epidemiological and clinical research: potential and pitfalls. *BMJ*. 2009;338:b2393.
- 26 7. Seaman S, Galati J, Jackson D, Carlin J. What is meant by "missing at random"? *Statistical Science*.
27 2013;28(2):257-68.
- 28 8. Mohan K, Pearl J. Graphical models for processing missing data. *Journal of the American Statistical*
29 *Association*. 2021;116(534):1023-37.
- 30 9. Moreno-Betancur M, Lee KJ, Leacy FP, White IR, Simpson JA, Carlin JB. Canonical Causal Diagrams to
31 Guide the Treatment of Missing Data in Epidemiologic Studies. *American Journal of Epidemiology*.
32 2018;187(12):2705-15.
- 33 10. National Research Council. *The prevention and treatment of missing data in clinical trials*. 2010.
- 34 11. Mackinnon A. The use and reporting of multiple imputation in medical research—a review. *Journal of*
35 *internal medicine*. 2010;268(6):586-93.
- 36 12. Tan P-T, Cro S, Van Vogt E, Szigeti M, Cornelius VR. A review of the use of controlled multiple
37 imputation in randomised controlled trials with missing outcome data. *BMC medical research methodology*.
38 2021;21(1):1-17.
- 39 13. Rabe BA, Day S, Fiero MH, Bell ML. Missing data handling in non-inferiority and equivalence trials: A
40 systematic review. *Pharmaceutical statistics*. 2018;17(5):477-88.
- 41 14. Fiero MH, Huang S, Oren E, Bell ML. Statistical analysis and handling of missing data in cluster
42 randomized trials: a systematic review. *Trials*. 2016;17(1):1-10.
- 43 15. Bell ML, Fiero M, Horton NJ, Hsu C-H. Handling missing data in RCTs; a review of the top medical
44 journals. *BMC medical research methodology*. 2014;14(1):1-8.
- 45 16. Powney M, Williamson P, Kirkham J, Kolamunnage-Dona R. A review of the handling of missing
46 longitudinal outcome data in clinical trials. *Trials*. 2014;15(1):1-11.
- 47 17. Ibrahim F, Tom BD, Scott DL, Prevost AT. A systematic review of randomised controlled trials in
48 rheumatoid arthritis: the reporting and handling of missing data in composite outcomes. *Trials*. 2016;17(1):1-8.
- 49 18. Rombach I, Rivero-Arias O, Gray AM, Jenkinson C, Burke O. The current practice of handling and
50 reporting missing outcome data in eight widely used PROMs in RCT publications: a review of the current
51 literature. *Quality of Life Research*. 2016;25(7):1613-23.

- 1
2
3 19. White IR, Horton NJ, Carpenter J, statistics rim, social, Pocock SJ. Strategy for intention to treat
4 2 analysis in randomised trials with missing outcome data. *BMJ*. 2011;342:d40.
- 5 3 20. Hayati Rezvan P, Lee KJ, Simpson JA. The rise of multiple imputation: a review of the reporting and
6 4 implementation of the method in medical research. *BMC medical research methodology*. 2015;15(1):1-14.
- 7 5 21. Tompsett DM, Leacy F, Moreno-Betancur M, Heron J, White IR. On the use of the not-at-random fully
8 6 conditional specification (NARFCS) procedure in practice. *Statistics in medicine*. 2018;37(15):2338-53.
- 9 7 22. Cro S, Morris TP, Kenward MG, Carpenter JR. Sensitivity analysis for clinical trials with missing
10 8 continuous outcome data using controlled multiple imputation: a practical guide. *Statistics in medicine*.
11 9 2020;39(21):2815-42.
- 12 10 23. Rezvan PH, Lee KJ, Simpson JA. Sensitivity analysis within multiple imputation framework using delta-
13 11 adjustment: Application to Longitudinal Study of Australian Children. *Longitudinal and Life Course Studies*.
14 12 2018;9(3):259-78.
- 15 13 24. de Vries BBP, van Smeden M, Rosendaal FR, Groenwold RH. Title, abstract, and keyword searching
16 14 resulted in poor recovery of articles in systematic reviews of epidemiologic practice. *Journal of Clinical
17 15 Epidemiology*. 2020;121:55-61.
- 18 16 25. Veritas Health Innovation. Covidence systematic review software. Melbourne, Australia.
- 19 17 26. Tricco AC, Lillie E, Zarin W, O'Brien KK, Colquhoun H, Levac D, et al. PRISMA extension for scoping
20 18 reviews (PRISMA-ScR): checklist and explanation. *Annals of internal medicine*. 2018;169(7):467-73.
- 21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

The handling of missing data with multiple imputation in observational studies that address causal questions: Protocol for a scoping review

Supplementary Material

Rheanna M. Mainzer*^{1,2}, Margarita Moreno-Betancur^{1,2}, Cattram D. Nguyen^{1,2}, Julie A. Simpson³, John B. Carlin^{1,2,3}, Katherine J. Lee^{1,2}

1. Clinical Epidemiology and Biostatistics Unit, Murdoch Children's Research Institute, Parkville, Victoria 3052, Australia
2. Department of Paediatrics, The University of Melbourne, Parkville, Victoria 3052, Australia
3. Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, The University of Melbourne, Parkville, Victoria 3052, Australia

*Corresponding author: Rheanna Mainzer; rheanna.mainzer@mcri.edu.au

Supplementary Table 1. Anticipated challenges with data extraction and how they will be handled.

Challenge for data extraction	Category of items affected	How challenge will be handled
Articles may have more than one publication date, for example, the date the article first appeared online and when it was published in-print.	Inclusion criteria	Only one publication date is required to be between January 2019 and December 2021. If two or more publication dates are between January 2019 and December 2021, the earlier date will be recorded.
There are multiple causal questions, exposures or outcomes.	Missing data	We will identify the primary causal question based on the research aims and conclusion. The proportion of missing data in the exposure, outcome and confounders used to answer this primary question will be recorded. This is expected to be acceptable in most cases. If the primary causal question cannot be identified due to multiple outcomes, we will report the missing data details for the first outcome listed in the methods section. (This is comparable to the strategy taken by Fiero et al. (1)) Similarly, if the primary causal question cannot be identified due to multiple exposures, we will report the missing data details for the first exposure listed in the methods section.
Multiple sets of covariates are used for adjustment.	Missing data	The largest adjustment set will be considered. The number of incomplete covariates will be recorded categorically (no incomplete covariates, 1 incomplete covariate, 2 or more incomplete covariates, not stated or unable to establish). This categorisation has been chosen to enable determination of multivariable missingness.

Not clear whether all variables in the target analysis were included in the imputation model.	MI implementation	If some (but not all) analysis variables were reported as being included in the imputation model then we will assume that the analysis variables not explicitly mentioned were excluded from the imputation model. If there was no description of the imputation model, then we will categorise this as “unclear”.
Not clear whether auxiliary variables or interactions were included in the imputation model.	MI implementation	If it is not explicitly stated that these were included in the imputation model, we will assume they were excluded. If there was no mention of the imputation model then we will categorise this as “unclear”.
Imputation method used not explicitly stated.	MI implementation	If the imputation method used (e.g. multivariate normal imputation or multiple imputation by chained equations) is not provided, we will infer the method used, where possible, from the statistical software procedures listed in the main paper or supplementary material. If the method is unable to be inferred, we will categorise this as “unclear”.

REFERENCE

1. Fiero MH, Huang S, Oren E, Bell ML. Statistical analysis and handling of missing data in cluster randomized trials: a systematic review. *Trials*. 2016;17(1):1-10.

Data extraction questionnaire.**Study characteristics****Authors**

First author last name, e.g., Mainzer

Publication date

Publication date (mm-yyyy).

Journal

Journal in which paper was published

1. International Journal of Epidemiology
2. American Journal of Epidemiology
3. European Journal of Epidemiology
4. Journal of Clinical Epidemiology
5. Epidemiology

Inclusion criteria

Select all that apply

1. Study authors stated they were estimated a causal effect
2. Study authors estimated an effect of an exposure on an outcome that was given (at least implicitly) a causal interpretation

Did the study use any of the following approaches (typical signals of a causal question)?

Select all that apply

1. Study used a directed acyclic graph (DAG) or m-DAG to illustrate causal assumptions made in the analysis
2. Study identified a set of variables that were used to control for confounding
3. Study estimated an effect of an exposure on an outcome using a regression model that was adjusted for a set of covariates

Causal interpretation

If the study estimated an effect that was given (at least implicitly) a causal interpretation, provide details of the text indicating this. (Copy and paste)

Type of study design

1. Prospective longitudinal study

2. Individual patient data (IPD) meta-analysis / pooled cohort analysis
3. Retrospective analysis of routinely collected data (e.g., administrative or EMR data)
4. Interrupted time series (ITS)
5. Case-control study
6. Case-cohort study
7. Cross-sectional study
8. Other

Missing data

Was the size of the inception sample* for the research question of interest available or able to be established?

*Inception sample: Participants who met eligibility criteria for inclusion in the study to answer the research question of interest, where eligibility criteria does not include any requirements for variables to be complete.

1. Yes
2. No, eligibility criteria required one or more variables to be complete
3. Other

What was the size of the inception sample?

Number or NA

Was there a reduction in participants from the inception sample to the analysis sample* due to non-response or missing data in a variable used in the analysis (exposure, outcome, covariates)?

*Analysis sample: participants who were included in the study to address the research question of interest, who may or may not having missing data for analysis variables

1. Yes
2. No
3. NA
4. Other

What was the size of the analysis sample?

Number of NA

Was the percentage of complete cases* available or able to be established?

*Cases with observed data for each variable included in the analysis that was used to answer the research question of interest. The denominator is the size of the analysis sample.

1. Yes
2. Able to establish an upper bound only
3. No

Percentage of complete cases / upper bound on the percentage of complete cases

Give number to nearest percent, e.g. 64, or NA. Use the size of the analysis sample as the denominator.

What was the exposure?

What/which exposure was considered for this review?

If there are multiple exposures: Identify the primary causal questions based on the research aims and conclusion and use the exposure in this question. If the primary causal question can not be identified due to multiple exposures, use the first exposure listed in the methods section.

Were there missing values in the exposure?

1. Yes
2. Yes, but only able to establish a lower bound on the percentage of missing values
3. Yes, but unable to establish the percentage of missing values
4. No
5. Unclear

Percentage of missing values in the exposure / lower bound on the percentage of missing values in the exposure

Give number to nearest percent, e.g. 64, or NA. Use the size of the analysis sample as the denominator.

What/which outcome was considered for this review?

If there are multiple outcomes: Identify the primary causal question based on the research aims and conclusion and use the outcome in this question. If the primary causal question can not be identified due to multiple outcomes, use the first outcome listed in the methods section.

Were there missing values in the outcome?

1. Yes
2. Yes, but only able to establish a lower bound on the percentage of missing values
3. Yes, but unable to establish the percentage of missing values
4. No
5. Unclear

1
2
3 **Percentage of missing values in the outcome / lower bound on the percentage of missing values**
4 **in the outcome**
5

6 Give number to nearest percent, e.g. 64, or NA. Use the size of the analysis sample as the
7 denominator.
8

9
10
11

12
13 **Were there missing values in the covariates?**

14 If multiple sets of covariates are used for adjustment, consider the largest adjustment set.

- 15
16 1. Yes, in 2 or more covariates
17 2. Yes, in 1 covariate only
18 3. No
19 4. Unable to establish
20
21

22
23 **Missingness assumptions**
24

25 Was a statement provided about what missingness assumptions were made?

- 26
27 1. No
28 2. Yes, authors invoked (either explicitly or implicitly) the missing at random assumption
29 3. Yes, authors provided a comprehensive description of assumptions made about the missingness
30 process for all variables subject to missing data, for example, using a m-DAG or a more simplified
31 causal diagram
32 4. Other

33
34
35 **Were missingness assumptions justified?**

36 For example, comparison of baseline data between responders and non-responders (to rule out
37 MCAR) or a substantive assessment using expert knowledge. Note, no analysis of data can rule out
38 MNAR.
39

- 40
41 1. Yes
42 2. No
43

44 **Details of justification for missingness assumptions**

45 For example, comparison of baseline data between responders and non-responders (to rule out
46 MCAR) or a substantive assessment using expert knowledge. Note, no analysis of data can rule out
47 MNAR. If missingness assumptions were not justified, enter NA.
48

49
50
51

52
53
54 **Did authors address the potential for data to be MNAR?**

- 55 1. Yes, using external evidence such as expert knowledge
56 2. Yes, but only as a study limitation
57 3. No, the possibility that data were MNAR was not addressed
58 4. Other

Analysis methods

What method was used to obtain the primary results?

1. MI using the full analysis sample
2. MI using a reduced analysis sample
3. CCA, weighted (e.g. using IPW)
4. CCA, unweighted
5. delta-adjusted MI
6. Other

Was the primary analysis justified on the basis of missingness assumptions?

1. Yes
2. No

Details of justification for primary analysis on the basis of missingness assumptions.

Examples include: (i) CCA was used because there was a small proportion of missing data that was unlikely to influence the results; (ii) CCA was used because a comparison of responders and non-responders did not rule out data being MCAR; (iii) MI was used because it was assumed that data were MAR; (iv) MI was used because comparison of responders and non-responders ruled out data being MCAR.

If the primary analysis was not justified on the basis of missingness assumptions, write "NA".

Was a secondary analysis that handles missing data differently used to answer the same causal question?

Select all that apply.

1. Yes, MI using the full analysis sample
2. Yes, MI using a reduced analysis sample
3. Yes, weighted CCA (e.g. using IPW)
4. Yes, unweighted CCA
5. Yes, delta-adjusted MI
6. No
7. Other

Was the secondary analysis justified?

1. No
2. Yes, as a sensitivity analysis (without further justification)
3. Yes, as a sensitivity analysis to examine the influence of missing data
4. Yes, as a sensitivity analysis to parametric modelling assumptions
5. Yes, as a sensitivity analysis to causal assumptions made about the missing data mechanism
6. NA
7. Other

If a delta-adjusted analysis was used, was external information incorporated in the analysis?

If not delta-adjusted analysis select NA

1. Yes
2. No or not stated
3. NA

If a delta-adjusted analysis was used, provide details of the delta-adjusted analysis

How was external information incorporated? What values of delta were considered? How was the analysis implemented? Etc. If no delta-adjusted analysis was used, enter NA.

MI implementation

What method was used for multiple imputation?

If the imputation method used (e.g. multivariate normal imputation or multiple imputation by chained equations) is not provided, we will infer the method used, where possible, from the statistical software procedures listed in the main paper or supplementary material. If the method is unable to be inferred, we will categorise this as “unclear”.

1. MICE
2. MVNI
3. Unclear
4. Other

What software was used for multiple imputation?

1. R
2. SAS
3. SPSS
4. Stata
5. Unclear
6. Other

Number of imputations used in the multiple imputation procedure

Were all analysis variables included in the imputation model?

If some (but not all) analysis variables were reported as being included in the imputation model then we will assume that the analysis variables not explicitly mentioned were excluded from the imputation model. If there was not description of the imputation model, then we will categorise this as “unclear”.

1. Yes
2. No
3. Unclear

Were auxiliary variables included in the imputation model?

1
2
3 If it is not explicitly stated that these were included in the imputation model, we will assume they
4 were excluded. If there was no mention of the imputation model, then we will categorise this as
5 “unclear”.

- 6
7 1. Yes
8 2. No
9 3. Unclear

10 11 **Were interactions included in the imputation model?**

12
13 If it is not explicitly stated that these were included in the imputation model, we will assume they
14 were excluded. If there was no mention of the imputation model, then we will categorise this as
15 “unclear”.

- 16
17 1. Yes
18 2. No
19 3. Unclear

20 21 22 **Reported results**

23
24 **If results were obtained using both a CCA and MI, did the authors observe any substantial**
25 **difference between these?**

26
27 Substantial difference: a difference that the authors acknowledged as important or significant (for
28 example, based on a clinical cut-off or a P values)

- 29
30 1. Yes
31 2. No
32 3. NA

33
34 **If results were obtained using both a CCA and MI, AND no substantial difference between these**
35 **two sets of results was observed, was any interpretation or explanation provided for the**
36 **similarities between the two sets of results? If so, what was the interpretation or explanation.**

37
38 If yes, add details. Otherwise: no or NA.

39
40
41
42

43 44 45 **Other**

46 47 **Funding**

48
49 How was the study funded?

50
51
52
53

54 55 **Any other comments?**

56
57
58
59
60